# Handbook of Surface and Nanometrology

David J Whitehouse

*University of Warwick*

# I*o*P

Institute of Physics Publishing
Bristol and Philadephia

# Dedication

*This book is dedicated with love to my wife Ruth, without whose support and practical help it would never been started, let alone finished. She has suffered with great forebearance the pressures attendant upon the embarkation of such a mammoth work—again!*

# Contents

# Preface

In the past few years there have been some significant changes in metrology brought about by miniaturization on the one hand, and the increasing use of simulations and packages which allow easy visualization but divorce from practice on the other.

The effect of miniaturization has been to make the engineer jump down many scales of size from automotive type components to those of mm or even much less. Pushing up from even smaller scales are the physicist anxious to build practical working devices. These roles are almost the exact opposite of the traditional positions. Where these disciplines meet is at the 'nano' level taken to be 100 nm down to about 0.1nm. So at the nano level there is an uncomfortable confrontation between two disciplines both in unfamiliar guises: often having different word meanings and certainly different strategies.

The objective of this Handbook is to bring the two metrologies together in this nano range. Because most of the new technology is surface oriented, the vehicle for bringing the engineering and physics together is surface metrology. It plays a significant role in conventional engineering as well as at the semiconductor level. It is the natural bridge. The title of the book 'Handbook of Surface and nanometrology' mirrors this contact.

One side defect of this somewhat incompatible stand-off is a dearth of proven and tested facts; The Handbook is therefore not strictly an ordinary Handbook. It has been written with factual information as expected but in addition there has been an attempt to surround the metrology with suggestions, with new ideas and proceedings to help gel the new discipline. Each chapter has this format.

It will be noticed that much of current practice in surface and nanometrology will be questioned, hopefully constructively, throughout the book. It is vital that new researchers and practitioners approach surface nanometrology or nano surface metrology with open minds. Emphasis is placed on systems wherever possible to avoid the 'parameter rash' which seems to accompany computer packages wherever they are used. Understanding the boundary conditions is essential for physicists as well as engineers. It is hoped that this book will be stimulating as well as interesting and useful.

**David Whitehouse**
Emeritus Professor, School of Engineering
University of Warwick

# Acknowledgments

In chapter 8 figure 8.7 was provided by Prof. Miyashita, 8.1 by Tanaguchi and 8.16, 8.20 and 8.34 from Inst. Phys. Publishing Nanotechnology. Table 8.3 and figure 8.36 showing details of the 'MMM' instrument courtesy of Dr Clayton Teague (NIST).

IOP Publishing Ltd. has attempted to trace the copyright holder of the material reproduced in this publication and apologizes to copyright holders if permission to publish in this form has not been obtained.

# Chapter 1
# General philosophy of measurement

## 1.1   Where does surface metrology fit in general metrology, and what about nanometrology?

Before this question can be answered it would be a good idea to define the word metrology. This is the science of measurement.

Surface metrology is the measurement of the deviations of a workpiece from its intended shape, that is from the shape specified on the drawing. It is taken to include such features as deviations from roundness, straightness, flatness, cylindricity and so on. It also includes the measurement of surface texture. This book is devoted to this subject which represents a very important sector of engineering metrology. Also the role of surface metrology on a broader front is explored. In particular, where surface metrology fits into nanometrology will be a theme of the book.

Perhaps the best way to place the role of surface metrology is to consider just what needs to be measured in order to enable a workpiece to work according to the designer's aim — one has to measure in order to be able to control.

Assuming that the material has been specified correctly and that the workpiece has been made from it, the first thing to be done is to measure the dimensions. These will have been specified on the drawing to a tolerance. Under this heading is included the measurement of length, area, position, radius and so on.

So, *dimensional metrology* is a first aim because it ensures that the size of the workpiece conforms to the designer's wish. This in turn ensures that the workpiece will assemble into an engine, gearbox, gyroscope or whatever; the *static* characteristics have therefore been satisfied.

This by itself is not sufficient to ensure that the workpiece will satisfy its function; it may not be able to turn or move, for example. This is where surface metrology becomes important.

Surface metrology ensures that all aspects of the surface geometry are known and preferably controlled. If the shape and texture of the workpiece are correct then it will be able to move at the speeds, loads and temperatures specified in the design; the *dynamic* characteristics have therefore been satisfied.

The final group of measurements concerns the physical and chemical condition of the workpiece. This will be called here *physical metrology*. It includes the hardness of the materials, both in the bulk and in the surface layers, and the residual stress of the surface, both in compression or in tension, left in the material by the machining process or the heat treatment. It also includes measurement of the metallurgical structure of the material, and its chemical construction. All these and more contribute to the durability of the component, for example its resistance to corrosion or fatigue.

Physical metrology therefore is the third major sector of engineering metrology: the long-term characteristics.

As a general rule, all three types of measurement must take place in order to ensure that the workpiece will do its assigned job for the time specified; to guarantee its quality.

This book, as mentioned earlier, is concerned specifically with surface metrology but this does not necessarily exclude the other two. In fact it is impossible to divorce any of these disciplines completely from the

others. After all, there is only one component and these measurements are all taken on it. Physical and chemical features of importance will be covered in some detail as and when they are necessary in the text.

## 1.2    Importance of surface metrology

Figure 1.1 shows where the block representing surface metrology can be placed relative to blocks representing manufacture and function. In the block 'manufacture' is included all aspects of the manufacturing process such as machine performance, tool wear and chatter, whereas in the block marked 'function' is included all functional properties of the surfaces of components such as the tribological regimes of friction, wear and lubrication. The measurement block includes the measurement of roughness, roundness and all other aspects of surface metrology.



**Figure 1.1**

The figure shows that the texture and geometry of a workpiece can be important in two quite different applications: one is concerned with controlling the manufacture (this is examined in detail in chapter 6) and the other is concerned with the way in which the surface can influence how well a workpiece will function. Many of these uses fall under the title of tribology — the science of rubbing parts — but others include the effect of light on surfaces and also static contact.

In fact the two blocks 'manufacture' and 'function' are not completely independent of each other, as illustrated by the line in the figure joining the two. Historically the correct function of the workpiece was guaranteed by controlling the manufacture. In practice what happened was that a workpiece was made and tried out. If it functioned satisfactorily the same manufacturing conditions were used to make the next workpiece and so on for all subsequent workpieces. It soon became apparent that the control of the surface was being used as an effective go-gauge for the process and hence the function. Obviously what is required is a much more flexible and less remote way of guaranteeing functional performance; it should be possible, by measuring parameters of the surface geometry itself, to predict the function. The conventional method is very much a balancing act and unfortunately the delicate equilibrium can be broken by changing the measurement parameters or the production process or even the function. This may seem an obvious statement but within it lies one of the main causes for everyday problems in engineering.

The block diagram of figure 1.1 will keep showing up in the text and in particular in chapter 9 where it will be updated by taking into account what has been revealed in previous chapters of the book.

It will soon be made clear that surface metrology cannot simply be regarded as an irritant to be added to the general size measurement of a component. The smallness of its magnitude does not infer the smallness of its importance. It will be shown that the surface geometry is absolutely critical in many applications and that it contains masses of information that can be invaluable to a manufacturer if extracted correctly from the mass of data making up the surface. In almost every example of its use the criticism can be raised that it is not well understood and, even where it is, it is not properly specified especially on drawings.

From figure 1.1 it can be seen that the importance of surface geometry should be recognized not just by the quality control engineer or inspector but also, more importantly, by the designer. It is he or she who should understand the influence that the surface has on behaviour and specify it accordingly. It is astonishing just how ill-informed most designers are where surface metrology is concerned. Many of them in fact consider that a 'good' surface is necessarily a smooth one, the smoother the better. This is not only untrue, in many cases, it can be disastrous.

Surface metrology has the ingredients of a number of widely dissimilar disciplines. Since the introduction of scanning probe microscopes (SPM) it has broadened even more. The first part of the book is concerned with the subject itself, that is its theory in chapters 2 and 3, and the instrumentation in chapters 4 and 5. The uses of surface metrology are contained in chapters 6 and 7. chapter 8 is concerned with the relatively new aspect of surface metrology, namely nanometrology. This is slightly misleading because the measurement of surfaces has been in the nanometre range 100 — 0.1nm for many years; it simply has not been called by the name. Although chapter 8 has been called nanometrology, aspects of the impact of the new generation of instruments have permeated every chapter.

To guide the reader through the subject in more detail, the content of each chapter will be outlined in what follows.

Chapter 2 starts with a broad classification of characteristics which is followed by a discussion of reference lines. These are an integral part of surface characterization because all metrology is a comparison between two things: a test or workpiece and a reference or standard. Some early bases such as the M and E systems are considered together with some recent variants. Justification for waviness separation is also inserted for those researchers who think that the geometry should all be lumped together.

Random process analysis is described and how it can be used in surface metrology. In particular, how correlation methods and spectral analysis can be best used in controlling manufacturing processes is discussed. It is demonstrated how some parameters relevant to tribological applications can be deduced from general surface statistics and multinormal distributions. One treatment due to Nayak (Longuett-Higgins), and a different approach using discrete analysis given by Whitehouse, are compared. Additionally, some new ways of utilizing peak definitions are examined. It is remarkable how some of the older work pre-empts recent developments such as fractal analysis.

The vexed problem of areal analysis is dealt with in some detail. It is still sometimes called 3D analysis despite there being only two independent variables. However, this does not have to be a problem as long as researchers and users agree to differ. What does emerge is that describing even some of the simpler manufacturing processes is difficult. Stout drops some of the analytical methods in favour of 'functional parameters,' which is an insurance against making no progress. As in profile assessment, random processes and some of the more complicated yet useful techniques to extend the analysis to include spatial terms as well as frequency are described. Wigner distributions, ambiguity functions and wavelets come into this category.

Surfaces whose geometry is a mixture of processes, such as plateau honing, are described as well as some of the parameters derived from them. Mixed process surfaces are growing in application because of the increasing demands being made.

Often, surface metrology is considered to be just surface texture but, in fact, form and deviations from all sorts of geometrical shapes have to be included, of which out-of-roundness is prominent. Deviations from

straightness, flatness and methods of assessing cylinders, tapers, cones and such shapes as epitrochoids are also included in this chapter. Finally, singular defects sometimes found on the surface are classified.

Chapter 3 is particularly concerned with methods of carrying out some of the operations used in the characterization of chapter 2.

Digital methods and sources of error are considered early in the chapter. Digital methods are applied to random process analysis dealing with surface parameters such as summits. The numerical model is considered in some detail, especially in dealing with areal analysis. Comparison with digital profile analysis is made with some surprising results. Three, four, five and seven point models for summits and their characteristics are compared.

Digital analysis using DFFT is considered, with particular reference to space frequency functions. Graphical methods and analogue methods of processing are removed from this version of the Handbook because of the dominance of digital methods.

Chapter 4 is concerned with aspects of instrumentation from design principles to performance tables. The chapter starts with static and dynamic considerations of instrument systems as well as the nature of the interactions between mechanical force and metrology loops.

Because of the pre-eminence of stylus methods in engineering they are dealt with in some detail. The behaviour of stylus systems to the input of random and sinusoidal signals is investigated. It is shown how the effective bandwidth can be increased. The scanning probe instruments are described here as an adjunct to conventional methods. Scanning tunnelling microscopes (STM) and atomic force microscopes (AFM) are discussed in some detail. In many ways the usual planar form of small objects makes their design and use simpler than the engineering of surface instruments.

Optical methods are then explored together with an additional comparison between stylus and optical methods. Optical followers which imitate the stylus instruments are examined. This is followed by other optical methods using phase detection systems and interference. Heterodyne methods with different polarizations and frequencies are compared.

White light interferometry and absolute measurement are discussed with 'stitching' followed by Moiré methods and holographic methods. Speckle and other diffraction methods are investigated. This shows some interesting advantages over conventional optics in machine tool monitoring and defect detection.

Other traditional methods such as inductance capacitance and pneumatic are considered and also the other main contender for surface measurement, i.e. ultrasonics.

Scanning electron microscopy (SEM) and transmission electron microscopy (TEM) are described with the less well known photon tunnelling microscope (PTM). The latter could have been included earlier but is inserted here because of the electromagnetic connection.

Because the pick-up system comprising the stylus and transducing system is so important, some emphasis has been put on it. Considerations include types of conversion, noise levels and limitations. Optical inductive and capacitance methods are compared with stylus methods. It is ironic that the stylus technique has had a new lease of life since the introduction of the scanning microscopes. Concepts like contact, friction and stiffness all have now to be queried as the scale of size reduces to the nanometre and even atomic levels.

One of the major problems with the STMs and AFMs is their calibration. Conventional methods do not satisfy the necessary requirements because the unit of scale — the wavelength of light — is three decades too large. This is a very relevant problem, considered in chapter 5 along with the usual problems in error propagation traceability etc.

The chapter begins by defining the various errors and then describes some statistical tests useful in surface metrology. Such tests involve variance and mean value tests such as the F test and T test. Some consideration is given to the design of experiments — all with the aim of improving or better understanding the input data.

Calibration takes up a large part of the chapter. In surface metrology it is not only lateral $x$, $y$, scales and the normal $z$ axis but also the dynamics of filters that have to be included. The separation of errors between

the test piece and the instrument reference is examined with the idea of improving the capability of the instrument. Some alternatives in roundness are described. In response to the need to calibrate SPMs, techniques using the lattice spacing of silicon using x-rays together with a cleavage method on topaz are described. Some remnant problems with SPM calibration are pointed out: in particular, the problems with "metrological" instruments.

The geometric product specification (GPS) of many geometrical factors such as texture is investigated with a view to developing a chain of standards to ensure traceability. The linking of the geometrical features to their appropriate international standard is given, along with some international standards.

Drawing symbols for the surface texture process marks are the final subjects for the chapter.

Chapter 6 is concerned with the way in which the manufacturing process affects the surface and, conversely, how problems in the process or machine tool can be identified and monitored by measuring the surface. The fact is that the surface generated by the process is very sensitive to changes so process control is easily possible.

In the first part of the chapter some traditional processes are examined. These include turning, milling, broaching etc. Some factors which affect the surface such as tool wear and chip formation mechanisms are examined. Also, recent trends in dry cutting are reported. This is an important development because optical in-process surface measurement becomes viable due to the absence of coolant mists.

Abrasive processes such as grinding and polishing are included in the chapter. Nanogrinding and ductile machining in such materials have different properties at the molecular scale of size. Roundness errors due to chatter and elastic effects of the machine tool are part of this section. It also includes the various different types of grinding, including centreless grinding and creep grinding.

Non-conventional machining is considered. These include physical and chemical machining e.g. ECM and EDM.

Large scale fabrication and forging are included, together with near atomic machining, e.g. ion beam milling. Under this part of the chapter, designer (structured) surfaces are illustrated with comments on future applications. Some chip formations and plastic flow constitute one aspect of nanotechnology of manufacture.

How to make use of some of the newer analytical methods for machine tool performance is briefly discussed.

Chapter 7 is concerned with the influence of the surface on performance. It is particularly difficult because of the scarcity of substantiated information. Targeting a surface to satisfy a particular function is the job of the designer, yet the texture aspect is often lost or neglected in the totality of the design. One reason is the complexity of the function. There is a macro view of function which the designer does see and this is whether the surface is important or not; surfaces are obviously important in bearing. How to quantify performance is the problem. This is the micro or nano issue. The fundamental problems are not seen at the macro level which is why, in this chapter, some time is spent trying to clarify the situation with the concept of function maps. These are used for the first time to link the manufacture to function via the metrology. As this permeates the whole of the chapter because of its usefulness it will be referred to often in the text.

The chapter starts with some examples of the macroscopic behaviour of simple geometries in contact which is followed by microscopic contact issues including some mechanisms of contact. The relevance of fractals and the way in which waviness affects the actual contact between bodies is included in this section.

Various models of elastic and plastic behaviour are discussed and compared. The areal (3D) model is also discussed. How the models help in understanding fundamental effects of contact comes next.

Stiffness, mechanical seals, adhesion, thermal and electrical contact come under this heading of what is basically static contact in the 'normal' direction. This is followed by those applications involving lateral relative movement. These include friction wear and lubrication; in effect all aspects of tribology and shakedown.

Surface influence when normal loading is severe as in elasto- and plasto-hydrodynamic lubrication regimes, as well as boundary lubrication, are investigated and the associated scuffing failure considered. Fatigue and its variants that have cyclic or random loadings is investigated. These include rolling fatigue, fretting and some effects of squeeze films.

Areal effects of surfaces such as the lay and how it affects lubrication are difficult to model, especially when macro as well as micro geometry have to be considered. Models for this are referred to in some detail with interesting conclusions about mixing transverse and longitudinal surface patterns.

The possibility of using Wiebull distribution in characterizing wear is introduced as is the possible application of space frequency functions.

The importance of single body application cannot be overestimated. One aspect where the surface is important is in the triggering of fatigue failure as well as in initiating corrosion. Optical performance such as gloss is greatly influenced by the surface.

Scatter from deterministic and random surfaces are discussed with sections on shadowing and caustics. Fractal models and their implication are included. Extra consideration is given to the generation of aspheric surfaces and how the wavefront is modified by roughness.

Scattering of other waves such as acoustic, ultrasonic waves, elastic and non elastic, Bragg, Rayleigh, Raman scattering and surface influence are compared.

Finally in the general text tolerances and fits are considered.

In the discussion some ways of trying to quantify surface performance, and in particular areal (3D) parameters, are reported and some extensions are made to the function map in which more than one surface is within the surface characterization. A new aspect, considered in chapters 2 and 7 is that of the surface system in which not just one surface is considered, but two. Novel methods of system characterization are described and discussed. These include the use of cross convolution.

Nanotechnology is considered to be one of the most important developments since the computer which 'enables' progress to be made in many disciplines such as biology, MEMS, cluster theory and so on. It is not usually realized that detection of surface texture has long been within the range of size usually ascribed to nanotechnology (i.e. 100nm to 0.1 nm).

Chapter 8 discusses some of the implications of the nanotechnology of surfaces, called here nano surface metrology rather than nanometrology. Aspects of nanometrology, in engineering and instrumentation are discussed; some examples are taken from the main body of the text.

Comments are made on how quantum effects are now being felt in engineering as well as in physics and chemistry. It is shown how the very nature of traditional components of surface metrology — texture and form — is changing. Disciplines which were thought to be straightforward are no longer so at the small nanotechnology scale.

Chapter 9 attempts to draw conclusions from the other chapters. It shows how the book has extended the traditional role of metrology to encompass integrated measurement, the changing role of surfaces, the characterization of striated surfaces and the system of surfaces. Function maps and their role in the future are discussed as well as some problems encountered in making progress.

# Chapter 2
# Surface characterization

## The nature of surfaces

Surface characterization, the nature of surfaces and the measurement of surfaces cannot be separated from each other. A deeper understanding of the surface geometry produced by better instrumentation often produces a new approach to characterization.

Surface characterization is taken to mean the breakdown of the surface geometry into basic components based usually on some functional requirement. These components can have various shapes, scales of size, distribution in space and can be constrained by a multiplicity of boundaries in height and position. Issues like the establishment of reference lines can be viewed from their ability to separate geometrical features or merely as a statement of the limit of instrument capability. Often one consideration determines the other! Ease of measurement can influence the perceived importance of a parameter or feature. It is difficult to ascribe meaningful significance to something which has not been or cannot be measured. One dilemma is always whether a feature of the surface is fundamental or simply a number which has been ascribed to the surface by an instrument. This is an uncertainty which runs right through surface metrology and is becoming even more obvious now that atomic scales of size are being explored. Surface, interface and nanometrology are merging.

For this reason what follows necessarily reflects the somewhat disjointed jumps in understanding brought on by improvements in measurement techniques. There are no correct answers. There is only a progressively better understanding of surfaces brought about usually by an improvement in measurement technique. This extra understanding enables more knowledge to be built up about how surfaces are produced and how they perform.

This chapter therefore is concerned with the nature of the geometric features, the signal which results from the measuring instrument, the characterization of this signal and its assessment. The nature of the signal obtained from the surface by an instrument is also considered in this chapter. How the measured signal differs from the properties of the surface itself will be investigated. Details of the methods used to assess the signal will also be considered but not the actual data processing. This is examined in chapter 3. There is, however, a certain amount of overlap which is inevitable. Also, there is some attention paid to the theory behind the instrument used to measure the surface. This provides the link with chapter 4.

What is set down here follows what actually happened in practice. This approach has merit because it underlies the problems which first came to the attention of engineers in the early 1940s and 1950s. That it was subsequently modified to reflect more sophisticated requirements does not make it wrong; it simply allows a more complete picture to be drawn up. It also shows how characterization and measurement are inextricably entwined, as are surface metrology and nanometrology, as seen in chapter 8.

It is tempting to start a description of practical surfaces by expressing all of the surface geometry in terms of departures from the desired three-dimensional shape, for example departures from an ideal cylinder or sphere. However, this presents problems, so rather than do this it is more convenient and simpler to start off by describing some of the types of geometric deviation which do occur. It is then appropriate to show how

these deviations are assessed relative to some of the elemental shapes found in engineering such as the line or circle. Then, from these basic units, a complete picture can be subsequently built up. This approach has two advantages. First, some of the analytical tools required will be explained in a simple context and, second, this train of events actually happened historically in engineering.

Three widely recognized causes of deviation can be identified:

1. The irregularities known as roughness that often result from the manufacturing process. Examples are (*a*) the tool mark left on the surface as a result of turning and (*b*) the impression left by grinding or polishing. Machining at the nanoscale still has process marks.

2. Irregularities, called waviness, of a longer wavelength caused by improper manufacture. An example of this might be the effects caused by a vibration between the workpiece and a grinding wheel.

3. Very long waves referred to as errors of form caused by errors in slideways, in rotating members of the machine, or in thermal distortion.

Often the first two are lumped together under the general expression of surface texture, and some definitions incorporate all three! Some surfaces have one, two or all of these irregularities [1]. Figure 2.1 shows roughness and waviness superimposed on the nominal shape of a surface.

A question often asked is whether these three geometrical features should be assessed together or separately. This is a complicated question with a complicated answer. One thing is clear; it is not just a question of geometry. The manufacturing factors which result in waviness, for instance, are different from those that produce roughness or form error. The effect of these factors is not restricted to producing an identifiable geometrical feature, it is much more subtle: it affects the subsurface layers of the material.

Furthermore, the physical properties induced by chatter, for example, are different from those which produce roughness. The temperatures and stresses introduced by general machining are different from those generated by chatter. The geometrical size of the deviation is obviously not proportional to its effect underneath the surface but it is at least some measure of it. On top of this is the effect of the feature of geometry on function in its own right. It will be clear from the section on function how it is possible that a long-wavelength component on the surface can affect performance differently from that of a shorter wavelength of the same amplitude. There are, of course, many examples where the total geometry is important in the function of the workpiece and under these circumstances it is nonsense to separate out all the geometrical constituents. The same is true from the manufacturing signature point of view.

From what has been said it might be thought that the concept of 'sampling length' is confined to roughness measurement in the presence of waviness. Historically this is so. Recent thoughts have suggested that in order to rationalize the measurement procedure the same 'sampling length' procedure can be adopted to measure 'waviness' in the presence of form error, and so on to include the whole of the primary profile. Hence $l_r$, the sampling length for roughness, is joined by $l_w$. For simplicity roughness is usually the surface feature considered in the text.

It is most convenient to describe the nature and assessment of surface roughness with the assumption that no other type of deviation is present. Then waviness will be brought into the picture and finally errors for form. From a formal point of view it would be advantageous to include them all at the same time but this implies that they are all able to be measured at the same time, which is only possible in some isolated cases.

## 2.1 Surface roughness characterization

Surface roughness is that part of the irregularities on a surface left after manufacture which are held to be inherent in the material removal process itself as opposed to waviness which may be due to the poor performance of an individual machine. (BS 1134 1973) mentions this in passing.

(*a*)

Roughness    Waviness

Nominal shape

(*b*)

Subsurface plastic damage
caused by roughness

Roughness

Waviness

Subsurface elastic effects
caused by waviness

(*i*)  Stress profile vs. Geometric profile.

Energy into surface

Roughness    Waviness

λ

Geometric amplitude

Roughness    Waviness

λ

(*ii*)  Energy vs. Geometry (function of wavelength).

Deformation

Waviness effect

Elastic deformation

Roughness effect

Plastic deformation

μ

(*iii*)  Mode of deformation.

**Figure 2.1** (*a*) Geometric deviations from intended shape.

In general, the roughness includes the tool traverse feed marks such as are found in turning and grinding and the irregularities within them produced by microfracture, built-up edge on the tool, etc.

The word 'lay' is used to describe the direction of the predominant surface pattern. In practice it is considered to be most economical in effort to measure across the lay rather than along it, although there are exceptions to this rule, particularly in frictional problems or sealing (see chapter 7).

Surface roughness is generally examined in plan view with the aid of optical and electron microscopes, in cross-sections normal to the surface with stylus instruments and, in oblique cross-sections, by optical interference methods. These will be discussed separately in a later section (4.3.2). First it is useful to discuss the scales of size involved and to dispel some common misconceptions.

Surface roughness covers a wide dimensional range, extending from that produced in the largest planing machines having a traverse step of 20mm or so, down to the finest lapping where the scratch marks may be spaced by a few tenths of a micrometre. These scales of size refer to conventional processes. They have to be extended even lower with non-conventional and energy beam machining where the machining element can be as small as an ion or electron, in which case the scale goes down to the atomic in height and spacing. The peak-to-valley height of surface roughness is usually found to be small compared with the spacing of the crests; it runs from about 50 $\mu m$ down to less than a few thousandths of a micrometre for molecular removal processes. The relative proportions of height and length lead to the use of compressed profile graphs, the nature of which must be understood from the outset. As an example figure 2.2 shows a very short length of the profile of a cross-section of a ground surface, magnified 5000 ×.



**Figure 2.2** Distortion caused by making usable chart length.

The representation of the surface by means of a profile graph will be used extensively in this book because it is a very convenient way to portray many of the geometrical features of the surface. Also it is practical in size and reflects the conventional way of representing surfaces in the past. That it does not show the 'areal' characteristics of the surface is understood. The mapping methods described later will go into this other aspect of surface characterization. However, it is vital to understand what is in effect a shorthand way of showing up the surface features. Even this method has proved to be misleading in some ways, as will be seen.

The length of the section in figure 2.2 from A to D embraces only 0.1mm of the surface, and this is not enough to be representative. To cover a sufficient length of surface profile without unduly increasing the length of the chart, it is customary to use a much lower horizontal than vertical magnification. The result may then look like figure 2.2(*b*). All the information contained in the length AD is now compressed into the portion A′D′, with the advantage that much more information can be contained in the length of the chart, but with the attendant disadvantage that the slopes of the flanks are enormously exaggerated, in the ratio of the vertical to horizontal magnifications. Thus it is essential, when looking at a profile graph, to note both magnifications and to remember that what may appear to be fragile peaks and narrow valleys may represent quite gentle undulations on the actual surface. Compression ratios up to 100:1 are often used. Many models of surfaces used in tribology have been misused simply because of this elementary misunderstanding of the true dimensions of the surface. Examination of an uncompressed cross-section immediately highlights the error in the philosophy of 'knocking off of the peaks during wear'!

The photomicrographs and cross-sections of some typical surfaces can be examined in figure 2.3. The photomicrographs (plan or so-called areal views) give an excellent idea of the lay and often of the distance (or spacing) between successive crests, but they give no idea of the dimensions of the irregularities measured normal to the surface.

The profile graph shown beneath each of the photomicrographs is an end view of approximately the same part of the surface, equally magnified horizontally, but more highly magnified vertically. The amount of distortion is indicated by the ratio of the two values given for the magnification, for example $15000/150$, of which the first is the vertical and the second the horizontal magnification.

In principle, at least two cross-sections at right angles are needed to establish the topography of the surface, and it has been shown that five sections in arbitrary directions should be used in practice; when the irregularities to be portrayed are seen to have a marked sense of direction and are sufficiently uniform, a single cross-section approximately at right angles to their length will often suffice. Each cross-section must be long enough to provide a representative sample of the roughness to be measured; the degree of uniformity, if in doubt, should be checked by taking a sufficient number of cross-sections distributed over the surface.

When the directions of the constituent patterns are inclined to each other, the presence of each is generally obvious from the appearance of the surface, but when a profile graph alone is available, it may be necessary to know something of the process used before being able to decide whether or not the profile shows waviness. This dilemma will be examined in the next section.

Examination of the typical waveforms in figure 2.3 shows that there is a very wide range of amplitudes and crest spacings found in machining processes, up to about five orders of magnitude in height and three in spacing. Furthermore, the geometric nature of the surfaces is different. This means, for example, that in the cross-sections shown many different shapes of profile are encountered, some *engrailed* in nature and some *invected* and yet others more or less random. Basically the nature of the signals that have to be dealt with in surface roughness is more complex than those obtained from practically any sort of physical phenomena. This is not only due to the complex nature of some of the processes and their effect on the surface skin, but also due to the fact that the final geometrical nature of the surface is often a culmination of more than one process, and that the characteristics of any process are not necessarily eradicated completely by the following one. This is certainly true from the point of view of the thermal history of the surface skin, as will be discussed later.

It is because of these and other complexities that many of the methods of surface measurement are to some extent complementary.

Some methods of assessment are more suitable to describe surface behaviour than others. Ideally, to get a complete picture of the surface, many techniques need to be used; no single method can be expected to give the whole story. This will be seen in chapter 4 on instrumentation.

The problem instrumentally is therefore the extent of the compromise between the specific fidelity of the technique on the one hand and its usefulness in as many applications as possible on the other.

**Figure 2.3** Photomicrographs showing plan view, and graphs showing cross-section (with exaggerated scale of height) of typical machined surfaces.

The same is true of surface characterization: for many purposes it is not necessary to specify the whole surface but just a part of it. In what follows the general problem will be stated. This will be followed by a breakdown of the constituent assessment issues. However, it must never be forgotten that the surface is three dimensional and, in most functional applications, it is the properties of the three-dimensional gap between two surfaces which are of importance. Any rigorous method of assessing surface geometry should be capable of being extended to cover this complex situation. An attempt has been made in chapter 7.

The three-dimensional surface $z = f(x, y)$ has properties of height and length in two dimensions. To avoid confusion between what is a three-dimensional or a two-dimensional surface, the term 'areal' is used to indicate the whole surface. This is because the terms 2D and 3D have both been used in the literature to mean the complete surface. There are a number of ways of tackling the problem of characterization; which is used is dependent on the type of surface, whether or not form error is present and so on. The overall picture of characterization will be built up historically as it occurred. As usual, assessment was dominated by the available means of measuring the surface in the first place. So because the stylus method of measurement has proved to be the most useful owing to its convenient output, ease of use and robustness, and because the stylus instrument usually measures one sample of the whole surface, the evaluation of a single cross-section (or profile) will be considered first. In many cases this single-profile evaluation is sufficient to give an adequate idea of the surface; in some cases it is not. Whether or not the profile is a sufficient representation is irrelevant; it is the cornerstone upon which surface metrology has been built. In subsequent sections of this chapter the examination of the surface will be extended to cover the whole geometry. In the next section it will be assumed that the cross-section does not suffer from any distortions which may be introduced by the instrument, such as the finite stylus tip or limited resolution of the optical device. Such problems will be examined in detail in section 4.3.1.

Problems of the length of profile and the reliability of the parameters will be deferred until chapter 5.

### 2.1.1 Profile parameters

The profile graph shown in figure 2.4 and represented by $z = f(x)$ could have been obtained by a number of different methods but it is basically a waveform which could appear on any chart expressing voltage, temperature, flow or whatever. It could therefore be argued that it should be representable in the same sort of way as for these physical quantities. To some extent this is true but there is a limit to how far the analogy can be taken.



**Figure 2.4**  Typical profile graph.

The simplest way of looking at this is to regard the waveform as being made up of amplitude (height) features and wavelength (spacing) features, both independent of each other. The next step is to specify a minimum set of numbers to characterize both types of dimension adequately. There is a definite need to constrain the number of parameters to be specified even if it means dropping a certain amount of information, because in practice these numbers will have to be communicated from the designer to the production engineer using the technical drawing or its equivalent. More than two numbers often cause problems of comprehension. Too many numbers in a specification can result in all being left off, which leaves a worse situation than if only one had been used. Of the height information and the spacing information it has been conventional to regard the height information as the more important simply because it seems to relate more readily to functional importance. For this reason most early surface finish parameters relate only to the height information in the profile and not to the spacing. The historical evolution of the parameters will be described in the introduction to chapter 4. The definitive books prior to 1970 are contained in references [1–7].

However, there has been a general tendency to approach the problem of amplitude characterization in two ways, one attempting in a crude way to characterize functions by measuring peaks, and the other to control the

process by measuring average values. The argument for using peaks seems sensible and useful because it is quite easy to relate peak-to-valley measurements of a profile to variations in the straightness of interferometer fringes, so there was, in essence, the possibility of a traceable link between contact methods and optical ones. This is the approach used by the USSR and Germany in the 1940s and until recently. The UK and USA, on the other hand, realized at the outset that the measurement of peak parameters is more difficult than measuring averages, and concentrated therefore on the latter which, because of their statistical stability, became more suitable for quality control of the manufacturing process. Peak measurements are essentially divergent rather than convergent in stability; the bigger the length of profile or length of assessment the larger the value becomes. This is not true for averages; the sampled average tends to converge on the true value the larger the number of values taken.

The formal approach to statistical reliability will be left to chapter 5. However, in order to bring out the nature of the characterization used in the past it is necessary to point out some of the standard terms governing the actual length of profile used. The basic unit is the *sampling length.* This is the length of assessment over which the surface roughness can be considered to be representative. Obviously the application of such a definition is fraught with difficulty because it depends on the parameter and the degree of confidence required. For the purpose of this subsection the length will be assumed to be adequate—whatever that means.

The value of the sampling length is a compromise. On the one hand, it should be long enough to get a statistically good representation of the surface. On the other, if it is made too big longer components of the geometry, such as waviness, will be drawn in if present. The concept of sampling length therefore has two jobs, not one. For this reason its use has often been misunderstood. It has consequently been drawn inextricably into many arguments on reference lines, filtering and reliability. It is brought in at this stage to reflect its use in defining parameters. Sometimes the instrument takes more than one sampling length in its assessment and sometimes some of the total length traversed by the instrument is not assessed for mechanical or filtering reasons, as will be described in chapter 4. However, the usual sampling length of value 0.03 in (0.8 mm) was chosen empirically in the early 1940s in the UK by Rank Taylor Hobson from the examination of hundreds of typical surfaces [2]. Also, an *evaluation length* of nominally five sampling lengths was chosen by the same empirical method. In those days the term 'sampling length' did not exist; it was referred to as the meter cut-off length. The reason for this was that it also referred to the cut-off of the filter used to smooth the meter reading. Some typical sampling lengths for different types of surface are given in tables 2.1 to 2.4.

**Table 2.1** Gaussian filter 1995
Sampling lengths for $R_a$, $R_z$ & $R_y$ of periodic profiles

| $S_m$ mm | | | |
|---|---|---|---|
| Over | Up to (inclusive) | Sampling length mm | Evaluation length mm |
| (0.013) | 0.04 | 0.08 | 0.4 |
| 0.04 | 0.13 | 0.25 | 1.25 |
| 0.13 | 0.4 | 0.8 | 4.0 |
| 0.4 | 1.3 | 2.5 | 12.5 |
| 1.3 | 4.0 | 8.0 | 40.0 |

In tables 2.2, 2.3, and 2.4 some notation will be used which is explained fully later. See glossary for details.

**Table 2.2** Roughness sampling lengths for the measurement of $R_a$, $R_q$, $R_{sk}$, $R_{ku}$, $R\Delta q$ and curves and related parameters for non-periodic profiles (for example ground profiles).

| $R_a$ $\mu$m | Roughness sampling length $lr$ mm | Roughness evaluation length $ln$ mm |
|---|---|---|
| $(0.006) < R_a \leqslant 0.02$ | 0.08 | 0.4 |
| $0.02 < R_a \leqslant 0.1$ | 0.25 | 1.25 |
| $0.1 < R_a \leqslant 2$ | 0.8 | 4 |
| $2 < R_a \leqslant 10$ | 2.5 | 12.5 |
| $10 < R_a \leqslant 80$ | 8 | 40 |

**Table 2.3** Roughness sampling lengths for the measurement of $R_z$, $R_v$, $R_p$, $R_c$ and $R_t$ of non-periodic profiles (for example ground profiles).

| $R_z$[1] $R_z1$ max.[2] $\mu$m | Roughness sampling length $lr$ mm | Roughness evaluation length $ln$ mm |
|---|---|---|
| $(0.025) < R_z, R_z1$ max. $\leqslant 0.1$ | 0.08 | 0.4 |
| $0.1 < R_z, R_z1$ max. $\leqslant 0.5$ | 0.25 | 1.25 |
| $0.5 < R_z, R_z1$ max. $\leqslant 10$ | 0.8 | 4 |
| $10 < R_z, R_z1$ max. $\leqslant 50$ | 2.5 | 12.5 |
| $50 < R_z, R_z1$ max. $\leqslant 200$ | 8 | 40 |

1) $R_z$ is used when measuring $R_z$, $R_v$, $R_p$, $R_c$ and $R_t$.
2) $R_z1$ max. is used only when measuring $R_z1$ max., $R_v1$ max., $R_p1$ max., $R_c1$ max. and $R_t1$ max.

**Table 2.4** Roughness sampling lengths for the measurement of $R$-parameters of periodic profiles, and $RSm$ of periodic and non-periodic profiles.

| $RS_m$ $\mu$m | Roughness sampling length $lr$ mm | Roughness evaluation length $ln$ mm |
|---|---|---|
| $0.013 < RS_m \leqslant 0.04$ | 0.08 | 0.4 |
| $0.04 < RS_m \leqslant 0.13$ | 0.25 | 1.25 |
| $0.13 < RS_m \leqslant 0.4$ | 0.8 | 4 |
| $0.4 < RS_m \leqslant 1.3$ | 2.5 | 12.5 |
| $1.3 < RS_m \leqslant 4$ | 8 | 40 |

The usual spatial situation is shown in figure 2.5.

Most amplitude parameters are referred to one sampling length in the standards although, in practice, more than one is used. Some specific typical parameters will now be given in more detail.

Attempts to quantify the roughness height have historically been broken into two camps: one measuring peak-to-valley heights, the other average heights. Peak measurement was advocated by the Germans, French and other Europeans in an effort to link surface measurements obtained from stylus instruments to fringe

deviations obtained using interferometers [3,5]. Although a very laudable concept, it ran into difficulties because of the difficulty of measuring or even finding the highest peak-to-valley height on the surface. The British and Americans adopted an average height measurement philosophy in order to get a workable parameter. It should also be realized that from the outset, many firms developed their very own surface parameters which sometimes bore little resemblance to those suggested as preferred parameters by the ISO and national institutions. These individual parameters were devised to satisfy an in-house need and as such were acceptable. The problem arose when subcontractors were expected to buy equipment to measure the same parameters. This situation was an instrument-maker's nightmare until the more versatile computer systems emerged. In fact, it was one of the reasons for developing computer systems along with the need to measure more complicated parameters for investigating functional behaviour.



**Figure 2.5** Length measures on a profile.

### 2.1.1.1 Amplitude parameters

Figure 2.2(*b*) shows a typical profile trace condensed in the horizontal direction for convenience of display [4]. Perhaps the most obvious way of putting a numerical value to its magnitude is to measure the heights of the peaks and valleys. The level from which to measure is taken here to be the bottom of the chart. However, this assumes that the profile is reasonably level relative to the bottom of the chart. One such parameter is shown as $R_t$ in figure 2.6. Incidentally, none of the definitions are sacrosanct—they keep changing.



**Figure 2.6** Amplitude parameters—peaks $R_t$.

This parameter is simply the difference in level between the highest peak and the lowest valley in the length of profile usually called the sampling length. A similar parameter called $R_z$ is shown in figure 2.7.



**Figure 2.7** Peak parameters—$R_z$.

The formula for $R_z$ is

$$R_z = \frac{(P_1 + P_2 + P_3 + P_4 + P_5) - (V_1 + V_2 + V_3 + V_4 + V_5)}{5}$$

$$= \frac{\sum_{i=1}^{5} P_i - \sum_{i=1}^{5} V_i}{5} \tag{2.1}$$

$R_z$ is, in effect, the average of the height difference between the five highest peaks and the five lowest valleys. The reason for taking an average value of peaks is to minimize the effect of unrepresentative peaks or valleys which occasionally occur and can give an erroneous value if taken singly. $R_z$ is used often without reference to a mean line. It has also been used to get some idea of the surface roughness on very short lengths of surface such as might be found on a lip or shoulder where not even one sampling length of surface is available. $R_z$ is an insurance against freaks. There is another way of insuring against the freak occurrence. This is given the name Rautiefe. It was originally proposed by Swedish engineers. Its derivation is shown in figure 2.8. It is the separation of two parallel lines cutting through the profile such that the upper one is in metal for 5% of its path and the lower one is in metal for 95% (or 90%) of its path. This parameter not only attempts to relate the profile to some semblance of mechanical function, but also contains the germ of the more advanced statistical methods described in section 2.1.2. It suffers from the drawback that it cannot be measured from the chart of the profile as easily as $R_t$ or $R_z$.



**Figure 2.8**   Peak parameter — Rautiefe.



**Figure 2.9**   Peak parameter — levelling depth $R_p$.

A slightly different peak measure is $R_p$, the levelling depth (figure 2.9). This is the average depth of the profile below a reference line drawn through the highest peaks taken through the profile.

How this reference line is formally positioned relative to the profile will become clear later on. Other definitions of $R_p$ exist. These depend on the precise form of definition of the reference line.

Another definition of $R_p$ can be given simply as the maximum peak height from a mean line positioned in the middle of the profile such that the area above the line is equal to that below the line. A corresponding maximum valley parameter $R_v$ can also be defined from this line.

Sometimes the peak-to-valley measurement has been limited to a single sample length, such as $R_p$, but another assessment has been used, called $R_{tm}$, in which some averaging of peaks is achieved by taking the $R_t$

values from a succession of adjacent (contiguous) sampling lengths and averaging them. Typically, five sampling lengths are used for the reasons given earlier (figure 2.10). Thus

$$R_{\text{tm}} = \frac{\sum_{i=1}^{5} R_{\text{t}i}}{5} .$$

(2.2)

This method, although different from $R_z$, illustrates the quandary faced by metrologists determined to use peak-to-valley or extrema criteria for the evaluation of surface features. Some degree of averaging has to be incorporated in order to make the method workable and capable of being calibrated. Many other methods have been tried, such as $R_{3z}$, the difference in height between the third highest peak and the third lowest valley within a sampling length. This has been used quite extensively in Japan.



**Figure 2.10** Ten-point height $R_z$.

Basically, at this level of technology, from a functional point of view, it does not really matter which peak or average peak parameter is used as long as it is consistently used and the definition agreed between the various people using the technique. Many problems arise because definitions that are assumed to be fixed are found to be different at a later date, due, perhaps, to a change in a standard.

Other attempts to bridge the gap between peak, average peak and average values will be discussed after a brief description of the average parameter.

The most obvious method of increasing the reliability of the amplitude measure involves using all of the profile signal rather than just the maximum and minimum values. The best known and most often used today is $R_a$, defined relative to a mean reference line. It is the mean departure of the profile from the reference line. This is shown in figure 2.11.



**Figure 2.11** Average parameters—$R_a$ and $R_q$.

Thus, if $z = f(x)$ is the profile measured from the reference mean line and $L$ is the length of the profile being assessed (this can be the sampling length), then $R_a$ is defined by

$$R_a = \frac{1}{L} \int_0^L |z| \, \mathrm{d}x .$$

(2.3)

Another parameter shown in figure 2.11 is $R_q$, the RMS deviation defined again relative to a mean line as

$$R_q = \sqrt{\frac{1}{L} \int_0^L z^2 \mathrm{d}x}. \tag{2.4}$$

Although $R_q$ is more statistical in nature it is difficult to measure from a chart. Early measurement in the USA was in fact an RMS value. The AA (arithmetic average) on drawings invariably implied measurement of an RMS value and dividing by a fixed constant of 1.11. This gave the true AA value for a sine wave but no other waveform.

Other names have been given to the $R_a$ value, namely the CLA (centre line average) value in the UK. Because of its wide acceptance and its usefulness it seems likely that the $R_a$ value will continue to hold its place as the foremost amplitude measurement. However, this does not imply that other mean value parameters are not forthcoming. One measure is called the $R$ value, which is derived by drawing two lines through the peaks and valleys respectively. This method is also called the 'motif' method and is interesting in that there are certain allowances in the rules for constructing both the superior envelope through the peaks and the inferior envelope through the valleys to preclude certain values, as shown in figure 2.12.



**Figure 2.12** Motif method—$R$ and $W$.

$R$ is taken to be the mean of several values of the separation of the envelopes. Mention of these envelopes will also be made in the next section on waviness. Also, the use of envelopes to provide boundaries from which height measurements can be made will be discussed in chapter 5. Parameters such as $R$ have only been in vogue since the advent of digital methods for the simple reason that exclusion conditions for some of the peaks and valleys are difficult to implement accurately with analogue devices.

The final conventional parameter to be discussed in this section is, by its nature, statistical. However, because it was introduced rather early in the natural evolution of the subject [8], it will be included in this section. This curve is fundamental to surface metrology and will keep arising in the text. This curve is the material ratio curve.

In figure 2.8 the Rautiefe [3] is shown defined relative to the air/metal ratio. A plot of the ratio of air to metal starting at the highest peak is referred to as the bearing ratio, material ratio or Abbott–Firestone curve (figure 2.13). In this book the term material ratio will be used as the preferred term.

The material ratio at height $z$ $MR(z)$ is given by $\sum t_p(z)/L = MR(z)$, where $\sum t_p(z)$ is the sum of the material $t$ found at level $z$ over a length $L$.

The height distribution function and the amplitude probability density function can be derived from this curve. Such transpositions and measurements taken from them have been proposed by Pesante [9] and Ehrenreich [10]. A great deal of additional interest in this parameter has been generated with the emergence of multiprocesses (or stratified processes). Note that at any one level the value of the material ratio curve (bearing ratio) is a ratio and not a length. It would give the same shaped curve if the horizontal scale of the profile were shrunk or expanded. Indeed it would give the same curve if the order of the $t_p$ values were exchanged or even pushed up against each other; the actual horizontal scale or order of the peaks horizontally is not relevant. A more detailed discussion of this method of assessing the material ratio curve is given in section 2.1.7.5.

**Figure 2.13** (*a*) Material ratio curve; (*b*) $R_{pk}$, $R_{vk}$, $R_k$.

It can be seen from this introduction that there is an additional criterion to be used for parameters. This concerns the ease of graphical measurement. It seems that any parameter which can be easily measured is not very reliable and *vice versa*. This is shown in table 2.5.

**Table 2.5**

| Parameter | Ease of measurement | Reliability |
|---|---|---|
| $R_q$ | Very difficult | Very high |
| $R_a$ | Difficult | High |
| $R_{tm}$ | Moderate | Moderate |
| $R_t$ | Easy | Low |
| $R_{max}$ | Very easy | Very Low |

Most of the parameters considered so far are very much a result of historical development and not very significant in themselves. For example, from the point of view of statistical importance, $R_q$ is much more valid than $R_a$ but in the early days it could not readily be measured. Today this is not true, but still the old parameter exists. The real exception to this rule is the material (bearing ratio) curve which existed almost at the start of roughness characterization and in fact contained many of the best features of any parameter; it is functional and reliable. It also has a reasonable theoretical basis. The parameters described in this book are nothing like exhaustive, neither should they be. In a healthy subject the parameters need to change to reflect circumstances.

As if this is not confusing enough some countries try to classify the parameters in terms of size. One typical classification is shown in Figure 2.14. This is called the N system. Another system uses a set of del symbols ∇∇∇ etc. The problem here is that standard documents can become cluttered up with what are in effect 'in-house' characterization methods used by large companies. Table 2.6 shows the hopeless confusion for some countries. High N means rough for some countries yet for others it means smooth. Unfortunately the N numbers are still used! Great care has to be taken to make sure which N scale is being used! Also remember that the letter N is used for sampling lengths in the assessment length.

**Table 2.6** N values

| CLASS | UK | USA | GERMANY | USSR | JAPAN | ISO |
|---|---|---|---|---|---|---|
| 1 | 1 | 0.25 | 400/630 | 200 | 200/280 | 200 (R) |
| | | 0.5 | | | | |
| | | 1.0 | | | | |
| 2 | 2 | 2 | 160/250 | 125 | 100/140 | 125 |
| 3 | 4 | 4 | 63/100 | 63 | 50/70 | 63 |
| 4 | 8 | 8 | 25/40 | 40 | 25/35 | 40 |
| 5 | 16 | 16 | 16 | 6.3 | 18 | 6.3 |
| 6 | 32 | 32 | 6.3 | 3.2 | 12 | 3.2 ($R_a$) |
| 7 | 63 | 63 | 4.0 | 1.6 | 6 | 1.6 |
| 8 | 125 | 125 | 2.5 | .8 | 3 | .8 |
| 9 | 250 | 250 | 1.6 | .4 | 1.5 | .4 |
| 10 | 500 | 500 | 1 | .2 | .8 | .2 |
| 11 | 1000 | 1000 | .63 | .1 | .4 | .1 |
| 12 | - | - | .4/ .25 | .05 | .2 | .05 |
| 13 | - | - | .16/ .1 | .12 | .1 | .025 |
| 14 | - | - | .06/ .04 | .06 | - | .012 |
| Unit | | | | | | |
| Standard | BS1134 | B46 | 4763 DIN | GOST2780 | JIS | |
| | 1950 | 1955 | 1954 | 1951 | 1955 | 1953 |

**Table 2.7**

| Turning | T | Boring | B |
|---|---|---|---|
| Diamond Turning | DT | Reaming | R |
| Grinding | G | Milling | M |
| Honing | H | Planing | P |
| Lapping | L | Scraping | S |
| Polishing | Po | Broaching | Br |

Obviously this shorthand is different in different languages.

In some documents the finish is referred to by process. Table 2.7 gives a few typical examples.

### 2.1.1.2    Spacing parameters

The height parameters, such as $R_a$ and $R_t$, have usually been specified on the drawing in an attempt to control the manufacturing process. To a large extent this in turn has been successful in controlling the performance of the part. Nowadays, however, with increasing efforts to optimize the function, the functional performance of a workpiece as a requirement has been growing. One way of supplementing the height information is to provide some index of crest spacings on the surface. As in the case of the height parameters many propositions have been made. One of the simplest, called the high spot count, is a measure of the number of crossings of the profile through the mean line $z = 0$. Various alternatives to the mean line crossings have been suggested. For instance, an index can be taken to be the number of crossings of the profile by a line parallel to the mean line but displaced vertically by one $R_a$ value. Much use has been made of this concept in the statistical methods described in the section 2.1.3 on statistical parameters.

Figure 2.14 Roughness grades.

Whereas the high spot count is referred to at a given height from the reference line within the profile, the peak count is not. It is just a measure of maxima irrespective of height. This count is very easy to measure from the record of the profile and hence it has had widespread use, especially in the sheet steel industry.

The big disadvantage of this method is that it is very dependent on the frequency response of the recording apparatus, a point which will be expanded upon later. The effect of this is to make comparisons between peak counts very difficult.

One technique which has been used in the past to minimize this effect is to have some form of amplitude discrimination. For instance, a peak is not registered unless the previous valley lies 50 microinches below it. This approach has also been used in the steel industry in the USA (figure 2.15). Various other methods of discrimination have been tried but they do and should relate to the particular function.

The peak density, the zero-crossing density and the peak count, in which an amplitude discriminating band is used, are shown in figures 2.15–2.17.

Later on it will be demonstrated that the apparently simplistic methods of assessing surfaces shown here should not be treated with contempt. Many current theories make use of a number of these concepts; for

**Figure 2.15** Other peak parameters—amplitude discrimination.



**Figure 2.16** Length parameters—peak density.



**Figure 2.17** Length parameters — zero crossings.



**Figure 2.18** Length parameter—peak count (amplitude discriminator).

example, the ratio of the high spot count to the peak count can be used as a measure of the process. This concept is used extensively in mathematical modelling to which a lot of importance will be attached.

Other names for peak spacing parameters have been advocated in international and national standards and, in particular, VNIIMS Moscow. [11]. One is the $S$ value which is the average distance between local (small) peaks (figure 2.16) and $S_m$ the average distance between positive crossings of the profile with the mean line (figure 2.17). The ratio of $S/S_m$ has been used in the former USSR for characterizing surfaces. For a very random surface this ratio is low, yet for a deterministic one it can approach unity. Both $S$ and $S_m$ are meant to be related to one sampling length. Other peak counts incorporate an amplitude discriminating factor (figure 2.18). The value of this factor depends on the application.

There is yet another wavelength parameter [12] which is almost the average distance between peaks—but not quite! This is the average wavelength of the surface. It takes into account the size of the peaks and the basic sinusoidal wavelengths making up the waveform. It will only be mentioned here but will be considered in more detail in a later section, 2.1.3.4. However, it is useful to introduce it because it indicates that some parameters are best expressed as a mixture of the height and wavelength dimensions; these are called hybrid parameters and some will be described here.

### 2.1.1.3 Hybrid parameters

Other suggested parameters have incorporated both height and spacing information. These parameters have usually been derived from the differentials of the waveform. A number of investigators have found this type of parameter useful, particularly in friction wear and optical behaviour.

Spacing and hybrid parameters were initially introduced because of the inability to specify the texture effectively by means only of height parameters. Their widespread use, however, was severely curtailed because of the inability to measure them properly. There were many attempts to use the hybrid parameters to characterize surfaces, notably by Myers [13]. The difficulty arose in measuring differentials by analogue means. By their very nature differentials are noisy. This meant that many early attempts were unreliable and so, instead of clarifying the characterization, they often clouded it. This situation has largely been overcome by the use of digital methods but it should be emphasized that great care still has to be taken.

Some examples are the average ($\Delta_a$) or RMS ($\Delta_q$) slope, the curvature ($\rho$), peak curvature, etc.

Thus the average slope ($\Delta_a$) is given by

$$\Delta_a = \frac{1}{L} \int_0^L \left| \frac{dz}{dx} \right| dx \tag{2.5}$$

where $dz/dx$ is the instantaneous slope of the profile. These parameters are fundamental in function and should be remembered here for reference in chapter 7 on function. The RMS slope ($\Delta_q$) is defined as

$$\Delta_q = \left[ \frac{1}{L} \int_0^L \left( \frac{dz}{dx} \right)^2 dx \right]^{1/2}. \tag{2.6}$$

Note that $\Delta_a$ and $\Delta_q$ are the same to slope measurement as $R_a$ and $R_q$ are to amplitude and $\lambda_a$ and $\lambda_q$ to spacing. The average curvature of the profile $\rho$ is given in equation (2.7). Signs are ignored.

$$\rho = \frac{1}{L} \int_0^L \frac{d^2z/dx^2}{[1 + (dz/dx)^2]^{3/2}} dx. \tag{2.7}$$

This formula can be simplified if account is taken of the fact that $dz/dx$ is usually small on a surface. (This is a result which is true in certain cases only, particularly rough surfaces.)

Thus

$$\rho \simeq \frac{1}{L} \int_0^L \left| \frac{d^2z}{dx^2} \right| dx. \tag{2.8}$$

Curvature is the reciprocal of the radius of the profile at any point. Other variations of this definition have been proposed. For instance, the average peak curvature ($\rho_\Delta$), the average valley curvature ($\rho_\nabla$), etc. As for peaks and valleys these definitions have to be qualified somewhat to minimize the errors due to variations in resolution, frequency response, etc. These curvature formulae are considered to be especially important in contact theory.

Another parameter sometimes used is the profile length $l_p$ (figure 2.19). This is given by



Figure 2.19 Profile length.

$$l_{\mathrm{p}} = \frac{1}{L} \int_0^L \left| \left( \frac{\mathrm{d}z}{\mathrm{d}x} \right)^2 + 1 \right|^{1/2} \mathrm{d}x \tag{2.9}$$

which for small $\mathrm{d}z/\mathrm{d}x$ reduces to

$$l_{\mathrm{p}} \simeq \frac{1}{L} \int_0^L \left[ 1 + \frac{1}{2} \left( \frac{\mathrm{d}z}{\mathrm{d}x} \right)^2 \right] \mathrm{d}x$$

$$\simeq 1 + \frac{\text{mean square slope}}{2}. \tag{2.10}$$

The directionality $D$ of the profile can also be defined using the differential parameters. Thus

$$D = \frac{L_1 - L_2}{L} \tag{2.11}$$

where $L_1$ is the length of profile along the horizontal scale having positive slope and $L_2$ is the length having a negative slope.

Interest has been shown in this directionality parameter in the context of journal bearings. Considerable differences in wear properties have been encountered depending upon how the directionality of the roundness profile of the shaft compares with the direction of rotation of the bearing. In this case the directionality of the circumferential surface roughness is the important factor.

The next section is concerned with methods of separating waviness from roughness. Part of the section is devoted to the theory of linear filtering, which can be omitted on a first reading. However, it is necessary to introduce the concepts here because so much of what follows is dependent upon them. Filtering is a form of spatial characterization. The filtering concepts discussed in this section should be linked with section 2.2 on waviness for completeness.

It should not be thought that the establishment of reference lines and the definition of so-called sampling lengths are purely instrumental considerations. They are not. These issues constitute a very important aspect of characterization of the surface into its most significant components. Often they provide the means whereby specific features can be highlighted in space. This is why they are included here; the reference helps to characterize the roughness. However, there are many occasions when the reference itself needs to be classified. This is when the reference line is considered to contain the waviness. In these cases the reference line has to be classified in much the same way as the roughness has been considered.

Waviness characterization is considered in section 2.2. Both roughness and waviness characterization appear in chapter 3 on processing and chapter 4 on instrumentation.

Reference lines therefore seem to occur all over the place! This is inevitable because the numerical values of many of the parameters depend on which one, if any, is used. The problem is this: surface roughness is a deviation—it needs a reference from which to quantify it. Ideally the external mechanical reference provided within the instrument should be positioned within the roughness marks. This is obviously impossible. All it can do is to enable, by means of the pick-up and transducer, the generation of the surface data, that is the profile. It is then necessary to use some mathematical procedure to position the reference within the profile from where the roughness can be measured. Because the problems associated with reference lines occurred historically before even areal considerations or random processes became important, they have to be considered here. Their involvement is deeply entwined with the development of the subject as a whole and is just as important today as it was when surface metrology was in its infancy.

## 2.1.2 Reference lines

Many of the parameters described above depend on the definition of some reference line from which to measure. This reference line has to have a fixed known height relative to the roughness, it has to have a pre-determined shape and its length must be known. As for the determination of the sampling length, this is a basic device for limiting the length of surface over which an individual measurement of roughness is made, as shown in figure 2.5.

If the roughness marks, that is the process marks, are repetitive, then the sampling length should include about 10 whole marks within it. This criterion was originally established empirically but it has stood the test of time and is still used. The early investigations led to the standard sampling length of 0.8 mm (0.03 in) being adopted as standard, this being about 10 times larger than the typical feed rate of turning in the 1940s [2].

For random marks such as found in grinding, the determination of the sampling length is more difficult because of the lack of obvious cycles of repetition. However, a rule of thumb can be given here and verified later, which is that if a line is drawn by eye through the profile to include the roughness marks only then the sample length should include between 20 and 30 crossings of the profile with the line. Small deviations of the reference line from being straight or misplaced vertically are not important. This corresponds to having a high spot count of about 10. Formally, the determination depends upon a knowledge of the statistics of the surface, which will be considered in the next section. The term sampling length used to be called the 'cut-off length' or sometimes the 'meter cut-off'.

In the past, by convention, there was a selection of sampling lengths available for use to suit whichever process was being measured.

For convenience and to provide a basis for instrumentation a series (originally) based on √10 differences has been standardized in the British Standard BS 1134, the US Standard B46 and ISO Recommendation R468. The accepted values with the Imperial equivalent are as follows:

| 0.08 | 0.25 | 0.8 | 2.5 | 8.0mm |
|------|------|-----|-----|-------|
| 0.003 | 0.01 | 0.03 | 0.1 | 0.3 in |

(0.03 in = 0.762 mm rounded to 0.8mm).

It should be said here that the specific values were chosen for convenience and are by no means sacrosanct. For many of the new machining processes it could be argued that any suitable sampling length could be chosen.

When the sampling length has been determined the reference has to be positioned in the profile in such a way that departures from it represent the roughness. For many parameters five sampling lengths are used (figure 2.20).



**Figure 2.20** Sampling length.

Four main types of line have been investigated in the past all over the world and each incorporates some wavelength discrimination such as the sampling length limitation. These methods have been achieved in a

number of different ways ranging from graphical, electrical, mechanical and now computing. Nowadays the graphical methods are rarely used but are held in reserve.

The four methods most used are straight lines, polynomials, low-pass filters and envelopes (or motifs).

### 2.1.2.1 Straight lines

Perhaps the simplest form of reference line that can be put onto a chart is a straight line extending over the sampling length. It can be positioned so that it touches the two highest crests within the sampling length as for measuring $R_p$, the levelling depth in figure 2.9, or derived similarly from the deepest valleys or it may be a mean line, in which case the profile area enclosed above the line is equal to that below it. Note that the mean line is often referred to as the centre line.

The problem is how to determine the most meaningful angle and position at which to draw the line relative to the profile; a number of lines having the equal-area criterion are possible. One traditional method is to qualify the definition of equal area with a statement that the line must follow the general direction of the surface. This rather subjective method has the effect of tending to minimize the modulus of the area enclosed between the profile and the line therefore minimizing the $R_a$ value. Another method which is not the same is that of applying a constraint to the slope and position by specifying a 'best-fit' least-squares straight line through the data. In this case the line will be unique, free from subjective positioning. From the point of view of simply fitting a best line this method is very good but it has two problems associated with it. One is that it is difficult to determine from a graph and the second is that it sometimes produces unrealistic results.

Whereas the former method of fitting a straight line is difficult to quantify theoretically the latter is not.

If the measured values of $z$ as a function of $x$ on the chart are compared with the values expected from a best-fit line $z'$ at the same $x$ then the problem is to minimize the sum $S_u$ where

$$S_u = \int_0^L \left(z - z'\right)^2 dx. \tag{2.12}$$

Strictly the deviations should be in a direction normal to the best-fit line itself, that is $z \cos \alpha - z' \cos \alpha$, as shown on the graph in figure 2.21. Thus $S_u'$ is minimized where

$$S_u' = \int_0^L \left(z - z'\right)^2 \cos \alpha^2 \, dx \tag{2.13}$$

$$S_u' = \int_0^L \left[z(x) - \left(mx + c\right)\right]^2 \cos^2 \alpha \, dx \tag{2.14}$$



**Figure 2.21**  Best-fit least-squares lines.

which, given that m = tan$\alpha$, is that

$$\int_0^L \left[ z(x)\cos\alpha - x\sin\alpha - c\cos\alpha \right]^2 dx \qquad (2.15)$$

is minimum, where $m$ is the tangent and $c$ is the intercept of the line on the $z$ axis, from which $\alpha$ is obtained as

$$2\alpha = \tan^{-1} 2\left( \frac{\int_0^L z(x)x\,dx - (1/L)\int_0^L x\,dx\int_0^L z(x)\,dx}{\int_0^L x^2\,dx - (1/L)(\int_0^L x\,dx)^2 - \int_0^L z(x)^2\,dx + (1/L)(\int_0^L z(x)\,dx)^2)} \right) \qquad (2.16)$$

which is in effect

$$\tan 2\alpha = 2\left( \frac{\text{covariance of } x \text{ and } z}{\text{variance of } x - \text{variance of } z} \right). \qquad (2.17)$$

(Both $\alpha$ and $c$ or $m$ and $c$ are obtained by setting $\partial S'_u / \partial\alpha$, $\partial S'_u / \partial c = 0$ or $\partial S_u / \partial m$, $\partial S_u / \partial c = 0$ and solving.)

This formula should be used if the horizontal and vertical magnifications are equal, as they often are in flatness measurements. In the case of roughness the vertical magnification is often 50 times bigger than the horizontal so that the direction of the deviation is distorted from the normal to the line and towards the normal to the base line. Under these circumstances equation (2.16) should be used to give

$$m = \frac{\int_0^L xz(x)\,dx - (1/L)\int_0^L x\,dx\int_0^L z(x)\,dx}{\int_0^L x^2\,dx - (1/L)(\int_0^L x\,dx)^2}. \qquad (2.18)$$

Digital equivalents of equations (2.17) and (2.18) will be given in section 3.6.14.

The advantage of using a linear mean line as opposed to a linear crest or valley line is that poor positioning of these latter lines can easily result if freak asperities are present in the profile.

Despite the undoubted advantages of the best-fit line approach it does run into difficulty when periodic waveforms are being measured, for instance a sine wave. The question arises as to what constitutes a realistic mechanical direction of the line.

In the case of the centre line, the difficulty arises only when the sample contains a single cycle of the waveform, for when there are two or more cycles, the crest line at once determines the inclination which leads to minimum values.

In the case of the least-squares mean line, however, the inclination is determined both by the starting point and by the number of cycles included in the sample.

Consider figure 2.22 where the full curve represents a single sample of a sinusoidal profile.

The best-fit least-squares line has a different slope depending upon its starting point in the cycle. The intercept of the line with the starting ordinate axis may easily be shown to be a function of the number of cycles within the sample.

For $n$ samples in the cycle the intercept $c$ as a function of the amplitude of the wave itself $A$ is given by

$$\frac{c}{A} = \frac{3}{\pi n}. \qquad (2.19)$$

**Figure 2.22** Least-squares line: $n$ is number of cycles in sampling length.

Obviously for $n = 1$ this angle is enormous and unrealistic; the mean line has to be parallel with the crests or valleys of the sine wave. This would have been the case if a cosine rather than a sine had been selected. So the least-squares line for a periodic wave has a slope dependent upon the starting point! This is nonsense, so in order to make the definitions (2.17) and (2.19) usable $n$ has to be large. For $n = 10$, $c/A < 10\%$, which is deemed reasonable and constitutes one reason why the 10 cycles of a periodic function are used as some basis for the determination of the sampling length in terms of machine tool feed.

The use of the straight reference line fitted within each of the individual sample lengths making up the evaluation length, whether by least squares or otherwise, has the advantage of simplicity but suffers from the fact that, at the end of each sampling length, the end of one line inevitably does not contact the start of the next. This means that a discontinuity is introduced into the waviness shape, both in level and in slope at each sampling length. This can cause problems in measuring the parameters considered in section 2.1. The errors are more severe for the peak parameters than for the average ones.

It has been known for a best-fit least-squares line to be used over the whole evaluation length and not be used for each individual sampling length making up the whole assessment. When this method is used obviously only the tilt of the instrument datum relative to the workpiece can be removed. Errors of form or curvature of the part cannot be removed at the same time. If, at the same time as removing tilt it is considered advantageous to remove the form, another method has to be used.

### 2.1.2.2    *Polynomial fitting*

It may be that the form can be adequately represented over the evaluation length by a cubic or quadratic equation. This would give a better reference than a straight line (figure 2.23). Then a best-fit polynomial could be used of the form $z = a_1 + a_2 x + a_3 x^2 + \ldots$ where $a_1$, $a_2$, $a_3$ are the least-squares coefficients as $m$ and $c$ were for the least-squares line.



**Figure 2.23** Best-fit limitation.

In this case they can be extracted from the equations

$$a_1 L + a_2 \int_0^L x\,\mathrm{d}x + a_3 \int_0^L x^2 = \int_0^L z(x)\mathrm{d}x$$

$$a_1 \int_0^L x\,\mathrm{d}u + a_2 \int_0^L x^2\mathrm{d}x + a_3 \int_0^L x^3\mathrm{d}x = \int_0^L xz(x)\mathrm{d}x$$

$$\vdots$$

$$a_1 \int_0^L x^{n-1}\,\mathrm{d}x + a_2 \int x^n\mathrm{d}x + a_3 \int x^{n+1} + \ldots = \int_0^L x^{n-1} z(x)\,\mathrm{d}x \qquad (2.20)$$

where $n$ is the order of the polynomial.

The big advantage of this polynomial approach is that all the data in the evaluation length can be used. There are no end effects! In other words one of the great advantages of the use of polynomials to represent tilt, waviness and form error through the roughness data is that all the data is utilized. The penalty for this, however, is that to get the most realistic reference line the order of the polynomial should be made to 'match' the geometry of the profile itself. For example, if the surface is aspheric. As an example of mismatching, consider a workpiece which is perfectly smooth and has a form error of three-quarters of a sine wave. If a sinusoidal waveform is fitted to the form error without knowing that only three-quarters of the full sine wave are present then considerable errors can result.

The imposition of an arbitrary period to the workpiece based only upon its length can actually introduce harmonic distortion as is shown in figure 2.24. It is only acceptable to make an error of judgement like this when there are considerably more than one cycle of the natural wave within the length of the profile. Under these conditions the errors are reduced to an acceptable level.



**Figure 2.24** Distortion-mismatched length.

There is an alternative to the least-squares or best-fit approach and this is that of minimum divergence of the reference curve from the test piece—given that the order (e.g. quadratic or cubic) is known. The obvious example of this sort of approach is to use the Chebychev polynomials; these, however, could generate a reference curve which, unless care is taken, will suffer from the same freak behaviour as is possible with crest or valley lines.

One very interesting feature of this method, however, is that some attempt can be made by using the Chebychev series to characterize fully the whole profile, that is roughness, waviness, errors of form and set-up errors, using just one type of function. This assumes that the idea of looking at the whole evaluation

length in one go, rather than breaking it down into sampling lengths, has been abandoned—which it does not necessarily have to be. What can be done to remove the discontinuities that occur between straight lines from each sampling length? A number of attempts have been made. The first is due to Reason [6] who postulated that instead of using a succession of adjacent sampling lengths with lines through each, a multiplicity of sampling lengths of the profile, staggered with respect to each other and consequently over-lapping, should be used. The reference line should then be made up from the locus of the mean profile within each sampling length and plotted at the mid-point. This method has therefore been referred to as the 'mid-point locus' reference line. It has the considerable advantage that it is unique and gives the same answer as if the centre points on centre lines or least-squares lines within the samples had been connected up as shown in figure 2.25.



**Figure 2.25** Mid-point locus line.

Another advantage of this method is that the centre point within each sampling length can easily be obtained from a graph simply by measuring the area between the profile and the bottom of the chart, dividing by the sampling length, and plotting the value in the middle of the sample. The accuracy of the locus is dependent on the extent to which sampling lengths of the profile are allowed to overlap, so obviously there must be a compromise between the amount of labour involved and the accuracy. Five shifts of the averaging box within a sampling length have proved to be reasonable. The shifts are not necessarily equally spaced.

The only disadvantage of this method is in its transmission characteristics: they are worst for periodic signals as seen in figure 2.26. For example, if the locus is plotted on the sine wave then it will undulate unrealistically if the sampling length is not an integral number of wavelengths [15].



**Figure 2.26** Mid-point locus characteristics.

It can be shown that if the amplitude of undulation of the locus is taken as a measure of performance and if it is plotted against the wavenumber (not wavelengths) of different sine waves, the resulting curve is that of the sine function $S(v)$ where $v$ is the ratio of the wavelength to sampling length:

$$S(1/\nu) = \frac{\sin\left[2\pi(1/\nu)\right]}{2\pi(1/\nu)} \tag{2.21}$$

which for $1/\nu$ having a value of 0.66, for instance, gives an undulation amplitude of $-19\%$ (see figure 2.26). Usually the errors are much smaller, especially for random waves. This undulation in the mean line for waves shorter than the sample length looks mechanically unrealistic. In any technique to determine a mean line it is preferable that waves shorter than the sampling length are not distorted.

### 2.1.2.3  Spline functions

Another possible way has been suggested which could reduce the shortcomings of the standard graphical method. In effect what this method does is to apply constraints to the behaviour of the reference line at the boundary of each sampling length. Discontinuities of neither position nor slope are allowed between successive samples, which means that parameters should suffer less distortion. The method makes use of a variation of a technique used in shipbuilding for many years in the construction of ship outline drawings. A piece of flexible wood called a spline was used to fill in the curved line between the ribs. It has the property that the curve which it takes up between the ends is that of minimum potential energy and whose least-squares curvature is a minimum. Hence there is the possibility of replacing the succession of discontinuous lines from each sampling length by a 'spline function' having well-defined properties. This spline function could even be made physically for use on graphs if necessary [14].

If the spline function between $M$ ordinates is $\varphi_i$, and the profile is $z_i$ then the summation $S_u$ defined as

$$S_u = \sum_{i=1}^{M} (\varphi_i)^2 + \lambda \sum_{i=1}^{M} (\varphi_i - z_i)^2 \tag{2.22}$$

is a minimum, where $\lambda$ is a constant which has to be fixed. If $\lambda$ is very large compared with unity then the spline function is nearly a least-squares line. $\lambda$ has to be chosen to give the correct balance between smoothness and closeness of fit. $\lambda$ corresponds roughly to the elastic modulus of the piece of wood.

Much recent work has been done in using spline functions, particularly in numerical interpolation problems involving the generation of continuous surfaces from point-to-point numerical data points, a situation occurring in CAD. Essentially the use of splines is beneficial because of the good behaviour of the curve at the sample junctions. In common with filtering and envelope methods there are conditions at the beginning and end of the evaluation which need to be considered when choosing the spline value $\lambda$ [14].

### 2.1.2.4  Filtering methods

Perhaps the easiest way to separate the components of a signal on a frequency basis is by the use of filters. They have one big advantage over polynomial-type curve fitting: within generous limits little prior knowledge of the input waveform is required. The filter takes the waveform 'as received' and then operates on it. What happens in order to pay for this apparent freedom is that a certain amount of the input information has to be used up by the filter before it can give useful outputs; it needs time (or data) to get the feel of the input signal. Therefore the amount of usable data is reduced at the expense of versatility of operation. Filters have been used since the earliest days of surface metrology as a prerequisite to making meaningful measurements of any parameter.

The first sort of filter network used in this respect was a single $CR$ network. Later this was changed to a $2CR$ network buffered in the middle. The characteristics of this filter are such that if $z$ is an input profile sinusoidal signal of wavelength $\lambda_s$ and if $\lambda_c$ is the cut-off wavelength taken to be at 75% transmission (figure 2.27) then the transmission is given by

$$\text{transmission (\%)} = \frac{1}{1 + \frac{1}{3}(\lambda_s/\lambda_c)^2} \qquad (2.23)$$



**Figure 2.27** Transmission characteristics.

The 2*CR* characteristic was chosen for convenience and simplicity and to allow measurement with a skid. The value of 75% adopted for the cut-off itself is a compromise value. In the earliest instruments the British chose 80%, the Americans 70.7%. In the end, to make the instruments compatible, 75% was used. This is purely arbitrary but it turned out to be reasonably practical. The earliest filters used electrical components such as resistors and capacitors to produce the desired effect. Later instruments used digital methods. This type of transmission curve and the 75% cut-off was referred to as the 'standard filter' and was the preferred filter for many years as set down in international and national standards. Other filters and cut-off percentage transmission values are recommended today but this does not detract from the usefulness and practicality of the original standard filter which is still in use in many instruments.

Because the filter has been designed to simulate the graphical method, $\lambda_c$ is taken to be equal to the sampling length mentioned in section 2.1.1.2. The three different cut-off wavelengths recognized have been 2.5 mm, 0.8 mm and 0.25 mm, these roughly representing the boundaries of waviness found by experience for manufactured surfaces.

In what follows the 75% transmission characteristic will be used initially because this characteristic is in most existing instruments. The new recognized transmission at the cut-off is 50%. This will be introduced later. (See for example figure 2.33.)

Comparison of many $R_a$ values worked out graphically and those obtained by filtering shows that the decision to make the cut-off wavelength equivalent to the sampling length was justified. Differences of more than 5% are rare. Such correspondence is not of the same order if other parameters such as $R_t$ or $R_p$, for example, are being measured. The filter characteristics given by equation (2.23) are meant to be those of the instrument as a whole and not of any constituent part. The characteristics constitute therefore an instrument transmission curve, and so for a stylus instrument the characteristic should be taken as the way in which sinusoidal displacements at the stylus are revealed at the output. All frequency cut-offs, for example in amplifiers, recorders, etc, should be taken into account.

For most modern instruments these all tend to fall well outside the range of wavelengths under consideration and as a result negligible error is introduced by incorporating the simple two-stage filter having the desired characteristic. The way in which filters work on the surface profile waveform has to be understood, especially for the earlier type of filter. This is because not only are amplitudes modified but a certain distortion of the wave occurs due to phase shift. This is not a problem in the later digital filters but nevertheless it has to be demonstrated so that discrepancies between old and new instruments can be understood and taken into account.

How such a filter produces a mean line can be seen in figure 2.28.

In this figure the upper graph shows the input signal and the bottom shows the output. Between these is the profile signal but this time it has had the difference between the other two graphs plotted as a broken line.

**Figure 2.28** Relation of mean line to profile.

This line represents the low-frequency signal being blocked by the filter; this blocked signal is the reference line from which the surface parameters can be measured.

There is a difference between filter techniques and other methods, which is that in its simplest analogue form the filtering technique is time dependent. The behaviour of the filter is determined by the time dependence of the signals passing through it. The wavelength on the surface of these signals has to be related to time by the scan rate of the instrument. Thus the frequency resulting from a wavelength $\lambda_s$ is given by $f = s/\lambda_s$ where $s$ is the speed of traverse. This time dependence used to be a serious problem but it is not so today because the behaviour of the filter can be exactly reproduced by a mathematical technique which, like all the previous reference lines, can be purely spatial. Ideally all procedures in surface metrology should be spatial and not temporal.

The benefit of using filters, apart from their impartiality in dealing with all incoming signals, is their ease of use and cheapness. The standard $2CR$ filter does, however, suffer from some disadvantages over other methods. An example of this is shown in figure 2.29

When a profile signal has significant components with a wavelength near to that of the cut-off of the filter a certain amount of distortion can arise. The output signal no longer closely resembles the input waveform. This is due to phase shift. What happens is that some components of the profile get shifted relative to each other in their passage through the filter [15, 16]. At the same time some get attenuated according to



**Figure 2.29** Effect of filtering on periodic shapes.

where their wavelengths fall on the transmission curve of figure 2.27. Some results on typical waveforms are shown in figures 2.29 and 2.30.



**Figure 2.30** Effect of phase correction on standard shapes: (*a*) standard filter, sine wave profile; (*b*) standard filter, triangular profile; (*c*) and (*d*) phase-corrected standard filter; (*e*) and (*f*) phase-corrected linear profile.

Apart from the visual distortion of the profile which can sometimes result from this effect, parameters measured from the output waveform can also be in error. Fortunately this effect is not serious in the case of an $R_a$ measurement but it is, or can be, important in the case of $R_p$ or $R_t$, for example. It is fundamental when peak curvature or similar parameters are being measured for tribological experiments. This is why phase distortion of surface data occupies such an important position in surface metrology.

Another problem associated with the standard filter emerges in the problem of calibration. Because the standard filter has a rather gradual transition from the nominal 'pass' band to the nominal rejection region and because of the phase distortion the calibration of instruments can be difficult.

In order to elucidate the nature of subsequent improvements to the assessment of parameters, the manner in which the standard filter responds to some idealized repetitive waveforms will be considered. This is because the standard wave filter is generally satisfactory for random roughness waveforms. Repetitive waveforms tend to accentuate any shortcomings.

Look first at figure 2.30(*a*) which shows how the standard filter responds to sine waves. The figure shows sine waves of different wavelengths. The full line shows the original wave and the dotted line is the mean line found by the wave filter, this being the line that accounts for the output behaviour. The difference at any instant between the two lines is the filtered profile value. The top four graphs show sine waves within the pass band; notice that the mean line is displaced to the right relative to the original sine wave in all cases, which means that it is delayed in time relative to the profile, and further that the mean line has a large amplitude of undulation even for wavelengths much smaller than the cut-off, for example for one-third of the cut-off (the top graph). Another point is that the amplitudes of the mean line can be of the same order as the original sine wave, even just outside the pass band. Despite these points the filtered profile, which is the difference at any position between the mean line and the profile, obeys exactly the specified amplitude transmission characteristic mentioned previously. As an example, at the cut-off the mean line has an amplitude of just over 90% of the original sine wave profile, but the filtered profile has the correct amplitude of 75% of the amplitude of the profile. The amount by which the mean line is shifted relative to the profile depends on the wavelength of the sine wave profile.

For most commonly used filters, specifying the amplitude transmission characteristics automatically fixes the phase characteristics. This class is known as minimum phase [16], of which the standard filter of 2*CR* networks is a typical member [15]. Specifying the transmission characteristics for the standard wave filter automatically implies a phase-shifted mean line.

Suppose that a filter is available which has the same transmission characteristic as the standard filter, but at the same time has a mean line which is not shifted in phase relative to the profile. How is the mean line amplitude affected?

This situation is shown in figure 2.30(*b*) for the same sine waves as in figure 2.30(*a*), and for triangular waves in (*d*) corresponding to (*b*). It can be seen that the mean line is nearly straight at one-third of the cut-off. Compare this with the same profile for the standard wave filter. There is a dramatic reduction in the amplitude of the mean line, and this is true for all the cases shown. At the cut-off, the filtered profile for this new filter has the required amplitude of 75% of the profile—the 25% attenuation is accounted for entirely by the amplitude of the mean line. In other words, if the mean line is kept in phase with the sine wave profile, then, for the value of wavelength, the maximum amplitudes of the filtered profile and the mean line add up to unity, a criterion which does not hold in the case of the phase-shifted mean line. A direct result of this is that the mean line in phase with the profile undulates less. Summarizing, it may be said that the mean line becomes straight much nearer to the cut-off in the filter whose mean line is in phase than it does in the standard wave filter, although the filters have precisely the same amplitude transmission characteristics. This fact is of fundamental importance. It means that the shape of the roughness can be preserved even in the presence of considerable waviness and form error. Hence the roughness parameters retain their credibility throughout the filtering process.

For a sine wave profile the phase distortion simply takes the form of a phase shift of the filtered profile relative to the original profile, but for any other profile that can be considered to be made up of a number of such sine wave components, the distortion is more complicated.

Consider now triangular waveform profiles of differing wavelengths. Remembering that the filtered profile is the difference at any point between the mean line and the profile waveform, it can be seen that the filtered profile for the zero phase-shifted mean line bears a much closer resemblance to the original waveform than it does for the standard wave filter [16].

The zero-phase filter has a more realistic mean line because the sine wave components making up the triangular waveform are not shifted relative to each other in the filter. Consequently, the components have the

same relative positions upon emerging from it. Hence, even taking account of those components that have been attenuated in the filter, the output still resembles the input. This is not so true in the case of the standard wave filter. Distortion of the filtered profile can make it difficult to assess numerically. This is the problem that can be encountered in practice when highly repetitive profiles just within the pass band are put into the standard wave filter. As an example, the triangular waveforms shown in figure 2.30(*c*) are a close enough approximation to a practical waveform to illustrate how the problem arises. Consider the filtered profile in figure 2.30(*a*); the concept of a peak output at the cut-off, say, is difficult to imagine—the peak shape has virtually disappeared. Now look at figures 2.30(*b*) and 2.30(*c*) where the peak is noticeable and unambiguous to measure.

So far only the phase characteristics have been considered. There remains the problem of deciding on a suitable transmission characteristic. This is not easy, for it has to be remembered that a good reference (or mean) line has to be realistic for waviness as well as for roughness. Another point is that transmission characteristics for surface filters are plotted as transmission percentage as a function of spatial wavelength, not frequency as is conventional. This adds some complications. Take as a start an ideal roughness filter, that is one where waviness is not considered except to be excluded. The original ideal for this type of filter, known as the Whitehouse filter, is outlined below. The mean line is straight for all wavelengths within the pass band of the filter, that is up to the cut-off. This would mean that the profile would suffer no attenuation for wavelengths up to the cut-off. From the point of view of surface roughness measurement this seems sensible, because it is natural to suppose that if a cut-off has been chosen for the filter of a value larger than the longest roughness wavelength on the surface, then all the roughness will be passed unattenuated. Another point about the transmission characteristics which can be mentioned concerns the behaviour outside the cut-off. Although the behaviour outside the cut-off is not as important as that within the cut-off, it still has some relevance to the measurement of the roughness. The standard wave filter tends to fall off too gradually for wavelengths longer than the cut-off, with the result that waviness components can be included in the roughness assessment.

From these factors it appears that an amplitude transmission characteristic which is unity up to the cut-off wavelength and which falls rapidly after the cut-off would be more suitable for roughness measurement. Excessively high rates of attenuation, however, could be unrealistic mechanically because a considerable variation is not expected in the functional behaviour of, say, two surfaces having a roughness of equal amplitude but slightly different wavelength. One such characteristic that has seemed practical is one having unity transmission up to the cut-off which then falls off to zero at three times the cut-off, the rate of attenuation being linear with equivalent frequency. This has the merit of being compatible with the long-established sequence of sampling length cut-off values.

Figure 2.30(*e*) and (*f*) shows how a filter having an in-phase mean line similar to the one mentioned previously, but having the new transmission characteristics, behaves with the sine and triangular waveforms. The figures show a straight mean line right up to the cut-off and no distortion of the filtered profile. This is what could justifiably be called a filter with a well-behaved mean line. This filter will be referred to as the Whitehouse phase-corrected filter [16].

So far only idealized waveforms have been shown. In fact these accentuate the difficulties encountered in practice. Mean lines, as a rule, do not oscillate with such a large amplitude for practical waveforms within the cut-off even for periodic profiles, because a random component is always present. In the case of profiles which are random, the undulation of the mean line and the distortion of the filtered profile from the standard wave filter are not so obvious. The majority of profiles of this type have realistic mean lines, providing that the longest spacings are short compared with the cut-off length. However, for the standard wave filter the distortion of the filtered repetitive profile and the undulation of the mean line have presented a serious enough problem in some instances to warrant correction.

Some practical profiles are shown in figure 2.31, together with a comparison of the mean line found by the standard wave filter and the phase-corrected wave filter. They show that the phase-corrected filter has advantages over the standard wave filter in a number of ways. The gain is usually greatest when other factors

apart from roughness measurement have to be considered, such as only a small amount of available traverse or waviness components lying near to the dominant tool marks on the surface. In fact, the advantage is usually greatest in cases where it is impossible, for one reason or another, to choose a cut-off for the standard filter long enough to make phase distortion negligible. Many of the parameters used to measure roughness, such as the peak value, the derivatives or the bearing ratio curve, are affected by phase distortion. Centre-line average $R_a$ is affected to a lesser extent, while RMS $R_q$ and some related parameters are not affected at all.

Turned

Ground

Blast



**Figure 2.31** Comparison between standard and phase-corrected filters.

The roughness which is transmitted by the phase-corrected filter actually looks like the roughness on the original profile and it is a sensible precaution to measure any desired parameter, even for waveforms containing components near to the cut-off. Another point is that the mean line for components within the cut-off is straight, with the result that all roughness within the cut-off is assessed. Yet another possibility offered by this filter is that the mean line of the roughness found by its use could properly be regarded as the waviness. The characterization of the surface geometry taken as a whole therefore becomes a realistic possibility. Having the capability just to measure the roughness alone does not allow this vital possibility. This is because the shape of the mean line is only affected by components outside the cut-off (which are usually due to waviness) and also because the mean line will be in phase with these components. The mean lines for the repetitive waveforms outside the cut-off can be seen in figure 2.30(c).

Figure 2.31 shows how the mean lines for practical profiles look convincing for the waviness.

From what has been said it is possible to set down some characteristics that would be suitable for a filter having these properties. It would have the transmission characteristic $F(\lambda)$ given by

$$F(\lambda) = \begin{cases} 1 & \lambda_s < \lambda_c \\ P(\lambda) & \lambda_c < \lambda_s < K\lambda_c \\ 0 & K\lambda_c < \lambda_s. \end{cases} \tag{2.24}$$

In the case shown above $F(\lambda)$ is usually expressed in terms of frequency for reasons which will become obvious.

Thus, if frequency $f_c$ corresponds to $1/\lambda_c$, $F(f)$ corresponds to $F(\lambda)$ and $K_a$ corresponds to $1/K$:

$$F(f) = \frac{K_a f - f_c}{(K_a - 1)f_c}. \tag{2.25a}$$

In the filter example chosen as suitable, $K_a = 3$. Over the range $f_c/K_a < f < f_c$ the characteristic is linear with change in frequency, or $F(\lambda) = (\lambda K - \lambda_c)/\lambda(1 - K)$ over the range $\lambda_a < \lambda < K\lambda_c$, i.e. $P(\lambda)$.

The phase shift is ideally zero as is seen from figure 2.31. This implies that the imaginary part of the filter characteristic should be zero, that is phase $\Phi(f) = X(f)/R(f)$ for a general filter; this is only zero if $X(f)$ is zero. What this means and how it is achieved will be seen in the next few pages.

The Whitehouse phase-corrected filter has characteristics which are virtually ideal for roughness but not for waviness. In many instances this does not matter. However, in order to make the roughness and waviness filters complementary (i.e. always add up to unity) the cut-off has been repositioned to be 50% rather than 75% and the transmission characteristic gradual rather than sharp as advocated by Whitehouse. This means that, at the cut-off, both waviness and roughness will be transmitted. The new preferred filter has nominally Gaussian characteristics rather than the standard phase-corrected filter or Whitehouse phase-corrected filter. In order to understand the implications of this some idea of how filters work will be considered. The concepts here will be referred to many times in the text because the decomposition of the surface geometry into components is functionally important and filtering is a preferred way of doing it.

*(a)  Filter theory*

The frequency characteristic (or Fourier transform) $H(\omega)$ of a filter can be expressed in terms of the angular frequency '$\omega$' as

$$H(\omega) = k(\omega)\exp(j\phi(\omega)) \tag{2.25b}$$

where $k(\omega)$ is the amplitude transmission characteristic and $\phi(\omega)$ is the phase. The problem in surface metrology is to get a useful form for $k(\omega)$ whilst maintaining a form for $\phi(\omega)$ which eliminates phase distortion. This is true for roundness, roughness and all other metrology issues.

The impulse and step responses of a filter are a means of explaining how the filter behaves when it is subjected to a signal. The output of a filter in general terms is given by the superposition integral which, in effect, says that the output at any time is given by the convolution of the input signal with the impulse response of the filter.

Thus

$$g(t) = \int_{-\infty}^{t} f(\tau)h(t - \tau)\mathrm{d}\tau \tag{2.25c}$$

or

$$g(t) = f(\tau) * h(t)$$

where $f(t)$ is the input signal, $g(t)$ is the output and $h(t)$ is the impulse response.

In equation 2.25(c), $h(t - \tau)$, the reversed impulse response, is sometimes called the weighting function of the filter.

If $H(\omega)$, $G(\omega)$ and $F(\omega)$ are the Fourier transforms of $h(t)$, $g(t)$ and $f(t)$ respectively

$$H(\omega) = \int_{-\infty}^{\infty} h(t)\exp(-j\omega t)\mathrm{d}t \tag{2.25d}$$

and

$$G(\omega) = F(\omega)H(\omega). \tag{2.26}$$

The idea of the weighting function is very important in the application of filters to surface metrology problems because it enables the behaviour of a filter in time to be visualized as if the filter was acting on the profile on the chart. It will be shown next that the time variable $t$ can equally be replaced by a spatial variable

(as indeed it can be replaced by a numerical variable in a computer) without loss of generality. The effect of different filters can therefore be visualized for the metrologist. It also relates more nearly to the concept of functional filtering in tribology where the weighting function can be a contact pressure distribution.

Equation (2.26) is very important because it shows that the frequency characteristics of a filter (i.e. its response to sinusoidal signals) are tied down to the shape of its impulse response. Altering the impulse response alters the frequency characteristics.

The impulse response can alternatively be expressed in terms of a fraction of the time constant of the filter which, in turn, can be expressed in terms of a fraction of the sampling length when referred to the profile. This has the advantage of making not only the weighting function but also the variable in the superposition integral non-dimensional as well as relating more directly to the surface.

*(b) Different impulse responses and weighting functions*

Consider the low-pass filter in figure 2.32(*a*). It attenuates high frequencies. When referred to the time domain it means that the impulse response actually spreads in time. Because of this the impulse response tends to average out the high frequencies in the input waveform during the convolution process.

For the high-pass filters in figure 2.32(*b*) the situation is different because an impulse is present in the impulse response, which is opposed by a low-pass component. If the transfer function of an ordinary high-pass filter is $\overline{H}(p)$ it can be written in the form

$$\overline{H}(p) = 1 - H(p) \tag{2.27}$$

where $H(p)$ is for a low-pass filter and $p$ is the Laplace operator.

Equation (2.27) inverse-transforms into an impulse response $h(t)$, where

$$h(t) = \delta - h(t) \tag{2.28}$$

$h(t)$ is the impulse response of the low-pass component and $\delta$ is an impulse at the origin of unit weight.



**Figure 2.32** Impulse response of linear phase and standard filter. Standard filter in (*c*) has phase distortion; linear phase equivalent in (*d*) has no phase distortion.

A signal $f(t)$ put into a high-pass filter gives an output

$$g(t) = \int_{-\infty}^{t} h(t - \tau) f(\tau) d\tau$$

$$= \int_{-\infty}^{t} \delta(t - \tau) f(\tau) d\tau - \int_{-\infty}^{t} h(t - \tau) f(\tau) d\tau \tag{2.29}$$

but $\int_{-\infty}^{t} \delta(t - \tau) f(\tau) d\tau = f(t)$ because of the sampling property of impulses. Hence equation (2.29) becomes

$$g(t) = f(t) - \int_{-\infty}^{t} h(t - \tau) f(\tau) d\tau = f(t) - m(t). \tag{2.30}$$

In practice the lower limit of the integral can be taken to be zero.

Electrically $m(t)$ is the signal blocked by the filter. In surface metrology $m(t)$, when referred to the profile graph, is the mean line. The removal of $m(t)$ from the input profile constitutes the high-pass filtering action; $h(t - \tau)$ is the weighting function of the mean line.

Equation (2.30) can be expressed in a form more suitable for surface metrology by changing the time axis to a non-dimensional fraction of the sampling length $\alpha$:

$$g(\alpha) = \int_{0}^{\alpha} \delta'(\alpha - \alpha') f(\alpha') d\alpha' - \int_{0}^{\alpha} h(\alpha - \alpha') f(\alpha') d\alpha' = f(\alpha) - m(\alpha). \tag{2.31}$$

In equation (2.31) $\alpha'$ is a variable similar to $\tau$ in equation (2.30), that is $h(\alpha) = \delta' - h(\alpha)$ where $\delta'$ and $h(\alpha)$ have unit weight when integrated with respect to $\alpha$.

For the standard wave filter

$$\bar{h}(t) = \delta - \frac{1}{RC}\left(2 - \frac{t}{RC}\right) \exp\left(-\frac{t}{RC}\right) \tag{2.32}$$

or

$$\bar{h}(\alpha) = \delta' - A(2 - A\alpha) \exp(-A\alpha) \tag{2.33}$$

where $A = \lambda/s\,RC$ ($\lambda$ is the sampling length and $s$ is the tracking speed) and $\alpha = x/\lambda$, where $x$ is the distance along the surface. In equation (2.32) both parts have the dimensions of reciprocal time whereas they are dimensionless in (2.33), which means that the ordinate scale does not change with the cut-off. The factor $1/T_c$ is taken into the variable $d\alpha'$ of equation (2.31) where $T_c$ is the equivalent time of the sampling length.

*(c) Linear phase filters ([figure 2.32(b)](#))*

The phase characteristics of a filter effectively show how sinusoidal signals get shifted in time in their passage through the filter. The real criterion for the filtered profile to be undistorted is that the constituent sinusoidal components making up the profile are not shifted relative to each other during the passage of the profile through the filter. One method of doing this would be to ensure that none of the components got shifted at all. This would mean that the phase shift is zero for all frequencies.

Now

$$H(\omega) = k(\omega)\,\exp(j\phi(\omega)) = R(\omega) + jX(\omega) \tag{2.34}$$

where $\phi$, the phase angle, is given by

$$\tan^{-1} X(\omega)\big/ R(\omega) \qquad (2.35)$$

So $X(\omega)$, the imaginary component, would have to be zero in equation (2.35) to make the phase zero, leaving

$$H(\omega) = k(\omega) = R(\omega) \qquad (2.36)$$

But, for the transmission characteristic of the filter to be real, the impulse response must be an even function (i.e. symmetrical about the time origin axis), which is impossible in practice because the impulse response or the weighting function of the filter cannot extend into the future as well as into the past. Hence, the only possibility of getting no distortion of the filtered profile is to arrange that the wavelength components of the profile are all shifted by the *same amount in time* by the filter. Suppose that this meant delaying all components by $t_0$. One component at angular frequency $\omega_1$, say, would have to be shifted in phase by $-\omega_1 t_0$ rad to get this delay. A component of angular frequency $\omega_2$ would have to be shifted by $-\omega_2 t_0$, and similarly for any component. In fact, to satisfy this delay of $t_0$ in general $\phi(\omega) = -\omega t_0$; the phase has to have a linear relationship with frequency. Therefore, the transmission characteristic for a filter having no phase distortion but a delay is of the form

$$H(\omega) = k(\omega)\exp(-\mathrm{j}\omega t_0) \qquad (2.37)$$

Such a filter is called a linear phase filter. It took a long time in surface metrology to realize the potential of such filters!

How do the impulse responses of the zero-delay and linear phase filters compare? It is easy to show that if $h_0(t)$ is the impulse response of the zero-delay filter, then the impulse response of the linear phase filter having the same amplitude transmission is $h_0(t - t_0)$. This means that they have the same shape, but the linear phase impulse response is shifted by $t_0$ along the positive time axis — they both have a symmetrical weighting function but the zero-delay impulse response has its axis of symmetry on the time origin, whereas the linear phase impulse response has the axis of symmetry at $t = t_0$. Shifting the axis of symmetry to $t_0$ makes the impulse response practically realizable. Summarizing, it may be said that it is possible practically to make a filter giving no phase distortion only by allowing a uniform delay of all components passing through it. This is achieved by a filter having linear phase characteristics, which implies that it has an impulse response which is symmetrical about an axis of $t = t_0$ on the realizable side of the time axis. In the non-dimensional form $t_0$ becomes $\alpha$ (figure 2.30(d)).

If

$$h_0(t) \Leftrightarrow H_0(\omega) \qquad (2.38)$$

then

$$h_0(t - t_0) \Leftrightarrow H_0(\omega) \exp(-\mathrm{j}\omega t_0). \qquad (2.39)$$

For the zero-delay case the high-pass impulse response $h_0(|t|)$ can be expressed as $\delta - h_0(|t|)$, where $h_0(|t|)$ is the low-pass impulse response and $\delta$ is an impulse at the origin. The corresponding linear phase high-pass impulse response $h_L(t)$ is therefore given by

$$h_L(t) = \delta(t - t_0) - h_0(|t - t_0|). \qquad (2.40)$$

Equation (2.30) becomes for a linear phase filter

$$g(t) = f(t - t_0) - \int_0^t h_0(t - t_0 - \tau) f(\tau) d\tau. \tag{2.41}$$

Equation (2.40) shows that the impulse component lies at the axis of symmetry of the low-pass impulse component, which means that in terms of equation (2.40) the signal itself has to be delayed by $t_0$ before taking it from the low-pass component. This has the effect of producing, at time $t$, the filtered output—without distortion corresponding to the profile at $t - t_0$.

Two points are worth noting concerning symmetrical impulse responses: one is that the step response has an axis of odd symmetry about $t = t_0$, and the other is that the operations of correlation and convolution become the same except for an averaging term.

*(d) Different linear phase filters*
The conditions for a symmetrical weighting function to be suitable are the following:

1. It has an axis of symmetry later than $t = 0$ to an extent such that no considerable part of the function crosses the $t = 0$ axis.
2. It is concentrated in a central lobe and dies away quickly on either side of the axis of symmetry.
3. Any negative portions should be small.
4. It must have an amplitude transmission characteristic suitable for use in surface metrology.

A number of different linear phase filters have been investigated, including those which have Gaussian and raised cosine impulse responses. Perhaps one of the most obvious is the linear phase filter having the same amplitude transmission as the standard filter. It has an impulse response

$$\bar{h}(t) = \delta(t - t_0) - \frac{\omega_c}{\pi\sqrt{3}} \exp\left(-\frac{\omega_c}{\pi\sqrt{3}}|t - t_0|\right) \tag{2.42a}$$

or alternatively

$$\bar{h}(\alpha) = \delta'(\alpha - \bar{\alpha}) - \frac{\pi}{\sqrt{3}} \exp\left(-\frac{2\pi}{\sqrt{3}}|\alpha - \bar{\alpha}|\right) \tag{2.42b}$$

where the latter expression puts the impulse response in non-dimensional form. $\omega_c$ is the angular frequency equivalent of the sampling length, and $\bar{\alpha}$ is the position of the axis of symmetry—equivalent to $t_0$.

Another alternative is the ideal linear phase high-pass filter where

$$\bar{h}(t) = \delta(t - t_0) - \frac{\sin \omega_c(t - t_0)}{\pi(t - t_0)} \tag{2.43a}$$

or

$$\bar{h}(\alpha) = \delta'(\alpha - \bar{\alpha}) - \frac{\sin 2\pi(\alpha - \bar{\alpha})}{\pi(\alpha - \bar{\alpha})}. \tag{2.43b}$$

These, however, did not seem to be as suitable a starting point as the Whitehouse filter, which has unity transmission for wavelengths up to the cut-off and zero transmission at three times the cut-off. The attenuation rate is proportional to the equivalent frequency.

The expression 'phase-corrected filter' has a slightly different connotation in communication theory, but it is useful in this context. The general equation for impulse responses having this linear amplitude form is

$$\bar{h}(t) = \delta(t - t_0) - \frac{2}{\pi\omega_c(1 - B)} \times \frac{\sin\left[\omega_c(1 + B)(t - t_0)/2\right]\sin\left[\omega_c(1 - B)(t - t_0)/2\right]}{(t - t_0)^2} \qquad (2.44a)$$

or

$$\bar{h}(\alpha) = \delta'(\alpha - \bar{\alpha}) - \frac{1}{\pi^2(1 - B)} \times \frac{\sin\left[\pi(1 + B)(\alpha - \bar{\alpha})\right]\sin\left[\pi(1 - B)(\alpha - \bar{\alpha})\right]}{(\alpha - \bar{\alpha})^2} \qquad (2.44b)$$

where $B$ is the ratio of wavelengths having unity to zero transmission, being equal to $1/3$ in this case.

Other possible filters have been examined such as those in which the attenuation characteristic is linear in wavelength rather than frequency as it is in equation (2.25a).

In these filters the formulae for the weighting functions and the transmission characteristics are very different. It seems plausible to consider a filter system which has virtually the same form of equation in both time and frequency. This is in fact easily done by making the shape of the weighting function Gaussian. The frequency characteristics of this in terms of amplitude are also Gaussian because the Gaussian shape is its own Fourier transform. Unfortunately surface metrologists work in wavelengths, so the transmission characteristic (in wavelength terms) of a Gaussian weighting function is not itself Gaussian. It is, however, recognizable as being derived from a Gaussian weighting function and has been selected as the new preferred filter. Also, in order to match the waviness and roughness roles the cut-off is at 50%.

Thus the weighting function is

$$h(\alpha) = \frac{1}{\alpha_1\lambda_c}\, \exp\left[-\pi\left(\frac{\alpha}{\alpha_1}\right)^2\right] \qquad (2.44c)$$

(see figure 2.33). Putting the transmission characteristics in terms of $\alpha$ gives $H(1/\alpha)$:

$$H(1/\alpha) = \exp[-\pi(\alpha_1/\alpha)^2]. \qquad (2.44d)$$

The transmission of the roughness profile is complementary to that of the transmission of the mean line because the roughness profile is the difference between the actual profile and the mean line. Thus

$$H_R(1/\alpha) = \exp[-\pi(\alpha_1/\alpha)^2]. \qquad (2.44e)$$



**Figure 2.33** Gaussian filter with attenuation curves.

This gives the amplitude of the sinusoidal signal of wavelength $\alpha$ where

$$\alpha_1 = \sqrt{\frac{\ln 2}{\pi}} = 0.4697.$$

Having examined the behaviour of filters it is possible to compute them using methods other than the obvious convolution method described above. Many variants on the standard filter and phase-corrected filter are possible. One of the major variables is the point on the characteristic at which the cut-off attenuation is assumed to occur. It could be unity, 75% or 50% but current thinking is based on 50%. Another factor is concerned with the problem of realizability. Can the weighting function be adequately represented in a computer? These factors will be considered in chapter 3. Robust filters will be in chapter 5.

Filtering using a phase-corrected method has its problems if characteristics such as those outlined in equation (2.25a) are to be used. This is because of the sharp attenuation drop-off. As this is faster than the Gaussian characteristic the impulse response tends to be slowly decaying, with the result that some degree of truncation of the weighting function has to be allowed for in the computer. This produces problems in standardization for instrumentation. The alternative is to limit the weighting function arbitrarily to a known extent, for example to use a box function (as in the mid-point locus method) or a triangle. The other problem is that such amplitude characteristics do not lend themselves well to recursive filtering techniques.

So far the methods considered have been confined to linear techniques in which the method of separating the roughness from the rest of the signal could be judged using Fourier analysis. This idea is fine for conventional wave filters and any technique that can be reduced to the basic convolution. However, there are alternative methods which depend on discrimination in domains other than spacing, such as in height, slope and curvature. Also multiprocesses such as plateau honing require special treatment.

### 2.1.2.5    Envelope methods

In the methods considered above all the whole-profile trace has been used in the positioning of the reference line. There are other methods, however, in which only selected parts of the profile like the extreme peaks or valleys contribute to the positioning of the reference, such reference lines highlighting the importance of the peaks and valleys. A simple straight line grazing the deepest valleys was originally proposed by Schmaltz whilst Nicolau decided on using the highest peaks [17], the roughness value being taken as the average or maximum profile distance from the line. Neither of these methods has been extensively adopted. Two different envelope methods have since been proposed, the first, the rolling circle envelope, by Professor von Weingraber and the other by M Scheffer and Professor Bielle. These will be considered in turn.

*(a) Rolling circle envelope (E system)*
In 1957, von Weingraber [18, 19] proposed a system in which a circle was imagined to be rolling over the crests of the profile graph. He advocated such a system because he believed that it mimicked very largely the vertical position on the surface that the anvil of a micrometer or similar gauge would take if the size of the component were being measured. The E system was a metrological attempt to link surface and dimensional metrologies as seen diagrammatically in figure 2.34.



**Figure 2.34**

During its path the circle would obviously not follow a straight course; sometimes it would drop into valleys, sometimes it straddled the peaks. It all depended on the radius of the circle relative to the average separation of the peaks. The reference line was taken to be the locus of the lowest point on the circle (figure 2.34).

The advantage of this system was that, with a suitable radius, the waviness and roughness could be separated out to some extent, the radius acting as an equivalent of the sampling length in the mean line methods (M system). It represented an ambitious attempt to provide a basis for a large number of geometrical features. Where it failed was its difficulty in instrumentation. This method has an advantage, at least in this form, of being quite simple to construct graphically and it does not suffer from the choice of origin, as does the least-squares method described earlier. Suggested radii for the assessment of roughness have ranged from 3.2mm to 50mm, the earliest being 25mm. It seems obvious that, in order to get the same degree of versatility as the choice of sampling length allows, it is necessary to provide the same sort of choice for the radius. In a relatively recent paper Radhakrishnan [20], examining the effect of the radius size on the measured roughness, showed that the rate of change tended to be greatest for 3.2 mm and concluded that this value would be best for distinguishing between roughness and waviness. This philosophy is questionable because it implies that some of the roughness is being attenuated. Perhaps a better criterion is to choose a radius corresponding to the minimum rate of change, which was the criterion adopted in deciding on the meter cut-offs in the M system. Using this criterion has led to the conclusion that the original 25 mm radius would be a better choice for general use, compared with the 0.8 mm cut-off in the M system.

One of the big problems with this method, apart from the instrumental ones, is that the technique, although being an effective filter, has transmission characteristics which depend on amplitude. For different amplitudes of waveform the transmission characteristic changes.

That the choice of a set of radii is not simple may be judged from the point that for a random surface the $R_p$ value has an approximate relationship to $R$ given by

$$R_p \propto \sqrt{\ln R} \qquad (2.45)$$

which will be shown in the chapter on instrumentation. In equation 2.45 $R$ is the ball radius.

It has been suggested that any system which is primarily concerned with a reference based on peaks rather than the profile as a whole must be functionally more significant in rubbing or contact processes.

Recently, Radhakrishnan [21] has been investigating the three-dimensional characteristics of the envelope system using a computer. What he has shown is that there are considerable divergences between results obtained from a circle on the profile graph and the three-dimensional results nominally obtained with a ball on a real surface. The E system was originally dropped because of instrumental and computational difficulties. However, it is now proving to be more useful in both function and process control. This resurgence has been a direct result of improvements in measurement techniques and data processing. However, the original concept — that of using the ball directly on the surface—has not been pursued. Now it is simply simulated in a computer once the original data has been obtained—by means of the conventional sharp stylus or optical probe!

*(b) Peak and valley envelopes—R & W*

A system of measurement being advocated in some countries, *R & W*, involves a peak and valley envelope line. In its original form the method was based on a subjective assessment of the profile. The graph is taken with a horizontal magnification of about $\times$ 20 and then a skilled operator traces two undulating lines more or less following the peaks and valleys but ignoring minor features separated by less than about 0.5 mm. The roughness value $R$ is then taken to be the average of 10 measurements of height between the two envelope lines.

Waviness can also be measured in this system by a system of tangents drawn between the uppermost points on the peak envelope.

Certain restrictions have to be placed on the drawing of these tangents. If the depression of the envelope between any two contacts A and B shown in figure 2.35(*a*) is only 10% of the roughness value, then the tangent

has to be redrawn ignoring the secondary tangent contact B and proceeding to the next contact C. How this technique is instrumented is described later in the chapter on instrumentation.

### *(c) Motif methods [6]*

The motif method is an envelope technique. It was suggested in France in the early 1970s; see Standard P.R.E.05-015 (1981).

The 'motif' method in principle is rather simple. The unfiltered profile is divided into useful geometrical characteristics in accordance with an empirically found algorithm which should be determined by the combined experience of people dealing with profile analysis. In one form the motif itself is two profile peaks with a valley in between them.

The motif is characterized by the two heights $H_1$ and $H_2$. The characteristic depth $T$ of the motif has a special significance in the evaluation: it is the smaller of $H_1$ and $H_2$. The mean depth is $(H_1 + H_2)/2$ and $AR_1$ is the distance between the peaks. Application of the motif method involves checking whether two adjacent motifs can be treated separately or whether one of the motifs can be regarded as insignificant and absorbed into the other one to make one large motif. Two motifs cannot be merged if any of the four conditions apply.

1. Envelope condition (figure 2.35(*b*) (ii)). Two adjacent motifs cannot be combined if their common middle peak is higher than the two outer peaks.
2. Width condition (figure 2.35(*b*) (iii)). Adjacent motifs can only be combined if the result is less than or equal to 500 μm. This is a very weak condition indicating that motifs up to 500 μm are to be included as roughness. It corresponds in effect to a cut-off. The width limit for evaluating waviness is 2500 μm. These arbitrary rules have been proposed with the car industry in mind.
3. The magnification condition (figure 2.35(*b*) (iv)). Adjacent motifs may not be joined if the characteristic depth of each of the adjacent depths $T_3$ of the result is less than the largest characteristic depth of each of the adjacent motifs. This condition is most important. It means that in each case an unequivocal evaluation can be reached independent of the tracing direction.
4. The relationship condition (figure 2.35(*b*) (v)). Adjacent motifs may be considered if at least one has a characteristic depth less than 60% of the local reference depth $T_R$. Thus, adjacent motifs of approximately the same size cannot be combined to form one single overall motif.

The reference depth $T_R$ may appear in two forms. It may be the largest depth found in any one suitable profile section approximately 500 μm wide; it is then called the local depth reference. Alternatively it may be the characteristic depth $T_3$ which a combined motif would have.

The basic principle to follow when deciding on an evaluation is to find the largest motif in the measured profile which fulfils all conditions, that is the profile should be described with the least number of the largest possible motifs.

### *(d) Motif procedure (according to Fahl [22], figure 2.35(b))*

First, the local reference depth over the relevant sections in the profile is found. Then within each individual section every single motif is checked against its neighbour using the four conditions above and whenever necessary combined with a new motif.

Every new motif must be compared with its neighbour until no further combination of two motifs is possible within that section. This is carried out for each section of the profile. When completed, all motifs that lie outside the limits of the section must be checked with each other until no further mergers are possible.

In theory this is the end result. In practice, however, there are other considerations. Examples of this include the presence of single isolated peaks or valleys. These have to be smoothed to a 'reasonable' level. After this the 'corrected' profile results. Then, $R_m$ the average motif depth, $A_R$ the average width and $p(R_m)$, $p(A_R)$ the distributions of $R_m$ and $A_R$ are found. After these calculations a near profile comprising the peaks just found is made. The whole procedure is repeated, except that, instead of 500 μm, 2500 μm is used. This

corresponds to the envelope above the corrected profile. This time $W_m$, $A_w$ and $p(W_m)$, $p(A_w)$ are found and taken as the waviness parameters.



**Figure 2.35** (*a*) Motif method — R and W. (*b*) The motif and the four conditions for motif combination: (i) definition of motif; (ii) envelope condition; (iii) width condition; (iv) magnification condition; (v) relationship condition.

The whole procedure can be repeated to give the form profile.

This represents the complete breakdown of the profile into roughness, waviness and form.

*(e) Intrinsic filtering*

The representation of the original profile in blocks is reminiscent of the methods of filtering using the sampling length described earlier. In this case, however, the blocks are determined from the surface itself and could therefore better be described as 'intrinsic filtering'. Any local effects on the surface which are not characteristic will automatically be taken into account. The sampling 'matches' itself to the local properties of the surface.

In this motif method $R_m$ is approximately $R_y$ or $R_z$ qualified by a varied sampling length. $A_R$ is the average width between significant peaks. This corresponds to $S_m$, the average distance between zero crossings of the profile with the mean line.

The advantage of the motif method is that it 'reacts' to the surface uniquely according to a set of rules. The disadvantage is that the rules are arbitrary. They could, in principle, be modified to suit any given function. Another disadvantage is that the motifs are peak oriented; quite major peaks near to a valley are ignored, yet small peaks near to a major peak are not. A different set of rules would be needed to make the algorithm symmetrical. The arbitrary selection of 500 μm and 2500 μm is a disadvantage but no more than the 0.8 mm cut-off selected for the filters.

It seems that the motif method is another way of looking at profiles, perhaps complementing the filtering methods or perhaps best suited for dealing with multiprocess surfaces where no standard characteristics occur.

Another possible use is in plotting the distribution of given sizes of motifs. This is almost the metrology equivalent of a Pareto curve. Using this curve motif size against frequency of occurrence can be plotted for roughness and waviness and can represent some interesting information. Exactly how it can be used is not obvious but it is a different type of breakdown which does not rely on a reference line but on adequate definitions of peaks.

### 2.1.2.6 Summary

In the foregoing the development of characterization methods for a profile graph has been described. These have been somewhat piecemeal in the sense that height and spacing parameters have been described. At the same time considerations of reference lines and sampling lengths or equivalent cut-off lengths have been given. This is all part of the attempt to classify the surface into its significant components and as such is pertinent to the chapter on the characterization of surfaces. The classification has been involved with many factors including dimensional metrology, functional and manufacturing considerations as well as height and spacing parameters. A rather more coherent, yet not necessarily more relevant, classification is based upon random process analysis of the surface. It may appear that random process analysis occurs just about everywhere in classification in different guises. This observation is true. Such is the potential of random process analysis that inklings of it have occurred over many years and for different purposes. This chapter reflects this situation. It would not be correct to isolate the subject of random process analysis from its various functional usages, for they add to its credibility. Hence vestiges appear in many places in this and other chapters.

### 2.1.3 Statistical parameters of surface roughness

The theory of random process analysis is now well established in the field of metrology research but has yet to make an impact in industry. This difficulty in penetration has been partly due to the expense of measuring the various parameters and partly due to their interpretation. The former disadvantage is rapidly disappearing with the advent of cheap computers. The latter problem still remains.

Many investigators have tried to formulate statistical rules with which to describe the geometric properties of surfaces. Basically the problem is that surfaces can be random, deterministic or usually a complex

mixture of both. High-order probability density functions are needed to specify the general surface statistics completely at one end of the scale of complexity, yet only a simple formula involving a few parameters is needed at the other end to determine the complete profile. Somewhere in between these extremes lies the practical characterization of surfaces. Such a compromise may exist by considering the autocorrelation function of the surface and the amplitude probability density function respectively. Taken together they provide a reasonable basis for the topographic analysis of surfaces, especially if only a profile is being considered, but they are by no means the only functions that can be and have been used as will be seen later in this section.

### 2.1.3.1    *Amplitude probability density function (APDF ) or p(z)*

It is possible for any continuous waveform, whether it is a profile signal or not, to be drawn as a graph as shown in figure 2.36.



**Figure 2.36** Statistical breakdown of profile.

The two axes of the graph are height $z$ and $p(z)$, the probability density that any measurement has the value of $z$. Some of the properties of the APDF are that it has unit area, it extends over the range $z_{max}$ to $z_{min}$ and it may take on a variety of shapes depending on the profile. In metrology it is a useful convention to plot the $z$ axis to the APDF in the same direction as the $z$ direction of the profile. This arrangement makes visualization easier. It is perpendicular to the method of display usually used in statistical textbooks.

Also used to describe amplitude information is the cumulative distribution function (CDF) which is perhaps more useful. This represents at any level $z$ the probability that a measurement of the waveform will be less than $z$. To find the probability of a measurement of the surface lying between levels $z_1$ and $z_2$, it is necessary to integrate the APDF between these two values.

Thus, if $p(z)$ is the APDF, then the CDF is given by $P(z)$ where

$$P(z') = \int_{-\infty}^{z'} p(z) \mathrm{d}z \qquad (2.46)$$

and the probability of a measurement lying between $z_1$ and $z_2$ is

$$\int_{z_2}^{z_1} p(z) \mathrm{d}z.$$

Similarly $p(z) = \mathrm{limit}_{\Delta z \to 0}$ (probability of measurement between $z$ and $z + \Delta z$) (figure 2.36).

### 2.1.3.2    *Material ratio*

Related to the CDF is the Abbott–Firestone (material ratio) curve which has been mentioned earlier [8]. (This is $1 - P(z)$ or the probability that a measurement is higher than a given level $z$.) In mechanical terms it gives

the percentage of material to air of the surface profile at any level. For this reason it is sometimes referred to as the 'bearing ratio', as mentioned in section 2.1.1.1, or material ratio which is the current terminology.

These curves only give height information. That is, they only give the frequency of occurrence with which measurements lie between certain limits. They do not give any information about the spacing between such measurements or the order in which they occur.

Ways of disseminating the height information of a profile based upon these techniques have been attempted for a number of years. The earliest seems to have been by Pesante [9], who advocated the use of APDF curves to control the different manufacturing processes. Ehrenreich [10] and Reason [23] have proposed different features of the material ratio curve. Ehrenreich suggested measuring the slope between the air/metal ratios of $1/3$ and $1/2$. This value can be determined also from the APDF because it is a form for the differential of the bearing ratio curve. Reason used a model of the surface comprising what he called a consolidated profile made up of a sequence of mirror-imaged bearing curves. He also incorporated the scale of size of the graph in both directions.

In formal terms the APDF is more of a tool for specifying the characteristics of unknown waveform statistics. It is not necessary to specify it in the case of deterministic waveforms such as a square wave, sawtooth, etc, simply because the curve can be worked out from a knowledge of the formula for the waveform.

This can easily be verified with reference to a sine wave. Thus $z = A \sin x$ where the angle $x$ takes values from 0 to $2\pi$— in other words $x$ is uniformly distributed within the interval $0 \to 2\pi$. The probability function for $x$ treats it as a random function because it does not matter (from the point of view of height distribution) in which order the $x$ values are taken.

The probability density function is

$$p(x) = \frac{1}{2\pi} \qquad 0 < x \leq 2\pi \tag{2.47}$$

$$= 0 \qquad \text{otherwise.}$$

If $p(x)$ and $z = g(x)$ are known—a deterministic function for all $z$, $z'_0 = g(x_0)$—the probability that $z$ lies between $z_0$ and $z_0 + dz$ must equal the probability that $x$ lies in the range $x_0$ to $x_0 + dx$. Hence

$$p(x_0)\,dx = p(z_0)dz \quad \text{or} \quad p(z) = \frac{p(x)}{dz/dx}. \tag{2.48}$$

If each value of $z$ has $n$ values of $x$ then $x = g^{-1}(z)$ is multivalued and $p(z) = np(x)/(dz/dx)$. This is the determining relationship for the APDF in terms of the deterministic relationship between $z$ and $x$.

For a sine wave

$$\frac{dz}{dx} = A \cos x = \sqrt{A^2 - z^2}$$

so $x = \cos^{-1}(z/A)$ which is double valued and

$$p(z) = \frac{2p(x)}{A \cos x} = \frac{1}{\pi \sqrt{A^2 - z^2}} \quad \text{for} - A < z < A. \tag{2.49}$$

The corresponding distribution function is $P(z)$:

$$P(z) = \int_{-A}^{z} p(z)dz = \frac{1}{\pi}\left[\frac{\pi}{2} + \sin^{-1}\left(\frac{z}{A}\right)\right]. \tag{2.50}$$

The material ratio curve is unity minus the distribution function, so for a sine wave it is

$$MR(z) = \frac{1}{\pi}\left[\frac{\pi}{2} - \sin^{-1}\left(\frac{z}{A}\right)\right].$$

Because many surface waveforms are neither truly deterministic nor random, the material ratio curve does have some use in describing the surface; the shape can give some idea of the nature of the surface relative to the purely random wave given by the Gaussian or normal distribution $p(z)$ (figure 2.37)

$$p(z) = \frac{1}{\sqrt{2\pi}}\exp(-z^2/2\sigma^2) \tag{2.51}$$

where $\sigma$ is the RMS value of the profile $R_q$. The reason for this will be given later.



**Figure 2.37** (*a*) Gaussian amplitude distribution; (*b*) sinusoidal amplitude distribution.

Rather than specify the whole shape of the amplitude probability density function it is often more convenient to break it down using the moments of the curve. Thus if the *k*th moment of the curve is $m_k$

$$m_k = \int_{-\infty}^{\infty} z^k p(z)\mathrm{d}z \tag{2.52}$$

and if the mean value of *z* is

$$\bar{z} = \int_{-\infty}^{\infty} z p(z)\mathrm{d}z \tag{2.53}$$

then the central moments are given by $\mu^k$:

$$\mu^k = \int_{-\infty}^{\infty} (z - \bar{z})^k p(z)\mathrm{d}z. \tag{2.54}$$

It is these central moments that give useful parameters about the shape. Because the mean level $\bar{z}'$ is arbitrary, it represents the mean distance of the surface from the reference level in the instrument measuring the surface. As such it therefore contains no useful information. The information is contained in the relative values of the central moments and in particular the skew $S_k$ and kurtosis $S_{ku}$.

The skew is defined as

$$S_{sk} = \frac{1}{\sigma^3} \int_{-\infty}^{\infty} (z - \bar{z})^3 p(z) \mathrm{d}z \qquad (2.55)$$

and the kurtosis as

$$S_{ku} = \frac{1}{\sigma^3} \int_{-\infty}^{\infty} (z - \bar{z})^4 p(z) \mathrm{d}z - 3. \qquad (2.56)$$

where $\bar{z}$ is the mean height of $z$. Both are normalized with respect to the RMS value of the profile waveform $\sigma = \sqrt{\mu_2}$ (where $\sigma$ is the standard deviation of the surface and is equivalent to $R_q$).

Essentially the skew shows the degree of symmetry of the profile while the kurtosis shows the degree of pointedness or bluntness of the waveform (figure 2.38).



**Figure 2.38**   Examples of skew and kurtosis in surfaces.

How these can be used will be seen in the subsection on characterization. Al-Salihi [24] first proposed the moments in the form of a Gram–Charlier series. However, it is only recently that they have started to be used seriously. The skew in particular has been used as a control parameter to judge the value of the conventional $R_a$, $R_t$ values etc. This control property is based upon experience of the typical shapes of profiles, the manufacturing processes and the functional effects of such waveforms. Take, for example, the profiles shown in figure 2.38. These are taken from surfaces produced by different methods.

The formulae expressed above are not the only ones for the central moments. They can all be expressed in terms of baseline moments. Remember also they can be defined along the $x$ axis.

For example, letting

$$\bar{z} = \int_{-\infty}^{\infty} z p(z) \mathrm{d}z \qquad v_2 = \int_{-\infty}^{\infty} z^2 p(z) \mathrm{d}z \qquad v_3 = \int_{-\infty}^{\infty} z^3 p(z) \mathrm{d}z \qquad (2.57)$$

where $\bar{z}$, $v_2$ and $v_3$ are respectively the first, second and third moments about the arbitrary baseline from which the $z$ values are measured, the skew $S_k$ is

$$S_k = \frac{1}{\sigma^3} (v_3 - 3\bar{z}v_2 + 2\bar{z}^3) \qquad (2.58)$$

where $\sigma^2 = v^2 - \bar{z}^2$.

Judging from these it seems reasonable to propose that, if the skew value is bigger in value than $\pm 2$, then a straightforward measure of $R_a$ etc will be questionable; it will not be quantifying the feature of the profile

considered to be functionally most important. This value of 2 is rather fortuitous because it can be shown that the skew of a random wave truncated at the mean level $\bar{z}$ gives a skew of 2.1 which is very close to 2. For control purposes, therefore, a value of skew greater than 2 would indicate that care must be exercised when taking simple measurements.

### (a) Other properties of the APDF

If, as is usually the case, a workpiece is made using more than one process the profile often shows evidence of it. The profile may comprise more than one identifiable characteristic, each one given by one or other of the processes. This produces what is called a stratified surface. If the processes are additive in nature so that $z$ at any $x$ can be written $z = z_1 + z_2$, then the APDF $p(z)$ is given by the convolution integral

$$p(z) = \int_{-\infty}^{\infty} p_1(z_1)p_2(z - z_1)dz_1 \tag{2.59}$$

where $p_1(z)$ and $p_2(z)$ are the individual APDFs for $z_1$ and $z_2$. In many cases in manufacturing this additive effect is present: one set of geometries is superimposed on another and in other cases it is not. Consequently the convolution does not hold. In the former case deconvolution enables the variables $z_1$ and $z_2$ to be identified statistically by their APDF providing one of the two is known.

If many interdependent variables are responsible for producing the final profile the resultant APDF $p(z)$ will usually have a form that is Gaussian in shape. This is because of the central limit theorem which can be stated as follows: *The resultant effect of repeated convolutions will produce a Gaussian output irrespective of the shapes of the signals being convoluted*. In practice any more than four or five convolutions will give this effect. This explains why many finished processes, such as grinding, have an APDF which is Gaussian. Because of this fortuitous phenomenon much analytical work has been made possible on models of surfaces for use in contact theory and optics as will be seen in chapter 7.

### 2.1.3.3 Autocorrelation function (ACF) and power spectral density (PSD)

A major breakthrough in recent years in the characterization of surfaces has been the use of some of the mathematical tools used in communication theory. Those used in random process analysis have been most prominent. In particular, the autocorrelation function has been especially important. First used in 1946 in surface roughness [25], the idea of correlation itself is well known in statistics and stems from the need to predict the influence that one set of numbers has on another set. One set comes perhaps from the input to an experiment and the other from the output. Dealing with one profile signal alone is hardly what is meant by a random process. What is usually implied is that many sets of information are being considered and the average relationships determined. In the case of surface metrology the many sets of data correspond to the whole surface including texture. One profile represents merely a very small subset of the total data. However, such is the nature of correlation functions that by investigating the statistical properties of one or at most a few profiles, valid information is obtained about the whole surface. This is not necessarily the case when using the deterministic parameters. Many of them give information which is dominated by the idiosyncrasies of the individual profile and hence is unsatisfactory. Because of this property the use of correlation methods is on the increase; the correlation function is bringing more predictability to the measurement of the random features of surfaces. Why this is so will be seen presently.

The property of relating to the whole surface rather than individual profiles makes the autocorrelation function (and its equivalent, the power spectral density) very attractive for use as a basis for characterization.

In fact random process analysis, which comprises the amplitude distribution and the autocorrelation function and power spectrum, has been the source of much research into characterization. The idea is to estimate many parameters of significance in the function of the surface from a few random process parameters, thereby making savings on the amount of measurement required as well as increasing the reliability of measurement.

Consider now characterization using random process analysis. This has followed two distinct paths: one uses the power spectral density and its moments, and the other uses discrete values taken from the autocorrelation function.

Both methods were initially devised to help to predict the many functional requirements to which surfaces are subjected.

The former approach originated from the work of Rice [26] in communication theory. This was expanded to cover other areas by Longuett–Higgins [27], in particular the sea. Nayak [28] then used his results to relate the moments directly to parameters of contact such as the mean peak height and the curvature of peaks. Nayak also coined the terms 'summit' for a three-dimensional (or areal) maximum as opposed to a profile 'peak'.

The other approach, the discrete technique, was originated by Whitehouse [29]. In its original form the autocorrelation function was constrained to be exponential. The characterization then became simply the RMS σ value of the surface and the correlation length—that length over which the autocorrelation length falls to a small value $\tau^*$. Later this discrete characterization was extended to an arbitrary autocorrelation function and three points $A(0)$, $A(h)$ and $A(2h)$ for a profile where $h$ was an arbitrary distance. Finally, four or more points were used for the areal characterization of the surface. From these the same functional parameters were derived.

The basic difference between these two approaches is that the Nayak model is worked from a continuous signal and hence can be regarded as exact. On the other hand the Whitehouse model, because it is discrete, can only be approximate yet it does mimic exactly what is measured. It is nevertheless a form of surface characterization in its own right. The Whitehouse model appears here in its role as a method of characterizing the surface. It also appears in chapter 3 in its role as a means of explaining the range of values which measured surface parameters can take, it being assumed that all calculations are nowadays carried out digitally. It may seem confusing to have two different methods of evaluating the same thing (namely, the functional parameters). Greenwood [30], who input much of the information on which the functionally important parameters were based, provides a critical view of both methods and sets them in perspective.

Before dealing with this some properties of the random process signals will be demonstrated. Later, in section 2.1.7, some other methods of characterization largely based on random process theory will be described.

The definition of the autocorrelation function is as follows. If $z_1$ and $z_2$ are two measurements taken on the surface a distance $\tau$ apart then the ACF is given by $A(\tau)$ where $A(\tau) = E[z_1, z_2]$, where $E[\ ]$ denotes the expected value of the product of these measurements taken over the entire surface:

$$E[z_1, z_2] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} z_1 z_2 p(z_1 z_2)\mathrm{d}z_1\,\mathrm{d}z_2. \tag{2.60}$$

Providing that two conditions are complied with, the formula $E[z_1, z_2]$ can be modified to a more useful form, that is, one incorporating a profile. The conditions to be satisfied are, first, that the random nature of the surface is uniform in distribution and, second, that the ergodic theorem holds. The first of these conditions is self-explanatory but the second is not. It must suffice here to comment that the real surfaces are close enough to uniformity and ergodicity to justify the use of the ACF. Some examples of processes not satisfying these conditions can be found in communications textbooks.

Under these circumstances the formula may be rewritten as

$$A(\tau) = E[z_1, z_2] = \lim_{L \to \infty} \frac{1}{L} \int_{-L/2}^{L/2} z(x)z(x + \tau)\mathrm{d}x \tag{2.61}$$

For $\tau = 0$, $E[z_1, z_2] = A(0) = \sigma^2$, the variance of the signal, or $R_q^2$ in surface terminology. Because $E[z_1, z_2] = E[z_2, z_1]$ it follows that the function has even symmetry about $\tau = 0$ (figure 2.39).

**Figure 2.39** (*a*) Typical profile; (*b*) autocorrelation function.

Basically, ergodicity means that averages of many surface profiles taken at the same time or *x* value are the same as averages taken of one record over all time. This in turn means that all phase is lost in the signal. It will be shown in section 7.10 that the phase can be retained to advantage in a systems approach. The phase angle $\phi$ is

$$\phi = \tan^{-1}\left(\frac{X(\text{function})}{R(\text{function})}\right) \tag{2.62}$$

but the function is real because the ACF is even. Hence $X(\text{function}) = 0$ and thus the phase is zero.

It is because of this that the ACF has the useful property of being more reliable than deterministic parameters. It is the phase effects and particularly the random phase effects which influence the individual characteristics of a profile such as the highest peak etc. The random amplitudes contain the essential information, the energy information. Consequently, removing, not measuring, the delinquent random element (the phase) improves the stability.

One property of the autocorrelation function not often realized is that it can reveal the 'average machining unit event'. To see how this happens consider shotblasting. The unit event here is the crater left on the surface after the bead has impinged on it. Such a crater is shown in figure 2.40. Alongside it is its ACF. The profile in the figure is a typical profile of a shotblast surface together with its ACF. Notice how little has been changed. The reason for this is that it is by the random positioning in space of the unit event that the profile is built up. This means that the beads hitting the surface are uniformly distributed in space. Heuristically it can be argued that the only difference between them (apart from some degree of geometric interaction) is spatial. This constitutes only phase information and, because the ACF is not responsive to phase, this is ignored. Hence the autocorrelation function itself is merely showing the ACF of the average unit event. The degree to which this is true can be seen from the figure. The only proviso is that the statistics of the manufacturing process have to be Poissonian, that is uniformly distributed in space, which is true for many finishing processes. More of this will appear in chapter 6. The first use of correlation in surface metrology was



**Figure 2.40**

suggested by Wormersley and Hopkins in 1946 [28], who advocated the use of methods previously used for examining turbulent air flow for looking at surfaces. (They were associates of R E Reason at Taylor Hobson. Hopkins joined with Reason in producing the first definitive book on surface metrology in 1944 [2].)

Linnik and Khusu [31] and Nakamura [32] looked at correlation later. However, it was Peklenik [33] who proposed the use of the ACF in terms of a typology for surfaces and, since this time, it has become acceptable to specify surfaces in terms of the statistical properties of surface height rather than the distribution of peaks and valleys, at least for certain applications.

### (a) Parameters of the ACF

Parameters taken from the single-profile graph will be examined first. This in itself will consist of two sets of measurements, one direct from the ACF and the other indirect via power spectra etc. The areal assessment of surfaces using statistical methods will be described.

### (i) Direct parameters

Two parameters emerge directly from the signal-discriminating property of the ACF. First is the measure of how quickly the random element decays and second is the wavelength of the periodicity of the more deterministic component if any. Also, there is the possibility of classifying the type of process by examining the shape of the function. These possibilities have been explored by Peklenik [27] and others. Another function, the structure function $S(\tau) = E(z(x) - z(x+\tau)^2) = 2\sigma^2(1 - A(\tau))$ is also used. It does not need a reference line. Its use is described in section 7.7.4.

The parameters themselves relate to the scale of size of the abscissa of the correlation function. The correlation length, that is the length over which the correlation function drops to a small fraction of its height at the origin, has been defined as that length of surface over which the ACF drops to between 10% and 50% of its original value (see figure 2.41).

That this value appears to be somewhat arbitrary is merely a reflection of the fact that the definition of the value of correlation required for independence depends largely on the proposed function of the workpiece.



**Figure 2.41** Autocorrelation function.

One definition that has been used is

$$\text{correlation distance} = \frac{1}{A(0)}\int_0^\infty |A(\tau)|d\tau \qquad (2.63a)$$

which, for $A(\tau) = \exp(-\tau/\tau^*)$, gives the correlation length to be $\tau^*$ which corresponds to a correlation value of $1/e$, that is 37%.

Although this definition is easy to carry out for a simple exponential function it is not so easy for more complex functions. For example, a typical second-order correlation function is

$$A(\tau) = \frac{\cos \alpha \tau}{1 + \gamma^2 \tau^2}.$$

This has a correlation length which is obtained from

$$\frac{1}{A(0)} = \int_0^\infty \frac{\left|\cos \alpha \tau\right|}{1 + \gamma^2 \tau^2} d\tau \tag{2.63}$$

which results in an expression

$$\frac{1}{\gamma} \left( \cosh \frac{\alpha}{\gamma} \right) \tan^{-1} \left( \frac{1}{\sinh \alpha/\gamma} \right) \tag{2.64a}$$

This is not a simple expression. It could also be argued that the real criterion for independence is where the *envelope* of the correlation function falls to a lower level than $1/e$ rather than the function.

For a purely random surface such as is often found in grinding and similar processes, use can be made of the correlation length to define the number of degrees of freedom contained within a given length of profile. Once this is known it is then possible to determine the average behaviour of the surface. An example of this is given in the chapter on instrumentation, where the degree of stylus integration is worked out using independent elements of the surface. The effect of the skid is also tackled this way. The degrees of freedom so obtained can also be used to estimate the variability of a parameter measured from the profile.

Using the autocorrelation function or power spectral density as part of the basis for a typology is useful. It is relevant to determine the type or family to which the generated surface belongs. Peklenik made the premise that every surface can be described by a basic autocorrelation function. Owing to the ability of the correlation method to separate the random from the periodic components of the surface he proposed a typology system which, by means of the basic autocorrelation functions and/or a combination of them, could meet the requirements. Any typology of the correlation functions or spectral densities has to address the size as well as the shape of the curve. Peklenik proposed five shapes, as shown in figure 2.42.

To augment this classification in terms of size he advocated the correlation length and also the correlation wavelength which relate to the unit of size for the dominant random component and any dominant periodicity respectively. The definition of the correlation wavelength is obvious; that of the correlation length is given in equation (2.63a) and tends to be more important because most finishing processes are random. The problem with such a typology is that to use it functionally requires a deep understanding of the way in which the surface configuration affects the function. This is not straightforward as will be seen in chapter 7. The typology as proposed, however, is more likely to be useful in a manufacturing context because it is much easier to relate periodic elements to machine tool problems of long wavelength and the process in single point cutting for short wavelength and the random element to the process. In any case the two types of shape can be readily tied in with the method of production of the surface. That this typology has not been taken up in practice is due in part to the fact that manufacturers are still reluctant to depart from the idea of a 'one number' surface, mainly because of cost and the practical difficulty of getting personnel to deal with any sort of statistical concept.

Returning to the determination of the 'size' element of the correlation function naturally brings in the concept of the correlation length.

From the previous section it is clear that the autocorrelation and the power spectrum are powerful tools of characterization. Methods based on them which are more speculative will be given later on in section 2.1.7.5. However, there is a word of caution. Under some circumstances (2.33) more than one type of surface can produce the same correlation function.

(a)

| Formula | Autocorrelation | Group |
|---------|-----------------|-------|
| $A(\tau)$=const. <br> $A(\tau)$=cos $\omega\tau$ | | I |
| $A(\tau)$=exp$^{-d\tau}$+cos $\omega\tau$ | | II |
| $A(\tau)$=exp$^{-d\tau}$+cos $\omega\tau$ | | III |
| $A(\tau)$=$\Sigma$(exp+ <br> sin+cos) | | IV |
| $A(\tau)$=exp$^{-\alpha\tau}$ | | V |

(b)

Random profile     Random telegraphic signal     Exponential

Amplitude modulated     Narrow-band frequency or phase modulation     Modulated

**Figure 2.42** (a) Peklenik classification; (b) example of ambiguity.

This example shows a frequency-modulated waveform

$$Z_{\text{Fm}} = \cos\,\omega_c t - D_m\sin\omega_s t\,\sin\omega_c t \qquad (2.64b)$$

which could well be generated by a tool which has yawing vibration.

Compare this with a surface exhibiting amplitude modulation of the form

$$Z_{Am} = \cos \omega_c t + M \cos \omega_s t \cos \omega_c t \tag{2.64c}$$

where $D_m = \Delta f / f_c$ is the modulation index for the case of frequency modulation (assumed to be small) and $M$ is the depth of modulation for the amplitude modulation signal.

Comparison of these two surface profiles shows that they could both have the same sidebands but different phase! They would give the same power spectrum yet they are produced by a different mode of vibration of the tool.

This ambiguous situation changes if $D$ is increased to greater than three because extra sidebands show up with the extra series of the Bessel function describing the mode of modulation.

To separate these two surfaces a classification of the amplitude density curve would have to be added to resolve the confusion. This is shown in the different shapes of APDF shown in figures 2.42(*b*) and 2.43. Kurtosis would make an ideal differentiator in these cases.

Other ambiguities can arise because of the phase insensitivity. A classic example is that of a genuine random surface having an exponential autocorrelation, function and the random telegraphic signal. Both signals in fact have Poissonian statistics but produce different waveforms. This is seen clearly in figure 2.42(*b*).

## (*ii*) *Indirect parameters*

One mathematical tool often used in the derivation of indirect statistical parameters is the joint probability density function (JPDF) and, in particular, the multidimensional normal distribution (MND).

The JPDF between variables $Y_1$, $Y_2$, $Y_3$, etc can be written as $p(Y_1, Y_2, \ldots, Yn,)$ and is taken to be the probability density that $Y_1$ takes the value $Y_1$ when variable $Y_2$ takes the value $Y_2$ and so on. $Y_1$ and $Y_2$ need not be the same type of variable, neither should they necessarily have the same dimensions. For instance, $Y_1$ can be a length and $Y_2$ an angle.

It is often the case that $Y_1$, $Y_2$, etc, are themselves the result of a large number of random variables and under these conditions the central limit theorem holds, that is, the variables can be considered to be Gaussian variates and, furthermore, over a wide range of conditions the joint probability density function will be Gaussian. Assuming each variate has zero mean value the JPDF will have the form $p(y_1, y_2, y_3, \ldots, yn)$ where

$$p(y_1, y_2, \ldots, y_n) = \frac{1}{(2\pi)^{1/2N} |M|^{1/2}} \exp\left( \frac{-\sum_{i,j=1} M_{ij} y_i y_j}{2 |M|} \right) \tag{2.65}$$

where $|M|$ is the determinant of $M$, the square matrix

$$M = \begin{pmatrix} d_{11} & d_{12} & d_{1N} \\ \vdots & \vdots & \vdots \\ d_{N1} & d_{N2} & d_{NN} \end{pmatrix} \tag{2.66}$$

and $d_{ij}$ is the second moment of the variables $y_i y_j$; $M_{ij}$ is the cofactor of $d_{ij}$ in $M$. *This multinormal distribution for random variables is absolutely crucial in the surface properties field* and is used in many places in typology and in chapter 7 on function.

Thus, the point probability density function of two height variables $z_1$, $z_2$ that are correlated by $p$ is given by $p(z_1, z_2)$:

$$p(z_1, z_2) = \left( \frac{1}{2\sqrt{\pi}} \right) \exp\left( -\frac{z_1^2}{2} \right) \left( \frac{1}{\sqrt{2\pi(1-\rho^2)}} \right) \exp\left( \frac{-(z_2 - \rho z_1)^2}{2(1-\rho^2)} \right) \tag{2.67}$$

**Figure 2.43** Modulated signals showing how ambiguity arises.

which is written in the form of a conditional probability density and an elementary probability density

$$p(z_1, z_2) = p(z_1)p(z_2 \mid z_2) . \tag{2.68}$$

That is $p(z_2|z_1)$ is the probability density that $z_2$ takes the value $z_2$ given that $z_1$ has the value $z_1$.
In the unconditional form equation (2.67) is

$$p(z_1, z_2) = \frac{1}{2\pi} \frac{1}{(1-\rho^2)^{1/2}} \exp\left[-\left(\frac{(z_1^2 - 2\rho z_1 z_2 + z_2^2)}{2(1-\rho^2)}\right)\right] \tag{2.69}$$

*(b) Parameters involving crossings of the profile*

Some investigators have derived statistical information from the number of crossings that occur at any level. Much of the original theory is derived from the work of Rice in communication theory [25] and relates to random or pseudo-random waveforms. Longuet–Higgins [26] has expanded the work into two dimensions.

Some typical parameters that have been derived are as follows.

Consider the profile shown in figure 2.44. It has been shown in terms of stationary random process theory by Peklenik [33] that the average number of crossings for unit length of profile is $1/\lambda_a$ where $\lambda_a$ here is not the average wavelength as such but related to crossings at '$a$'

$$1/\lambda_a = \int_0^\infty z' \rho(a, z') \mathrm{d}z' \tag{2.70}$$

and

$$z' = \frac{\mathrm{d}z}{\mathrm{d}x}$$

and the average thickness of the profile at height $a$, is given by $\gamma_a$

$$\gamma_a = \frac{\int_a^\infty p(z)\mathrm{d}z}{\int_0^\infty z' p(a, z')\mathrm{d}z'} . \tag{2.71}$$

$1/\lambda_a$ corresponds to the high-spot count sometimes used as a parameter of the surface. Notice that the term for $\gamma_a$ is normalized with respect to a unit length.

**Figure 2.44** Profile-crossing parameters.

For the case when the profile is Gaussian in form and of zero mean

$$p(z_1, z_2) = \frac{1}{\sigma\sqrt{2\pi}}\exp\left[-\left(\frac{z_1^2}{2\sigma^2}\right)\right]\frac{1}{\sigma\sqrt{2\pi}}\exp\left[-\left(\frac{z_2^2}{2\sigma^2}\right)\right]$$

(2.72)

where $z_1$ and $z_2$ are considered to be independent and hence $\rho = 0$.

From this Gaussian form $\sigma'$ is the standard deviation of the slope and $\sigma$ is the standard deviation of the profile:

$$\gamma_a = \frac{\sigma'}{\sigma\sqrt{2\pi}}\exp\left[-\left(\frac{a^2}{2\sigma^2}\right)\right]$$

(2.73)

$$\gamma_a = \pi\frac{\sigma'}{\sigma}\exp\left[-\left(\frac{a^2}{2\sigma^2}\right)\right]\left[1 - \Phi\left(\frac{a}{\sigma}\right)\right].$$

(2.74)

From these

$$\sigma' = c_a\sigma\lambda_a$$
$$\gamma_a = c_a'.1.\lambda_a$$

(2.75)

where $c_a$ and $c_a'$ are constants for any given level at which the crossings have been counted and $\gamma_a$ is the average number of crossings per unit length at level $a$.

Information about the correlation function is inherent in $\sigma'$ because it is well known that the variance of $z'$ is related to the correlation function by the formula.

$$\sigma'^2 = -\frac{d^2 A(0)}{d\tau^2}.$$

For $a$ equal to the mean of the surface, $c_a = 2\pi$ and $c_a' = 1/\pi$. For $a$ equal to the $\sigma$ level above the mean line, $c_a = 10.3$ and $c_a' = 0.16$. For instrumental reasons some sort of offset from the mean line is usually made. This is in case electrical noise is present which can happen at high magnifications.

The advantage of such an approach is that deterministic surfaces can be treated in a similar way.

A similar approach has been used by Sankar and Osman [17] who proposed a parameter called the MCE, the mean crest excursion. In addition they proposed the RMSCE, the root mean square crest excursion. Both of these are derived from $p(\lambda_a)$, the probability density that the profile has a crest (or high-spot) width less than or equal to the specified value of $\lambda_a$ about $z$. Thus

$$\text{MCE} = E|\lambda_a| = \int_{-\infty}^{\infty}\lambda_a p(\lambda_a)d\lambda_a$$

$$\text{RMSCE} = \left(\int_{-\infty}^{\infty}(\lambda_a - \text{MCE})^2 p(\lambda_a)d\lambda_a\right)^{1/2}.$$

(2.76)

Similar expressions have been developed for the valleys but only give the same information for a random wave. (Note that $\lambda_a$ is not the average wavelength as defined earlier and described below.)

As in the previous example they proposed counting profile crossings at a height above the mean line. This time, however, they chose a height of the $R_a$ value rather than the RMS $R_q$ value above the mean line.

### 2. 1. 3. 4  *Power spectrum*

Continuing the theme of a wavelength-conscious parameter naturally leads to methods which do not attempt to characterize the autocorrelation function but to condense the information in the power spectral density.

One method due to Spragg and Whitehouse [13] involves the estimation of the moment of gyration of the power spectrum (the RMS frequency) by a direct measurement of the RMS slope and the profile RMS value.

It has been shown that if $\omega_{RMS}$ is the RMS angular frequency

$$\omega_{RMS}^2 = \frac{\int_0^\infty \omega^2 P(\omega)\mathrm{d}\omega}{\int_0^\infty P(\omega)\mathrm{d}\omega}. \tag{2.77}$$

Similarly, it can be shown that

$$E(z'^2) = \sigma'^2 \tag{2.78}$$

$$= (R_q')^2 = \frac{1}{2\pi}\int_{-\infty}^\infty \omega^2 P(\omega)\mathrm{d}\omega. \tag{2.79}$$

The power spectral density is related to the autocorrelation function by means of Fourier transforms. Thus the well-known relationship is

$$P(\omega) = \int_{-\infty}^\infty A(\tau)\cos(\omega\tau)\mathrm{d}\tau \tag{2.80}$$

$$A(\tau) = \frac{1}{2\pi}\int_{-\infty}^\infty P(\omega)\cos(\omega\tau)\mathrm{d}\omega. \tag{2.81}$$

The power spectrum $P(\omega)$ is related to the Fourier spectrum of the surface by

$$A(\tau) = \lim_{L\to\infty} \frac{1}{L}\left|F(\omega)\right|^2. \tag{2.82}$$

The value of $A(\tau)$ at the origin is the zero moment $\sigma^2$, the variance of the surface equal to $R_q^2$.

Equations (2.80) and (2.81) are cosine transforms because both $P(\omega)$ and $A(\tau)$ are even—the phase has been destroyed. This is one of the advantages of using random processes—one less random variable is present. Consequently, $P(\omega)$ the power spectral density is stable.

Hence

$$\frac{1}{\omega_{RMS}} = \frac{\lambda_{RMS}}{2\pi} \tag{2.83}$$

and

$$\lambda_q = \lambda_{RMS} = 2\pi\frac{R_q}{R_q'} \tag{2.84}$$

where $R_q = \sigma$, $R_q' = \sigma'$.

This technique, as in the level-crossing methods for autocorrelation, is completely general, providing that the characteristics are stationary over the evaluation length. For convenience the $\lambda_q$ value is often replaced by the $\lambda_a$ value, the average wavelength alternative.

Some examples of the parameter $\lambda_a$ value are shown in figure 2.45. The essential practical detail that needs to be taken into consideration in autocorrelation as well as PSD measurement is the high-frequency content of the signal allowed through by the instrument. The change in parameter value as function of frequency cut is profound. As an example the way in which $\bar{\omega}_{RMS}$ changes with high cut $c$ is given by [12]:

$$\frac{\mathrm{d}}{\mathrm{d}c}(\bar{\omega}_e^2) = P(c)\int_0^c (c^2 - \omega^2)P(\omega)d\omega / (SP(\omega)(\mathrm{d}\omega)^2. \tag{2.84b}$$



**Figure 2.45** Average wavelength and power spectrum.

It will be demonstrated in a later section (2.4) how it is possible to build up a comprehensive measuring philosophy in roundness based upon such a system.

### 2.1.3.5  Peak and valley definitions and characteristics

Some characteristics of interest to the metrologist have already been studied in different disciplines, for example, the distribution of maximum values in statistics as in communication theory (Rice [12]).



**Figure 2.46** Surface contact.

The conventional way to calculate average peak height from a surface profile is to split the profile vertically into equal height intervals and then to count the number of peaks within each interval over the profile length. These peak counts, sometimes called peak frequencies, are plotted as a function of height; the resulting graph is called the peak distribution histogram. If these counts at each level are divided by the total number of peaks within the profile length assessed the plot becomes a peak probability density. This is possible because the area of the plot is unity. The average peak height is obtained from this by working out the height of the centre of gravity of the graph from the height datum.

In the above approach the count of peaks within each height interval can be regarded as an attribute. Each peak is only registered by its presence. In what follows it is suggested that there is another way to arrive at an estimate of peak probability. This method does not derive the peak probability curve using the peak count. Instead the 'extent' of profile surrounding each peak which satisfies certain discrete peak criteria is measured and summed for each height. These sums of 'extents' are normalized within each interval for all heights by the total length of assessed profile. In this way another way of generating a peak probability curve is realized in which the attribute of peak presence has been replaced by the measurement of a peak property (i.e. the extent). It turns out that this 'extent' of peak-like property is related very strongly to peak curvature (i.e. each peak is 'weighted' by its curvature). It has to be pointed out that, in the direct probability approach, neither the actual 'extent' of peaks or their curvature are measured. These peak properties are intrinsically contained within the definition of probability (see for example equation 2.85 below) and the limits of the integration. The fact that there is *no need* to measure curvature of the peaks and yet still include the effect of curvature (i.e. the extent along the profile describing the 'thickness' of the peak) is very attractive (see equation 2.88). In the following sections it will be shown that the values of average peak heights (and RMS heights) are different using this measure of probability. It is suggested in this paper that the two methods are complementary. The peak count method may be more useful in electrical and thermal contacts because the actual numbers of contacts are involved whereas the latter could be useful in loaded contacts (figure 2.46) [1].

The starting point is the joint probability density of the surface height, the slope and the second differential $z''$, (here approximated as the local curvature). This is represented as $p(z, z', z'')$. At this stage the joint probability density function is quite general.

The probability of peaks occurring between $zs$ and $z_s + \delta z$ having a curvature of between $z''$ and $z'' + \delta z''$ and a slope between o and $\delta z^{l}$ is therefore

$$P(\hat{z}) = \int_{z_s}^{z_s + \delta z} \int_0^{\delta z'} \int_0^{z_s + \delta z''} p(z, z', z'') \mathrm{d}z'' \mathrm{d}z' \mathrm{d}z \qquad (2.85)$$

$$= \delta z' p'(0) . p(z_s, z'') \delta z'' \, \delta z \qquad (2.86)$$

In equation (2.85) $\delta z$ is small when compared with $z$ and $\delta z''$ is small when compared with $z''$ (i.e. much more than an order of magnitude smaller) and $p'(o)$ is the probability density of the slope term at between zero and $\delta z'$.

The local length of surface over which this peak condition holds is given by $\delta z / z''$. The term $\delta z / z''$ is in effect the 'thickness' of a peak of curvature $z''$ between height $z$ and $z + \delta z$. and having a slope of between $o$ and $\delta z\, z'$ (figure 2.47).

Given the probability of a peak, the number of peaks per unit length, $N(z_s)$, is obtained by dividing equation (2.86) by $\delta z' / z''$ in the same way as Bendat [34]. Thus

$$N(z_s) = \delta z' p'(o) p(z, z'') \delta z'' \delta z / \delta z' / z''$$
$$= p'(o) p(z, z'') z'' \delta z'' \delta z \qquad (2.87)$$

The peak count is basically the *number of times* that the profile height falls between $z_s$ and $z_s + \delta z$ when the slope of the profile is between 0 and $\delta z'$ and the curvature is between $z''$ and $z'' + \delta z''$ per *unit length*. The distance information is not omitted. For a given $z$ and slope range (i.e. $\delta z'$ the count for unit distance) is determined for

**Figure 2.47** Peak weighting.

a given $z''$ by the extent of $x$ in which $z''$ is in the o to $\delta z'$ range (figure 2.47). The count is remembered but the $x$ values are not. So for the *peak probability* the $x$ values are remembered and the peak count is not and for *peak count* equation (2.87) the count is remembered but the $x$ values are not. The two parameters are not the same: they are absolutely complementary. It may be that a parameter derived from both approaches is possible.

It is useful to compare equation (2.87) with equation (2.86). At first glance this is a comparison of a count of peaks with a probability of peaks. There is, however, a completely different interpretation of the probability equation (2.86); it can be construed as the sum of the count of peaks each weighted by their 'thickness.' Take for example (2.88) for the probability of a peak at a height between $z$ and $z + \delta z$,

$$\underbrace{\delta\, z' p'(o) p(z,\ z'') \delta\, z''.\delta z}_{probability} \Rightarrow \underbrace{(p'(o).p(z,\ z''))z''\delta\, z''}_{count}\ \underbrace{\delta\, z.(\delta\, z'/z'')}_{thickness} \tag{2.88}$$

In words, the probability of a peak $P(\hat{z})$ does not only depend on the number of peaks but it also automatically takes into account their physical dimension as revealed by the radius of curvature $1/z''$. The total peak probability between $z$ and $z + \delta z$ is given by $P(\hat{z})$ where

$$P\,(z_s) = \delta z' p'(o) \int_{z_s}^{z_s + \delta z} \int_{-\infty}^{0} p(z, z'') \mathrm{d}z'' \mathrm{d}z \tag{2.89}$$

which represents what is in effect the count of all peaks at $z$ each weighted according to their radius of curvature $1/z''$. This is a measure of the total length over which the profile at height between $z$ and $z + \delta z$ satisfies the peak criterion. When divided by the profile assessment length, this equation (2.89) represents the peak probability.

For any height between $z$ and $z + \delta z$ there will be a family of probabilities each made up of counts of peaks, and each family having associated with it a different radius of curvature. Exactly how large the

probability is depends on the form of $p(z, z'')$. No assumptions about the nature of $p(z)$, $p(z, z'')$ or $p(z, z'\ z'')$ are needed.

A point to be aware of is that the mechanism of contact is made up of two elements, the behaviour at a 'typical' contact and the distribution of such contacts in space. The *probability* of peaks at $z$ gives in effect the total length of 'peak-like' surface at $z$ which has the capability of supporting normal loads. It does this because the definition has first isolated the peaks and then lumps them together after multiplying by $1/z''$. This weights them in importance according to $1/z''$. The peak count by contrast is an estimate of the distribution of peaks (i.e. how the number of peaks per unit length varies with height). The total length of $x_s$ values mentioned above is expressed as a fraction of the whole length of profile.

So the peak probability and peak count taken together can satisfy a functional requirement for surface parameters, (i.e. represent both the unit functional event (contact points) and their distribution in space).

Consider now the peak count. The total number of peaks per unit distance at height $z$ is given by $N(z_s)$ where

$$N(z_s) = p'(o) \int_{z_s}^{z_s+dz} \int_{-\infty}^{0} p(z, z'')z''\,\mathrm{d}z''\mathrm{d}yz \tag{2.90}$$

and is a simple count of all peaks at $z$ irrespective of their curvature divided by the total profile length. Equation (2.90) should be compared with equation (2.89). It is evident that probability and frequency of occurrence are not the same. This can be seen in equation (2.88).

The standard way of evaluating the average peak height is via the count (or frequency of occurrence). Thus the average peak height obtained conventionally, $\bar{\bar{z}}_c$ is given by

$$\bar{\bar{z}}_c = \delta\,zp'(0) \int_{-\infty}^{\infty} \int_{-\infty}^{0} z.p(z,z'')z''dz''dz \left/ \delta\,z'p'(o) \times \int_{-\infty}^{\infty} \int_{-\infty}^{0} p(z,z'')z''dz''dz \right. \tag{2.91}$$

An alternative way using probabilities of peak rather than the frequency of occurrence of peaks is suggested here. This gives $\bar{\bar{z}}_p$ where

$$\bar{\bar{z}}_p = \delta\,z'p'(0) \int_{-\infty}^{\infty} \int_{-\infty}^{0} zp(z,z'')dz''dz \left/ \delta\,z'\,f'(o) \times \int_{-\infty}^{\infty} \int_{-\infty}^{0} p(z,z'')dz''dz \right. \tag{2.92}$$

Only if the value of curvature $z''$ is fixed or its range severely restricted or if $z''$ is independent of height $z$ can equation (2.91) be equal to equation (2.92). As none of these are probable in practice, $\bar{\bar{z}}_p \neq \bar{\bar{z}}_c$ which is a general result and not dependent on the form of $p(z, z'')$. $\bar{\bar{z}}_p$ could be used as an alternative measure of average peak height rather than $\bar{\bar{z}}_c$, and peak probability as explained above could perhaps be more meaningful than peak count in contact situations. The fact is that for any given $z$ there are a number of peaks with curvature $z''_1$, say $n_1$, a number with curvature $z''_2$, say $n_2$, and so on which are taken into account with peak probability but not peak count. The number of peaks at $z$ is $n_1 + n_2 + \ldots$ . This is just the number of peaks at $z$. Nothing else is recorded.

If the waveform is $z$ and its derivatives are $z'$ and $z''$

$$z' = \frac{\mathrm{d}z}{\mathrm{d}x} \qquad z'' = \frac{\mathrm{d}^2z}{\mathrm{d}x^2}. \tag{2.93}$$

The probability for a peak occurring at a height between $z_1$ and $z_1 + dz$ is $dz_1$

$$\mathrm{d}z_1 = \int_{-\infty}^{0} p(z_1, 0, z'')z''\,\mathrm{d}z''. \tag{2.94}$$

This equation shows the restraints for a maximum, namely that $z' = 0$ and $z'' < 0$. Similar equations hold for the valleys. Note that all curvatures are lumped together!

In the case of random waveforms this distribution approaches a Gaussian form due to the central limit theorem.

Notice the similarity between this formula and the one for the average number of crossings per unit length of the profile.

The coefficients in the correlation matrix **M** of the multinormal distribution of equation (2.65) become

$$E(z^2) = \sigma^2 \qquad E(z')^2 = \left(\frac{-\mathrm{d}^2 A(\tau)}{\mathrm{d}\tau^2}\right) = \sigma'^2 = A''(0) \qquad E(z, z') = 0$$

$$E(z', z'') = 0 \qquad E(z, z'') = \left(\frac{\mathrm{d}^2 A(\tau)}{\mathrm{d}\tau^2}\right) = -A''(0) \qquad E(z''^2) = \left(\frac{\mathrm{d}^4 A(\tau)}{\mathrm{d}\tau^4}\right) = A^{\mathrm{iv}}(0). \tag{2.95}$$

For simplicity the notation below is used:

$$\left.\frac{\mathrm{d}^n A(\tau)}{\mathrm{d}\tau^n}\right|_{\tau=0} = A^n(0) \tag{2.96}$$

Then the multivariable matrix **M** becomes

$$\begin{vmatrix} \sigma^2 & 0 & A''(0) \\ 0 & -A''(0) & 0 \\ A''(0) & 0 & A^{\mathrm{iv}}(0) \end{vmatrix} \tag{2.97}$$

which yields cofactors

$$m_{11} = -A''(0)A^{\mathrm{iv}}(0) \qquad m_{13} = (A''(0))^2 \qquad m_{33} = -A''(0)\sigma^2 \tag{2.98}$$

from which the probability density becomes

$$p(z_1, 0, z'') = \frac{1}{2\pi^{3\psi 2}|\mathsf{M}|^{1\psi 2}} \exp\left(-\frac{1}{2|\mathsf{M}|}(m_{11}z_1^2 + m_{33}z0^2 + 2m_{13}z_1 z'')\right). \tag{2.99}$$

The expected number of peaks per unit length at a height in between $z_1$ and $z_2 + \mathrm{d}z$ is given by $p_p(z_1).\mathrm{d}z_1$

$$p_{\mathrm{p}}(z) = \frac{\mathrm{d}z_1}{2\pi^{3/2}M_{33}}\left\{|\mathsf{M}|^{1/2}\exp\left(-m_{11}\frac{z^2}{z|\mathsf{M}|}\right) + m_{13}z_1\left(\frac{\pi}{2m_{33}}\right)^{1/2}\exp\left(\frac{z^2}{2\sigma^2}\right)\left[1 + \mathrm{erf}\left(\frac{m_{13}z_1}{z|\mathsf{M}|^{1/2}m_{13}}\right)\right]\right\} \tag{2.100}$$

where

$$\mathrm{erf}(x) = \frac{2}{\sqrt{\pi}}\int_0^x \exp(-t^2)\mathrm{d}t. \tag{2.101}$$

For large $z$ it can be shown that the expected number of peaks per unit distance lying above the height $a$, $P_p(a)$, is given by

$$P_p(a) = \frac{1}{2\pi}\left(\frac{-\sigma'^2}{\sigma^2}\right)^{1/2}\exp\left(\frac{-a}{2\sigma^2}\right) \qquad (\sigma'^2 = E(z'^2)). \tag{2.102}$$

This corresponds to the number of positive crossings at height $a$. The total number of peaks per unit distance of any height is obtained from equation (2.94) by integrating with respect to $z$ and $z''$, yielding

$$\frac{1}{2\pi}\left(\frac{A^{\mathrm{iv}}(0)}{-A''(0)}\right)^{1/2} = \frac{1}{2\pi}\left(\frac{\int_0^\infty \omega^4 P(\omega)\mathrm{d}\omega}{\int_0^\infty \omega^2 P(\omega)\mathrm{d}\omega}\right)^{1/2} = \frac{1}{2\pi}\left(\frac{m_4}{m_2}\right)^{1/2}. \tag{2.103}$$

Furthermore it can also be shown that the number of zero crossings of the profile $N$ is given by

$$N = \frac{1}{\pi}\left(\frac{\int_0^\infty \omega^2 P(\omega)\mathrm{d}\omega}{\int_0^\infty P(\omega)\mathrm{d}\omega}\right)^{1/2} = \frac{1}{\pi}\left(\frac{-A''(0)}{A(0)}\right)^{1/2} = \frac{1}{\pi}\left(\frac{m_2}{m_0}\right)^{1/2} \tag{2.104}$$

and in general the average number of zero crossings of the $k$th derivative of the profile is given by

$$N_k = \frac{1}{\pi}\left(\frac{\int_0^\infty \omega^{2k+2} P(\omega)\mathrm{d}\omega}{\int_0^\infty \omega^{2k} P(\omega)\mathrm{d}\omega}\right)^{1/2} = \frac{1}{\pi}\left(\frac{-A^{2k+2}(0)}{A^{2k}(0)}\right)^{1/2} \tag{2.105}$$

from which, by using Maclaurin's series,

$$A(\tau) = A(0) + \frac{\tau^2}{2!}A''(0) + \frac{\tau^4}{4!}A^{\mathrm{iv}}(0) + \cdots. \tag{2.106}$$

In principle, by measuring $N_i$, $i = 1, 2, \ldots, k$, the autocorrelation function can be found from this equation.

This assumes that the profile is Gaussian and that the autocorrelation is differentiable at $\tau = 0$, which is not the case for exponential correlation functions, nor for fractal surfaces where all the derivatives at the origin are undefined.

Other peak characteristics of particular significance in surface topography can be found in terms of the ratio of maxima to crossings on the profile.

### (a) Peak density above mean line

For example, the peak density above the mean line is given by $\int_0^\infty p_p(z)\,\mathrm{d}z$ where $p_p(z)$ is as given in equation (2.100). Here the peak density as defined by [35] simply counting peaks.

After some manipulation, the peak probability density $p_p$ above the mean line [24] is

$$p_p = \frac{1}{4\pi}\left[\left(\frac{m_4}{m_2}\right)^{1/2} + \left(\frac{m_2}{m_0}\right)^{1/2}\right] = \frac{1}{4\pi}\left[\left(\frac{A_0^{\mathrm{iv}}(0)}{-A''(0)}\right)^{1/2} + \left(\frac{A''(0)}{\sigma_2}\right)^{1/2}\right] \tag{2.107}$$

where $m_4$, $m_2$, $m_0$ are the fourth, second and zeroth moments of the power spectrum, as in equation (2.104), which gives the result that $p_p$ = density of maxima over mean line = $\frac{1}{2}$ total maxima density + $\frac{1}{4}$ density of crossings. Also, density of maxima below mean line = $\frac{1}{2}$ total maxima density − $\frac{1}{4}$ density of crossings.

A corollary of this is that

$$p_{p+} - p_{p-} = \tfrac{1}{2} \text{ density of crossings} \tag{2.108}$$

$$p_{p+} - p_{p-} = 1/\lambda_q. \tag{2.109}$$

This definition of RMS wavelength shows how it may be obtained without having recourse to high-spot counting which may not be as accurate as peak counting on a profile obtained from a real surface.

*(b) Mean peak height $\overline{z}_p$*

$$\overline{z}_p = \frac{\int_{-\infty}^{\infty} z p_p(z) \mathrm{d}z}{\int_{-\infty}^{\infty} p_p(z) \mathrm{d}z}. \tag{2.110}$$

This can be evaluated using equation (2.100) to give

$$z_p = \frac{\sigma k}{2} \sqrt{\frac{\pi}{2}} \tag{2.111}$$

where $k$ is the ratio of density of zero crossings to maxima in unit length.

The mean peak height *above* the mean line can be similarly found. Thus

$$\overline{z}_p = \frac{\int_0^{\infty} z p_p(z) \mathrm{d}z}{\int_0^{\infty} p_p(z) \mathrm{d}z}. \tag{2.112}$$

$$= \sigma \sqrt{2/\pi} \frac{k}{2+k} \left[ \frac{2}{k} \sqrt{1-(k/2)} + \frac{\pi}{2} + \sin^{-1}\left( \frac{k}{2} \right) \right]. \tag{2.113}$$

Also

$$\overline{z}_p = -\sigma \sqrt{2/\pi} \frac{k}{2-k} \left[ \frac{2}{k} \sqrt{1-(k/2)} - \frac{\pi}{2} + \sin^{-1}\left( \frac{k}{2} \right) \right]. \tag{2.114}$$

Figure 2.48 shows how the respective peak height properties change with $k$. The physical interpretation of $k$ is that of a bandwidth parameter. It tells the *shape* of the spectrum. When $k = 2$ it corresponds to a very narrow bandwidth surface; when $k = 0$ a wideband spectrum occurs as seen in figure 2.49.

By changing $k$ many surfaces can be simulated. Therefore the tribological behaviour of many types of surface varying from the purely random $k \sim 0$ to nearly deterministic $k \sim 2$ can be investigated. The $k$ value is a characterization of the surface and is especially useful in contact problems involving peak and valley behaviour. This is in effect a typology based upon random process theory. This concept of taking the ratio of the zero crossings to the peak count was used extensively in the USSR [36] by Lukyanov.

As can be seen with reference to the definitions given so far the value $k$ can be related to $S$ and $S_m$

$$k = 2S/S_m.$$

This $k$ is not to be mistaken for the $k$ surfaces of sea reflection discussed later in the context of optical scatter in chapter 7.

**Figure 2.48** Properties of peaks as a function of the type of surface.



**Figure 2.49** Shape of power spectrum as function of surface characterizer $k$.

The use of $k$ (or $\alpha$ described in areal parameters) is a means of making a typology from the shape of the power spectral density rather than using the Peklenik typology based on the correlation function. They are, however, both doing the same thing only in a slightly different way. In the foregoing analysis the importance of the zero-crossing approach should be noted. The reason for this is the ease with which crossings (and peak counts) can be measured from a chart. The count does not have to be very accurate in order to get some idea of the density of crossings.

To pursue this crossing initiative further, consider the possibility of estimating the correlation function itself from crossing theory. It can be done quite easily as will be seen. Apart from being useful as an estimate such techniques can be useful as checks on computed results.

### 2.1.3.6  *Cumulative distributions and peaks counts*

The well known curve describing the way in which the material to air ratio changes with height has many names. It started as the Abbott–Firestone curve and now is called the material ratio curve $MR(z)$. It is defined as

$$MR(z) = \int_0^\infty p(z)dz \qquad (2.115)$$

where $p(z)$ is the probability density of the profile taking the value $z$. This parameter is one of the most used of all parameters and is the basis for many modern parameters namely $R_{pk}\, R_{vk}$ [37].

The problem with the material ratio curve is that it relates to the whole profile and does not consider peaks as such. What is proposed here is a variant of MR based on peak probability rather than the whole profile probability.

Thus the peak probability ratio $PPR(\bar{z})$ is given by

$$PPR(\bar{z}) = p'(o)\delta z' \int_0^\infty p(y) \int_0^\infty p(z''/z)\mathrm{d}z''\mathrm{d}z \qquad (2.116)$$

by expressing $p(z, z'')$ in the conditional form $p(z).p(z''/z)$.

Normalizing equation (2.116) gives equation (2.117) in which it is seen that MR information is not lost. It is in effect modulated by the peak information.

$$PPR(\bar{z}) = \underbrace{\int_{\bar{z}}^\infty p(z)}_{MR(z)} \underbrace{\int_\infty^0 p(z)\mathrm{d}z''\mathrm{d}z}_{peak\ factor} \qquad (2.117)$$

Obviously, the mechanism of measuring $PPR(\bar{z})$ from a profile would be slightly more complicated than that of MR but not prohibitively so. As each peak between $z$ and $z + \delta z$ is recognized and counted it is weighted (i.e. multiplied by the measured value of the radius of curvature for that peak) before being added to the accumulator representing the probability value at that level. In practice this weighting is not usually necessary. It is taken account of automatically, as will be seen later.

It is suggested here that this curve of $PPR(\bar{z})$, could possibly be more meaningful than the MR curve in some contact situations and possibly in wear and friction because it is highlighting exactly that part of the waveform which is capable of supporting normal load. $PPR(\bar{z})$, in effect pools the advantages of a peak (which is the latter part of (2.117)) with that of the material ratio (which is the first part of the same equation).



**Figure 2.50** Apparent peak loading.

The corresponding cumulative peak count ratio $PCR(\bar{z})$ gives only the number of peaks down to level $(\bar{z})$ as a fraction of the total number of peaks i.e.

$$PCR(\bar{z}) = \int_0^\infty p(z) \int_0^\infty p(z''/z)z\mathrm{d}z''\mathrm{d}z \qquad (2.118)$$

Figure 2.50 shows the mechanism of the loading of a peak on the surface. It is seen that the effect on the larger peaks is straightforward. An elastic compliance is shown. what is also shown is a sympathetic compliance with peaks not being contacted. This possibility is discussed later when considering the Gaussian case.

**Figure 2.51** Discrete model of profile.

Figure 2.51 shows that if the quantization interval is $\delta z$ then the probability of peaks at $z$ is obtained by adding the width of *peaks* between the interval $z$ and $z + \delta z$ rather than the width of the profile at $z$ as in the material ratio evaluation. The $PCR(z)$ is therefore easy to measure.

In order to get some feel for the proposal it is useful to consider the case when the probability densities can be considered to be Gaussian. This is not a bad assumption in practice because many finishing processes have this property and, in any case, the gap between contacting surfaces is more than likely to have a Gaussian characteristic than either individual surface.

*Gaussian statistics (for profile)*

The multi-normal distribution is used to calculate the relevant statistics (after Rice) [26]. Thus

$$p(z, o, z'') = \frac{1}{(2\pi)^{3/2}} \frac{1}{|M|^{1/2}} \exp\left(-\frac{1}{2|M|}(m_2 m_4 z^2 + 2m_2^2 z z'' + m_0 m_2 z''^2)\right)$$

$$(2.119)$$

where $|M| = m_2 (m_0 m_4 - m_2^2)$, $m_0$, $m_2$ and $m_4$ are the variances of $z$, $z'$, and $z''$ respectively.

*Average peak heights (Gaussian Analysis)*

Using equation (2.119) in equations (2.91) and (2.92) and noting that $p'(0)$ is removed from the distribution by normalization the well known result for the average peak height is found. Thus $\bar{\bar{z}}_c$ is given by

$$\bar{\bar{z}}_c = \sqrt{\frac{\pi}{2}} \frac{m_2}{\sqrt{m_4}}$$

$$(2.120)$$

This result has been well documented for many years. See for example Bendat [34]. If, however, the average is found by probability of peaks and not frequency of peaks $\bar{\bar{z}}_p$ results.

Where

$$\bar{\bar{z}}_p = \sqrt{\frac{2}{\pi}} \frac{m_2}{\sqrt{m_4}}$$

$$(2.121)$$

Notice that $\bar{\bar{z}}_p$ is much lower than $\bar{\bar{z}}_c$ in the profile waveform illustrating that more emphasis has been given to the 'wider' peaks found lower in the profile than $\bar{\bar{z}}_c$ which assumes all peaks are of equal weight (i.e. each peak has the value unity).

The ratio of equation (2.120) to equation (2.121) is the considerable value of $\frac{\pi}{2}$! This represents a very considerable difference in height. It could be interpreted, therefore, that the average height of peaks weighted

by their radius of curvature is a factor of $\pi/2$ lower than the average height of all peaks. However there is another interpretation.

Letting
$$b = \frac{m_2}{\sqrt{m_0 m_4}}$$
(2.122)

where $b$ is a measure of the bandwidth of the surface spectrum gives the value of $\bar{\bar{z}}_c$ as

$$\bar{\bar{z}}_c = \sqrt{\frac{\pi}{2}} b \sqrt{m_0}$$
(2.123)

and
$$\bar{\bar{z}}_p = \sqrt{\frac{2}{\pi}} b \sqrt{m_0}$$

If $b = 1$ which corresponds to a narrow spectral band surface such as might be found in general turning or diamond turning.

then
$$\bar{\bar{z}} = \sqrt{\frac{2}{\pi}} \sqrt{m_0}.$$
(2.124)

Under these circumstances when $b = 1$, (i.e. nominally a sine wave surface) the average height of the 'significant or weighted' peak is equal to the $R_a$ value of the surface—the arithmetic average for the profile above the mean line. In this interpretation obviously the symmetrical $R_a$ level below the mean line is invalid. In equation (2.124) the value of $\sqrt{m_0}$ is $R_q$ the root mean square deviation of the profile about the mean line.

Using the same example, $\bar{\bar{z}}_c$, the 'all peak' average height, can be interpreted as the amplitude of the dominant waveform in the narrow band spectrum (i.e. the amplitude of the feed mark) because it is $\pi/2$ larger than the average ($R_a$) of such a waveform; it corresponds in fact to $R_p$.

The interpretations of $\bar{\bar{z}}_p$ and $\bar{\bar{z}}_c$ can in this specific case of $b = 1$ be identified with the traditional surface parameters $R_a, R_q$ and $R_p$. It is interesting to note that for this type of surface the well known and often used $R_a$ value has a peak interpretation as well as the traditional profile representation [1, 2].

One question that has not been answered in the above section is how the cumulative distributions of peak probability and peak count compare with the material ratio curve for the Gaussian case. This will be addressed in the next section.

*Cumulative distributions (Gaussian analysis)*

The material ratio curve is well known and often used. It has the form given in equation (2.125) and equation (2.93). Thus

$$MR(\bar{z}) = \int_{\bar{z}}^{\infty} p(z) \mathrm{d}z$$
(2.125)

which, for a Gaussian distribution, becomes

$$MR(\bar{z}) = \int_{\bar{z}}^{\infty} \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{m_0}} \exp\left(-\frac{z^2}{2m_0}\right) \mathrm{d}z$$
(2.126)

$$= \frac{1}{2}\left[1 - erf\left(\frac{z}{\sqrt{2m_0}}\right)\right] \tag{2.127}$$

This distribution is obtained from all of the profile signal and does not separate peaks.

The comparable cumulative peak probability for a Gaussian distribution of the profile is $PPR(z)$

where
$$PPR(z) = \int_{z}^{\infty} P(\hat{z})\mathrm{d}z$$

where the probability of a peak within the interval $z$ and $z + \delta z$ is

$$\frac{1}{4\pi\sqrt{m_0 m_2}}\exp\left(-\frac{z^2}{2m_0}\right)\left(1 + erf(HKz)\right)\delta z \tag{2.128}$$

The cumulative version of equation (2.128) appears not to have a closed solution. However, some idea of its form can be obtained by taking just the first term of the series for the error function.

This gives the cumulative peak distribution $PPR$

$$PPR(z) \approx MR(z) + \sqrt{\frac{2}{\pi}}\frac{m^2}{m_0 m_4}M^{1/2}\exp(-P^2 K^2 z^2) \tag{2.129}$$

where

$$H = \sqrt{\frac{m_2}{m_0}},\ P = \sqrt{\frac{m_4}{m_2}},\ K = \frac{m_2}{\sqrt{2\mid M \mid}}\ \text{ and } MR(\text{-}) \text{ is the material ratio.}$$

Equation (2.129) clearly shows the two components making up the $PPR$; the material ratio term plus a term involving peaks only. ($P$ is a measure of closely correlated peaks above $z$.)

The equation therefore shows in figure 2.52 that any load could be resisted at a higher level using the peak probability than it would for the material ratio curve. It is the contention in this paper that the $PPR$ approach could be more likely to satisfy the practical situation of contact and wear than $MR(\text{-})$! The important point is that there is more information in the statistics of the surface profile and no doubt also the areal geometry than has hitherto been extracted.

The value of the peak count between $z$ and $z + \delta z$ is well known [34] but it can be expressed in a slightly different form.

Thus if $M$ is the determinant of $p(z, z', z'')$ then,

$$N_c(\hat{z}) = \frac{1}{(2\pi)^{3/2}}\frac{M^{1/2}}{m_0 m_4}\exp(-P^2 K^2 z^2) + \frac{H}{4\pi}\cdot\frac{-z}{m_0}\exp(\frac{-z^2}{2m_0})(1 + erf(HKz)) \tag{2.130}$$

In equation (2.130) $P = \sqrt{\frac{m_4}{m_2}}$ and $H = \sqrt{\frac{m_2}{m_0}}$ as in equation $\tag{2.129}$

The count ratio down to $z$ is given by $\int_{z}^{\infty} N_c(z)dz = PCR(z)$ $\tag{2.131}$

$$PCR(\not{z}) \;\; = \underbrace{\frac{1}{4\pi}\, P(1 - erf\,(PK\not{z}))}_{A} + \underbrace{\frac{1}{4\pi}\, H\, \exp(-\frac{\not{z}^{\,2}}{2m_0})(1 + erf\,(HK\not{z}))}_{B} \qquad (2.132)$$

Equation (2.132) breaks down in a simple way. The components $A$ and $B$ are shown in figure 2.54 and are clearly higher than the $PPR$.



**Figure 2.52** Cumulative peaks.

Two features of the cumulative peak count given in equation (2.132) stand out. First is the fact that there is no identifiable $MR(\text{-})$ term within it, which is not surprising considering the way in which $PCR$ is constructed from the profile; spacing information in the form of the $P$ and $H$ multipliers is preserved. The second feature is that the equation clearly splits into two components: the $A$ part in which $P$ occurs and the $B$ part in which the parameter $H$ occurs. In random theory $P$ is a measure of peak spacing and $H$ is a measure of zero crossing spacing. It can be shown that the $A$ term is a measure of the count of peaks above $\not{z}$. These peaks must have valleys also higher than $\not{z}$, and the $B$ term is a measure of those peaks above $\not{z}$ which have valleys below $\not{z}$. These two components of peaks can roughly be identified in Figure 2.53. The $H$ terms in $B$ in equation (2.132) can be thought of as the large peaks which directly deform and the $P$ terms in $A$ those peaks which indirectly deform. The situation is more complicated than this but the foregoing description is a possibility.



**Figure 2.53** Types of peak.

The two different types of peak relative to $\not{z}$ (i.e. the $P$s and $H$s) must have different functional properties. This suggests that the cumulative peak count $PCR$ should be split into two component curves (instead of just one as in figure 2.52) and that both of these curves should be used. A typical case is shown in figure 2.54.

**Figure 2.54** Breakdown of cumulative peak count.

A physical interpretation of these components is that for any given $z$ the *A* curve (in terms of *P*) represents peaks which are highly correlated with each other whereas the *B* curve represents a count of peaks which are independent [125].

*Discussion*

The section above indicates that obtaining peak parameters by using the peak 'count' distribution used extensively in communication theory may be somewhat misleading to tribologists. Using the peak count implies that every peak in the count is the same. This is not true for surfaces. It is only by resorting to the probability of peaks that some measure of the different peak properties is possible. In this case the peaks are effectively weighted by their radius of curvature. This reflects to what extent the surface surrounding the actual maxima lies within the boundaries set by $\delta z'$ and $\delta z''$ which are arbitrarily small and only need to be applied consistently throughout the waveform to be entirely valid.

It is shown how large differences between estimated parameters can result depending on whether peak 'count' probability or just the peak 'probability' is used. A displacement of $\pi/2$. $R_q$ downward when assessing the average peak height would probably be significant in contact situations such as mechanical seals for example.

Another observation is that the peak probability distribution could be more realistic than the material ratio curves used extensively in industry which completely ignore the influence of peaks when carrying load. $R_{pk}$, the highest parameter of the material ratio curve, is loosely described as a 'peak-like' parameter, which is a highly optimistic description. The cumulative peak probability curve could well be a better indication of the load carrying capability and wear properties of a surface than the material ratio curve or the cumulative peak count. However, the latter is useful for a different reason; it not only provides spatial information in the form of peak separation, it can also give an insight into spatial correlation between peaks. This aspect of possible contact behaviour is not reported in this paper. It is explained elsewhere [127].

In practice not all peaks make contact; it is far better to consider contact as a 'gap' property between two rough surfaces rather than a rough surface touching a smooth as in this paper. Nevertheless the argument remains valid.

As a next step some practical experiments need to be carried out to see if direct probability measures, as recommended for tribologists in this paper, can be carried out. Finally, areal measures rather than profiles should be investigated. However, it seems obvious that equal if not greater divergence will occur in parameters obtained by the two paths, that is peak probability versus peak frequency (count) of occurrence. Ideally, it should be possible to contrive a metrology scheme which uses the benefits of both methods.

### 2.1.4 *Areal texture parameters, isotropy and lay (continuous signal)*

Methods of measuring and specifying the surface over an area are now often investigated. Areal information from a surface falls into a number of categories: one is statistical average values, another is the structure and a third is defects or flaws. The last two can be found from purely qualitative information, that is a visualization of the surface. The former is concerned with quantifying the information in the $z$ direction primarily but increasingly in the $x$ and $y$ directions also.

The structural problem will become more and more important in non-engineering fields as will be seen, for instance, in biological applications. The form of areal investigation has of two types: the first has been to measure deterministic parameters in a large number of places and then specify a spread which the parameter might take, and the second has been to try somehow to predict from a few profiles or a single profile what the areal properties really are without measuring the whole surface. It is the areal characteristics which are important in function. The biggest problem is that if the area is completely measured the amount of data needed may be too large to handle. Because of the amount of data the most obvious way of attempting to specify areal characteristics is to use the statistical methods described above. This is precisely what has happened. Because the work done in these areas is as yet relatively small, emphasis is given to short-term work which will be likely to be of practical use. However, some consideration of new tools in the analysis is relevant which, although restricted to random surfaces, does enable some light to be thrown onto those parameters having tribological significance. Later, some practical cases of techniques, rather than analysis, will be discussed which may ultimately provide a more general basis for future work. In particular the Wigner distribution methods will be considered.

Ultimately, the worth of either theoretical or practical methods will be judged by their use in the field. This has yet to be demonstrated fully where random processes are concerned.

### 2.1.4.1 *Direct methods of statistical assessment over an area*

One way of doing this has been suggested by Peklenik and Kubo [35]. They have proposed a two-dimensional mapping of correlation functions split broadly into two fields; one where the surface is basically anisotropic and the other where it has only small variations from isotropic symmetry. Peklenik tackled the first and Kubo [36] the second. Fundamentally the distinction between the two approaches is the coordinate system used; Peklenik uses Cartesian coordinates (figure 2.55) and Kubo uses polar (figure 2.56). Neither of these methods has attracted much attention.

Other practical methods of determining areal properties will be dealt with in section 3.7 for the digital implications and in chapter 4 on instrumentation.

The Peklenik and Kubo methods are just techniques for extending the use of autocorrelation to areas, not for determining new parameters for three dimensions. This was achieved principally by Longuett–Higgins and later Nayak.

The advantage of the Peklenik–Kubo method is that it provides some tools for the description of surfaces which are not necessarily random in nature. Other techniques have been advocated which consider only random and similar surfaces.

Such a piece of work has been illustrated by Nayak [30], who transposed some of the work of Longuet–Higgins to manufactured surfaces rather than sea waves.

The importance of this original work was that it shows the errors which can result when extrapolating the peak and valley heights obtained from a profile to that of the whole surface. Longuet–Higgins quantifies the errors, pointing out that by using a simple profile the higher crests and troughs will sometimes be missed.

The analysis hinges again on the multivariate distribution and is similar to the Rice technique except that it is now in two independent dimensions rather than one. (This is defined as the areal case.) Once the

**Figure 2.55** Correlation length map (rectilinear).



**Figure 2.56** Correlation length map (polar).

areal case is considered the formulae become correspondingly more complicated. For example, the two-dimensional autocorrelation function becomes $A(\tau, i)$ where

$$A(\tau, i) = \lim_{L_1 L_2 \to \infty} \frac{1}{L_1 L_2} \int_{-L_1/2}^{L_1/2} \int_{-L_2/2}^{L_2/2} z(x, y), z(x + \tau, y + i) \mathrm{d}x \ \mathrm{d}y \tag{2.133}$$

$y$ being used as the other independent variable and $z$ the variable in the plane out of the surface. Similarly the spectrum becomes

$$P(u, v) = \frac{1}{4\pi^2} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} A(\tau, \ i) \exp[-j(\tau u + iv)] \mathrm{d}\tau \ \mathrm{d}i. \tag{2.134}$$

These formulae are used extensively in the optical work on Fourier analysis in chapter 4. Another useful formula is the moments of the spectrum, $m_{pq}$. Thus

$$m_{pq} = \int_{-\infty}^{\infty}\int P(u,v)u^p v^q \mathrm{d}u, \mathrm{d}v \qquad (2.135)$$

from which

$$\sigma^2 = m_{00}. \qquad (2.136)$$

It is possible to extract the information about a single profile from these formulae by noting that

$$P(\omega) = \int_{-\infty}^{\infty} P(u,v)\mathrm{d}l \qquad (2.137)$$

where

$$l = (u^2 + v^2 - \omega^2)^{1/2}. \qquad (2.138)$$

Some features of importance in tribology can be estimated analytically by using the MND and the parameters

$$z, \frac{\mathrm{d}z}{\mathrm{d}x}, \frac{\mathrm{d}z}{\mathrm{d}y}, \frac{\partial^2 z}{\partial x^2}, \frac{\partial^2 z}{\partial xy} \text{ and } \frac{\partial^2 z}{\partial y^2}. \qquad (2.139)$$

Using these the summit curvature, slopes, etc, can be found for the surface, at least on a microscopic scale. For example, the distribution of summit heights can be found by imposing the constraints that

$$\frac{\mathrm{d}z}{\mathrm{d}x} = 0, \quad \frac{\mathrm{d}z}{\mathrm{d}y} = 0, \quad \frac{\mathrm{d}^2 z}{\mathrm{d}x^2} < 0, \quad \frac{\mathrm{d}^2 z}{\mathrm{d}y^2} < 0 \quad \text{and} \quad \frac{\mathrm{d}^2 z}{\mathrm{d}x^2}\frac{\mathrm{d}^2 z}{\mathrm{d}y^2} - \left(\frac{\mathrm{d}^2 z}{\mathrm{d}x\mathrm{d}y}\right)^2 > 0. \qquad (2.140)$$

Then the probability of finding a summit at a height between $z$ and $z + \delta z$ in an area $\mathrm{d}A$ is

$$\delta z \iiiint f\left(z, \frac{\mathrm{d}z}{\mathrm{d}x}, \frac{\mathrm{d}z}{\mathrm{d}y}, \frac{\mathrm{d}^2 z}{\mathrm{d}x^2}, \frac{\mathrm{d}^2 z}{\mathrm{d}y^2}\right) \mathrm{d}\left(\frac{\mathrm{d}z}{\mathrm{d}x}\right), \left(\frac{\mathrm{d}z}{\mathrm{d}y}\right), \left(\frac{\mathrm{d}^2 z}{\mathrm{d}x^2}\right) \text{ etc} \qquad (2.141)$$

where the limits of integration have to obey the constraints given.

A key parameter that Nayak highlights is the ratio of some of the moments. In particular, for a Gaussian isotropic surface the parameter $\alpha$ is given by

$$\alpha = \frac{m_0 m_4}{(m_2)^2} \qquad \text{which in effect} = \left(\frac{\text{density of maxima}}{\text{density of positive crossing}}\right)^2 \qquad (2.142)$$

and is obviously closely related to $k$ used in the previous section 2.1.3.5. The parameter $\alpha$ (or $k$) can be useful in determining the reliability of parameters.

In equation (2.142)

$$m_0 = m_{00} = \sigma^2 \quad m_2 = m_{02} = m_{20} \qquad \text{the mean square slope}$$
$$m_4 = m_{40} = m_{04} = 3m_{22} \qquad \text{the mean square second differential.}$$

The moments with double subscripts represent the non-isotropic case.

It is known that all the relevant characteristics of stationary profiles can be expressed by means of the spectral moments. (Later it will be shown how these moments can be estimated by methods other than crossing and peak counts.)

For the areal surface the moments *mpq* are given by

$$m_{pq} = \iint \omega_i^p \omega_2^q p(\omega_1, \omega_2) \mathrm{d}\omega_1 \ \mathrm{d}\omega_2 \qquad (2.143)$$

and can be found by relating them to profile moments $m_r(\theta)$ at angle $\theta$:

$$m_r(\theta) = m_{r0} \cos^r \theta + \binom{r}{1} m_{r-1,1}, \cos^{r-1} \theta \sin \theta + \ldots + m_0, r \sin^r \theta \ldots \qquad (2.144)$$

using equation (2.143).

The second-order surface spectral moments $m_{20}$, $m_{11}$, and $m_{02}$, can be found from three second-order profile moments in arbitrary directions $\theta_1$, $\theta_2$, $\theta_3$,. Thus

$$\begin{pmatrix} m_{20} \\ m_{11} \\ m_{02} \end{pmatrix} = \mathsf{M}^{-1} \begin{pmatrix} m_2(\theta_1) \\ m_2(\theta_2) \\ m_2(\theta_3) \end{pmatrix} \qquad (2.145)$$

where

$$\mathsf{M} = \begin{pmatrix} \cos^2\theta_1 & 2\sin\theta_1 \cos\theta_1 \sin2\theta \\ \cos^2\theta_2 & 2\sin\theta_2 \cos\theta_2 \sin2\theta \\ \cos^2\theta_3 & 2\sin\theta_3 \cos\theta_3 \sin2\theta \end{pmatrix}. \qquad (2.146)$$

If $\theta_1 = 0°$, $\theta_2 = 45°$, $\theta_3 = 90°$, equation (2.145) reduces to

$$\begin{aligned} m_{20} &= m_2(0) \\ m_{11} &= m_2(45) - \tfrac{1}{2} m_2(0) - \tfrac{1}{2} m_2(90) \\ m_{02} &= m_2(90). \end{aligned}$$

Thus, if $m_{20}$ is the variance of the slope in the zero direction, $m_{02}$ is the variance of the slope in the 90° direction then $m_{11}$ is the covariance of the slopes in those two directions. It should be noted that the angles 0°, 45° and 90° were chosen for convenience only; however, equation (2.145) may be used with any three angles.

Again employing equation (2.144), the fourth-order surface spectral moments $m_{40}$, $m_{31}$, $m_{22}$, $m_{13}$ and $m_{04}$, can be calculated from five fourth-order profile spectral moments in arbitrary directions $\theta_1$, $\theta_2$, $\theta_3$, $\theta_4$ and $\theta_5$. These are

$$\begin{pmatrix} m_{40} \\ m_{31} \\ m_{22} \\ m_{13} \\ m_{04} \end{pmatrix} = \mathsf{M}^{-1} \begin{pmatrix} m_4(\theta_1) \\ m_4(\theta_2) \\ m_4(\theta_3) \\ m_4(\theta_4) \\ m_4(\theta_5) \end{pmatrix} \qquad (2.147)$$

where

$$M = \begin{pmatrix} \cos^4\theta_1 & 4\sin\theta_1\cos^3\theta_1 & 6\cos^2\theta_1\sin^3\theta_1 & 4\cos\theta_1\sin^3\theta_1 & \sin^4\theta_1 \\ \cos^4\theta_2 & 4\sin\theta_2\cos^3\theta_2 & 6\cos^2\theta_2\sin^3\theta_2 & 4\cos\theta_2\sin^3\theta_2 & \sin^4\theta_2 \\ \cos^4\theta_3 & 4\sin\theta_3\cos^3\theta_3 & 6\cos^2\theta_3\sin^3\theta_3 & 4\cos\theta_3\sin^3\theta_3 & \sin^4\theta_3 \\ \cos^4\theta_4 & 4\sin\theta_4\cos^3\theta_4 & 6\cos^2\theta_4\sin^3\theta_4 & 4\cos\theta_4\sin^3\theta_4 & \sin^4\theta_4 \\ \cos^4\theta_5 & 4\sin\theta_5\cos^3\theta_5 & 6\cos^2\theta_5\sin^3\theta_5 & 4\cos\theta_5\sin^3\theta_5 & \sin^4\theta_5 \end{pmatrix}.$$

The values of $m_{40}$ and $m_{04}$ are the variances of the curvature in two perpendicular directions, and $m_{22}$ is the covariance between these two curvatures. The values of $m_{31}$ and $m_{13}$ are the weighted covariances between the curvatures in two perpendicular directions defined as follows:

$$m_{31} = E\left[\left(\frac{d^2z}{dx^2}\right)^{3/2}\left(\frac{d^2z}{dy^2}\right)^{1/2}\right] \tag{2.148}$$

$$m_{13} = E\left[\left(\frac{d^2z}{dx^2}\right)^{1/2}\left(\frac{d^2z}{dy^2}\right)^{3/2}\right].$$

These expressions are considerably simplified when the surface is assumed to be isotropic, which implies that the stochastic properties of the surface in general and $m_r(\theta)$ in particular are independent of 6. Under this assumption, the surface and profile spectral moments are related by

$$\begin{aligned} m_{00} &= m_0 \\ m_{20} &= m_{02} = m_2 \\ m_{11} &= 0 \\ m_{40} &= m_{04} = 3m_{22} = m_4 \\ m_{31} &= m_{13}. \end{aligned} \tag{2.149}$$

Therefore, for this case the desired characteristics can be obtained from a single profile.

It can be shown, for instance, that using this nomenclature and with the associated constraints, that the density of summits is given by

$$D_{\text{sum}} = \frac{1}{6\pi\sqrt{3}}\left(\frac{m_4}{m_2}\right). \tag{2.150}$$

The parameter $\alpha$ is similar to the $Q$ factor of the power spectral density. Large $\alpha$ implies a wide spectrum. For $\alpha \to 1.5$ the spectrum becomes narrow.

Letting $z' = z/\sigma$ (normalized height)

$$\text{for } \alpha \to 1.5: \quad p_{\text{sum}}(z') = \begin{cases} \dfrac{2\sqrt{3}}{\sqrt{2\pi}}\left(-\dfrac{z'^2}{2}\right)[z'^2 - 1 + \exp(-z'^2)] & \text{for } z' \geqslant 0 \\ 0 & \text{for } z' < 0 \end{cases} \tag{2.151}$$

$$\text{for } \alpha \to \infty: \quad p_{\text{sum}}(z') = \frac{1}{\sqrt{2\pi}}\exp(-\tfrac{1}{2}z'^2)$$

which imply that for wideband surfaces the summit distribution is indeed Gaussian, whereas it is not for narrowband surfaces where it becomes more nearly Rayleigh or Maxwellian.

Similarly the expected value of the mean curvature for summits of height $z'$, expressed in a normalized form by dividing by $\sqrt{m_4}$, is

$$p(\text{sum. curv. at } z')\sqrt{m_4}$$

$$\frac{p(\text{sum. curv. at } z')}{\sqrt{m_4}} = \sqrt{2/3}z', z' \geq 0 \quad \text{for } \alpha \to 1.5$$

$$= (8/3).\sqrt{\pi} \quad \text{for } \alpha \to \infty$$

which indicates that for wideband surfaces the mean curvature varies linearly with peak height whereas for narrowband surfaces it is very nearly independent.

Estimates of $m_2$, $m_0$ and $m_4$ can be obtained from the number of maxima, minima and zeros on the profile (as the $k$ value of profiles). It has been shown by Longuet–Higgins that

$$\alpha = \frac{m_0 m_4}{m_2^2} = \left(\frac{\text{density of extremes}}{\text{density of zeros}}\right)^2 \tag{2.152}$$

thereby showing that $\alpha$ is in fact $k^2$ which is obtained from the profile. Hence the $\alpha$ value can be found by a simple counting routine.

Also, the density of peaks (or valleys) is as proposed earlier:

$$D_{\text{peak}} = \frac{1}{2\pi}\left(\frac{m_4}{m_2}\right)^{1/2} \tag{2.153}$$

which can be compared with the density of summits

$$D_{\text{sum}} \simeq 1.2(D_{\text{peak}})^2. \tag{2.154}$$

Note this is not the simple factor of $(D_{\text{peak}})^2$ as one might expect. Similarly the peak and summit heights differ. This is due to the presence of saddles and cols in the areal case.

The probability density $p_{\text{peak}}(z')$ is given by

$$p_{\text{peak}}(z') = \frac{k_1}{\sqrt{2\pi}}\left[\exp\left(-\frac{z'^2}{2k_1^2}\right) + \sqrt{\pi} \cdot k_2 \exp\left(-\frac{z'^2}{2}\right)(1 + \text{erf } k_2)\right]$$

where

$$k_1 = \left(\frac{\alpha - 1}{\alpha}\right)^{1/2} \quad \text{and} \quad k_2\left(\frac{1 - k_1^2}{2k_1^2}\right)^{1/2} z'. \tag{2.155}$$

It has also been shown that the probability density for a summit is

$$p_{\text{sum}}(z't) = \frac{\sqrt{3}k_3}{2\pi}\exp(-k_3 z'^2)[t^2 - 2 + 2\exp(-t^2/2)]\exp[-\tfrac{1}{2}(k_3 t^2 + k_2 t z')]$$

where $k_3 = \alpha/(2\alpha - 3)$ and $k_4 = k_3 (12/\alpha)^{1/2}$, and

$$t = \frac{1}{2}\left(\frac{3}{m_4}\right)^{1/2}\left(\frac{\mathrm{d}^2 z}{\mathrm{d}x^2} + \frac{\mathrm{d}^2 z}{\mathrm{d}y^2}\right) \tag{2.156}$$

from which $p_{\mathrm{sum}}$ at the limiting values $\alpha = 1.5$ can be found:

$$p_{peak}(z') = \begin{cases} z' \exp(-z'^2/2) & z' > 0 \\ 0 & z' > 0 \end{cases} \Big\} \alpha = 1.5$$

$$p_{peak}(z') = \frac{1}{\sqrt{2\pi}} \exp(-z'^2/2) \qquad \alpha \rightarrow \infty. \tag{2.157}$$

A comparison of these distributions reveals that the profile shows far fewer higher peaks and more low peaks than actually exist. The discrepancy is greatest for $\alpha \rightarrow 1.5$ and zero for $\alpha \rightarrow \infty$ (figure 2.57).



**Figure 2.57**  Peak and summit distributions.

The important point here is that tribologically important parameters can be found by means of the moments of the spectra. Whitehouse earlier approached this in a completely different way using the correlation function rather than the power spectrum and using *discrete* points from it rather than the continuous analysis of Nayak.

Similar effects are revealed for curvatures etc. Practical verification of these results has been obtained previously by Williamson [41] by making use of relocation profilometry. He shows divergences of the same order. See also Sayles and Thomas [42]. The conclusion reached from Longuet–Higgins' and Nayak's work is that the high-order statistics of the two-dimensional surface depends on the variances of the height, slope and second differential, $m_0$, $m_2$, $m_4$. which can be obtained from a single profile for Gaussian isotropic surfaces and from three non-parallel profiles for Gaussian non-isotropic surfaces—a technique checked by Wu *et al* [43].

It should be noted that there are fundamental assumptions made in much of the preceding work, which are that the height distribution is random (i.e. Gaussian in nature) and also that asperities and valleys follow reasonably well-behaved shapes.

Unfortunately neither has to be true for the surface roughness produced by many single-point and multiple-point cutting processes such as milling; the amplitude distribution is anything but random so that the results obtained above would hardly be applicable. Similarly the asperities are often not simple peaks; the very nature of cutting, especially with a blunt edge tool or grit, produces a wave-like effect which introduces overhang and re-entrant features into the geometry and gives multivalued $z$ values for $x$ and $y$ — not at all easy to deal with mathematically. (Incidentally, the latter factor is also true for sea waves, a factor which limits the usefulness of Longuet–Higgin's approach to a relatively calm sea.)

Other models have been used, some of which presume on the nature of the waviness present. For example, the work of Lukyanov [36] makes the assumption that the roughness is random and the waviness is not. Thus

$$z(x, y) = h(x, y) + w(x, y). \tag{2.158}$$

This approach will be considered under waviness.


### 2.1.4.2 Practical approach to areal (3D) measurement

Much of the previous work on areal (3D) assessment has been of a theoretical nature and has been based on the work of Longuet–Higgins [26] transposed from Oceanography by Nayak [30]. This work, coupled with the inability to map surfaces quickly and with high fidelity, put off practicing engineers. Both of these problem areas are being eroded by the more pragmatic approach of Stout *et al* [53], who, working with EC support have produced some plausible practical alternatives to the purely theoretical approach. It must be said that once the direct link with the theory has been broken the constraints imposed by rigour are lifted, allowing some useful 'functional' parameters to be put forward. The justification for these new parameters is powerful; they fit in better with practical engineering. A summary of these and similar parameters is given below.

One good idea in the development of a new parameter was to change the letter prefixing profile parameters from *R* to *S* indicating 'surface' parameters. This is not an international standard but it is a step forward in itself. There should be no ambiguity between the 'surface' and profile parameters. The abbreviations used below have not been adopted formally but, as they are straightforward, they are used here.

The primary parameter set is split into four groupings of parameters: amplitude parameters, spatial parameters, hybrid parameters and functional parameters.

Missing from the 17 parameters listed above is the symbol $S_a$ representing the surface average. It would probably be sensible to include it because it does represent the most tangible link with the part.

However, the original reasons for using $R_a$, namely that it could be checked from the profile chart whereas the RMS value could not, simply no longer applies because of the digital takeover. It certainly makes no sense to put both in—there are more than enough parameters anyway.

The amplitude parameters follow the profile definitions but the areal ones do not necessarily.

Some specific areal parameters result from the nature information, lay and isotropy can now be included. $S_{ae}$, the length of the shortest autocorrelation function, is one such parameter following on from the proposal of Kubo [36].

The hybrid parameters again follow on from the profile equivalent.

The functional parameters listed below use the definitions proposed by Stout [53]. They are basically an extension of the automotive parameters $R_{kv}$, $R_k$, $R_{ks}$ used for describing the shape of the material ratio curve. In the first instance these were obtained from the profile but in the 'surface' material ratio $S_{dr}$ they better describe what is taking place when the surface is being used. These surface parameters are provisionally listed in EUR 15178EN *Surface Bearing Index*. $S_{bi}$ is given by

$$S_{bi} = \frac{S_q}{z_{0.05}} = \frac{1}{h_{0.05}} \tag{2.159}$$

In equation (2.159) $z_{0.05}$ is the height of the surface at a height of 5% bearing (material) ratio.

A large value of $S_{bi}$ indicates good load carrying capability and would indicate a good bearing.

For a Gaussian surface $S_{bi}$ is about 0.6 and for a range of surfaces, including multiprocess finishes such as plateau honing, the value is between 0.3 and 2 which is a useful range. A run-in surface tends to have a larger value than an unworn surface so the $S_{bi}$ parameter could be used as a measure of wear.

**Table 2.8** Primary set of 3D surface roughness parameters.

| Amplitude parameters | |
|---|---|
| $S_q$ | Root-mean square deviation of the surface ($\mu$m) |
| $S_z$ | Ten point height of the surface ($\mu$m) |
| $S_{sk}$ | Skewness of the surface |
| $S_{ku}$ | Kurtosis of the surface |
| **Spatial parameters** | |
| $S_{ds}$ | Density of summits of the surface (mm$^{-2}$) |
| $S_{tr}$ | Texture aspect ratio of the surface |
| $S_{al}$ | Fastest decay autocorrelation length (mm) |
| $S_{td}$ | Texture direction of the surface (deg) |
| **Hybrid parameters** | |
| $S_\Delta$ | Root-mean square slope of the surface ($\mu$m/$\mu$m) |
| $S_{dr}$ | Arithmetic mean summit curvature ($\mu$m$^{-1}$) |
| $S_{sc}$ | Developed surface area ratio (%) |
| **Functional parameters characterizing bearing and oil retention properties** | |
| $S_{bi}$ | Surface bearing index |
| $S_{ci}$ | Core oil retention index |
| $S_{vi}$ | Valley oil retention index |
| $S_m$ | Material volume ($\mu$m$^3$/mm$^2$) |
| $S_c$ | Core valley volume ($\mu$m$^3$/mm$^2$) |
| $S_v$ | Deep valley volume ($\mu$m$^3$/mm$^2$) |

The primary parameter set is split into four groupings of parameters: amplitude parameters, spatial parameters, hybrid parameters and functional parameters.

*Surface material ratio $S_{bc}$*

$$S_{bc} = \frac{S_q}{z_{0.05}} \tag{2.160}$$

where $z_{0.05}$ is the height of the surface at 5% material ratio.

*Core fluid retention index $S_{ci}$*

$$S_{ci} = (v(0.05) - v_v(0.08))/S_q \text{ (unit area)} \tag{2.161}$$

where $v$ represents valley.

If $S_{ci}$ is large it indicates that the surface has good fluid retention properties.

The range is

$$0 < S_{ci} < 0,95 - (h_{0.05} - h_{0.08}) \tag{2.162}$$

when $h$ is $z$ normalized by $S_q$.

*Valley fluid retention index $S_{vi}$*

$$S_{vi} = (v_v(h = 0.8))/S_q \text{ (unit area)} \tag{2.163}$$

Here

$$0 < S_{vi} < 0.2 - (h_{0.8} - h_{0.05}) \tag{2.164}$$

For a Gaussian surface $S_{vi}$ is about 0.1 in value. The values 0.05, 0.8 etc are arbitrary but accepted as reasonable.

The use of functional parameters is not new. As early as 1936 Dr. Abbott proposed what is now the material ratio curve as a parameter useful for load carrying.

There is nothing wrong with functional parameters except for the fact that they are matched to specific applications and can give misleading values for unrelated functions. Mixtures of functional parameters could be used to describe different applications but then the parameter rash could well be the outcome.

Another problem is that the functional parameters mentioned above have been designed for multi-process finishes. However, the real gains to be made by measuring over an area are in predicting or modelling the path of the tool and not fine tuning the roughness. This is what is missing from today's list of parameters—all 17 of them according to table 2.8.

Using parameters such as $S_{ae}$ gives very little useful information about the 'lay' especially in multi-tool processes such as face milling. It could be argued that using a very general parameter to estimate a very specific lay pattern is questionable. The ironic part of this exercise is that the most important information (for example the tool path) in milling is completely ignored at the measurement stage only to be replaced by an almost irrelevant correlation procedure.

Undermining all the attempts to classify the areal properties of surfaces is the fact that the actual 'system', of which the surface is only a part, is neglected. The 'system' is comprised of at least two surfaces making contact at perhaps part of a bearing. The way these interact is of prime importance. One reason why the 'systems' approach is neglected is because the function of the system is difficult to achieve. Short of carrying out the function with real parts, which is impractical, the possibilities are to simulate the mechanics of the system in a computer given the measured data from both surfaces or to analyse the system theoretically. The latter is acceptable providing that the assumptions made are realistic.

The message is that it is pointless having 14 or 17 or 23 parameters for one or even both surfaces unless the interaction mechanism is itself characterized. At present, at best, the average distance between the mating surfaces is known and possibly any tilt. Usually even these factors are not known together with roundness error, cylindricity, etc. It seems inconceivable that most of the major systems parameters are neglected, yet more and more detail is sought for individual surfaces. The one redeeming factor is that Stout [130] has moved away from the theoretical parameters (found in communication theory and applied to surfaces) in favour of 'functional' parameters which should be at least partially relevant.

*Nature of Areal Definitions*

Description of areal properties corresponding to say a peak or valley and in particular dominant cases [45].

One of the first attempts to describe a surface in terms of hills and valleys was made by Maxwell [46]. According to him, a hill is an area from which maximum uphill paths lead to one particular peak and a dale is one in which maximum downhill paths lead to one particular dale. Also boundaries between hills are water courses and between dales watersheds; course line and ridge lines respectively. Both ridge and course lines emanate from saddle points and terminate at peaks and valleys.

Scott extends the ridge and course concepts to major hills and valleys [45], i.e. 'a hill is a single dominant peak surrounded by a ring of course lines connecting pits and saddle points, and a dale is a single dominant valley surrounded by a ring of ridge lines connecting peaks and saddle points'.

**Figure 2.58** Contour map showing critical lines and points: peaks *P*, pits *V*, saddles *S*.



**Figure 2.59** Full change tree.

This 'motif' method has some merit because functions such as mechanical seal and flow in general depend for their efficiency on a knowledge of the probability of escape routes. Pathways between peaks and valleys in the areal model can provide one way to estimate the flow of liquid within the zone. In the case where contact has been made, as in a seal, it may be a relevant parameter.

The pathway from summit to dale is shown in figure 2.58. This is called a 'change tree'. Clearly the tree can be made for hills and vales and for the typical joint situation.

The critical paths drawn on the contours of a surface together with the hills and vales trees give an idea of the way that flow of liquid or air passes across the surface. Rather than just saying that there must be leakage between two surfaces in, say, a mechanical seal it could be possible using this methodology to determine the actual pathway.



**Figure 2.60** Dale change tree.



**Figure 2.61** Hill change tree.

### 2.1.5    *Discrete characterization*

#### 2.1.5.1    *General*

Whitehouse based a typology on what was measured, that is the actual ordinates of the surface itself. Thus instead of using joint probability distributions $p(z, m, c)$ involving surface height $z$, slope $m$ and curvature $c$, all taken as random variables, he used the joint probability distributions $p(z_1, z_2, z_3, z_4., ... )$ of measurements of the surface. These were used to examine the average joint behaviour of the measurements using the multinormal distribution in such a way as to predict functionally important properties or manufacturing process properties.

This completely different way of using the multinormal distribution was tackled by considering the average behaviour of discrete ordinates of the surface rather than the continuous variables used by Nayak. This involved discrete values of the autocorrelation function rather than moments of the continuous power spectrum (figure 2.62(*a*)).

**Figure 2.62** (*a*) Three-point autocorrelation function. (*b*) Areal correlation coefficient positions.

In a series of papers from 1969, he and co-workers [29, 47–49] developed a characterization of profiles and areal (or 3D) surfaces totally in terms of discrete values of the autocorrelation function of the profile and areal coverage of the surface (figure 2.62(b)).

A summary of the way in which this typology of surfaces can be used will be given here. The theory is demonstrated in chapter 3, where exactly the same procedure can be used to explain digital measurement results. This is the innovation. It formally relates the characterization of the surface to the measurement possibilities. One important note is that, whereas the results for the discrete autocorrelation function of a profile converge to that of Nayak's profile as the spacing (*h*) between correlation points is taken to zero, the discrete areal results do not converge to the continuous areal results.

The results for profiles by Whitehouse and Nayak have been compared experimentally and theoretically by Bhushant [50]. Also a characterization by Greenwood [30] will be given. This latter treatment is an attempt to simplify and clarify the Whitehouse results.

All the methods of Nayak, Whitehouse and Greenwood rely on the multinormal distribution as a basis for the calculations. It has been found that a certain amount of non-Gaussian behaviour can be tolerated. It has been suggested by Staufert [51] that a skew value of up to $\pm 1$ still allows good results.

The Whitehouse parameters for a profile are given in terms of the profile correlation function at the origin $R_q^2$ and at positions $h$ and $2h$ from the origin, that is $A(h)$ and $A(2h)$ called here $\rho_1$ and $\rho_2$. Thus the mean peak height is

$$R_q \left( \frac{1-\rho_1}{\pi} \right)^{1/2} \Big/ 2 \left( 1/\pi \tan^{-1} \left( \frac{(3-4\rho_1+\rho_2)}{(1-\rho_2)} \right)^{1/2} \right) \tag{2.165}$$

the density of peaks is

$$D_p = \frac{1}{\pi h} \tan^{-1} \left( \frac{(3-4\rho_1+\rho_2)}{(1-\rho_2)} \right)^{1/2} = \frac{1}{\pi h} \tan^{-1} \left( \frac{A_3}{A_2} \right)^{1/2} \tag{2.166}$$

the density of zero crossings is

$$n_0 = \frac{1}{\pi h} \cos^{-1} \rho_1 \tag{2.167}$$

the mean slope is

$$m = \frac{R_q}{h}\left(\frac{1-\rho_2}{\pi}\right)^{1/2} = \frac{R_q}{h}\left(\frac{A_2}{\pi}\right)^{1/2} \tag{2.168}$$

the variance of peak heights is

$$1 + (R_q^2 \{[(A_1/2\pi D_p)(A_2/A_3)]^{1/2} - (A_1/4\pi D_p^2)\}) \tag{2.169}$$

the mean curvature of peaks is

$$(R_q^2 A_3)/[2 D_p h^2 (\pi A_1)^{1/2} \tag{2.170}$$

the variance of curvatures is

$$\sigma_c^2 = [(R_q^2 A_3)/[4 h^4 \pi A_1)](8\pi A_1) + (2h/D_p)(A_3 A_2)^{1/2} - (A_3 h^2/D_p^2)] \tag{2.171}$$

and the correlation between curvature and peak height is

$$\text{correlation} = [1 - (1 - 2\rho_1^2 + \rho_2)/(A_3 \sigma_p^2)]^{1/2} \tag{2.172}$$

where $A_1 = 1 - \rho_1$, $A_2 = (1 - \rho_2)$, $A_3 = (3 - 4\rho_1 + \rho_2)$ and $\sigma_p^2$ is the peak variance normalized with respect to $\sigma^2$, i.e., $R_q^2$.

This set of parameters derived from three points on the correlation can always be augmented by using more than three points, say five or seven.

As it stands this model permits a best spectral fit of the form P($\omega$):

$$P(\omega) \sim 1/[(j\omega)^2 + j\omega(\rho_1^2 - \rho_2)/(1 - \rho_1^2) + \rho_1(\rho_2 - 1)/(1 - \rho_1^2)]^2 \tag{2.173}$$

of which the exponential is a special case. Fitting four or more points to the autocorrelation function allows second-order spectra to be modelled and higher degrees of freedom. This has been found to be largely unnecessary.

In extending this to areal (or 3D) surfaces a number of options arise depending on how the area is sampled. The comparison of the behaviour of surface properties between continuous and discrete models is not straightforward. This is because of the options available in the discrete definition. Here, the straightforward tetragonal pattern of sampling will be used in which a summit is defined if it is larger than four orthogonal adjacent points as in figure 3.21(*c*). Digonal, trigonal and hexagonal patterns will be considered in addition to this in chapter 3 on digital processing.

The problem is that none of the patterns is wrong or right; they are just different. This difference matters because it is not just what is characterized here as a discrete model, it is actually what everybody measures. This discrete modelling provides the link between the theory and the actual experiments.

As an example of how the continuous method differs from the discrete, consider the density of summits. For the continuous case given by Nayak

$$D_{\text{sum}} = (6\pi\sqrt{3})^{-1}(D_4/-D_2) \tag{2.174}$$

and for the discrete tetragonal model for the limiting case when $h \to 0$

$$D_{ds} = \frac{\pi + 2\sin^{-1}(1/3) + 4\sqrt{2}}{4\pi^2}\left(\frac{D_4}{D_2}\right)h^2$$
$$= \sqrt{3}\,\frac{\pi + 2\sin^{-1}(1/3) + 4\sqrt{2}}{4\pi}\,D_{sum}$$
$$= 1.306\,D_{sum}. \tag{2.175}$$

Hence the discrete summit density is 30% larger than the continuous. This is because the discrete model identifies what it thinks are summits but which turn out not to be! This is easily possible because any discrete model cannot include all the possibilities of exclusion from the summit definition in the analogue case. Taking this further, the argument would maintain that the fewer the number of points making up the definition, the higher the density would appear to be. Indeed, for trigonal models in which three adjacent points to a central one are considered, the density is 73% higher. In equation (2.175) $D_2$ and $D_4$ are the second and fourth differentials of the autocorrelation function at the origin, that is $A''(0)$ and $A^{iv}(0)$: and

$$\eta = -D_2(D_4)^{-1/2} \simeq -\frac{A''(0)}{\sqrt{A^{iv}(0)}}$$

which is the correlation between the height and curvature for a random waveform.

For summit height the continuous summit mean height is

$$E[z\,|continuous] = \frac{4}{\sqrt{\pi}}\eta\ . \tag{2.176}$$

For tetragonal characterization, which is the standard pattern used by most researchers, $E[z|\text{discrete}]$ is given by

$$E(z\,|discrete] = \frac{8\sqrt{2\pi}}{\pi + 2\sin^{-1}(1/3) + 4\sqrt{2}}\eta$$
$$= (0.938)\left(\frac{4}{\sqrt{\pi}}\eta\right) = 0.938 \times E[z|\text{continuous}]. \tag{2.177}$$

The mean summit height is lower by 6%.

These disturbing differences are taken further in the next chapter but the point here is that, even if it is agreed that correlation (or spectral) methods form the basis of a good model for the characterization of the surface, much has to be agreed between researchers either before they will agree with each other's results on the same surface or on whether they ever get to the 'true' value of the surface even if it were possible. The problem has been masked by dealing mainly with profiles. Discrete (measured) values of the profile will gratifyingly converge to the theoretical case as the sampling gets higher. But in the areal case they never do! The difference is dependent on the coverage of the discrete model.

### 2.1.5.2 Alternative discrete methods

In an effort to clarify the situation Greenwood [49] reworked the discrete results of Whitehouse and changed the nomenclature to simplify the results. For example, the probability of an ordinate being a peak in Greenwood's theory becomes

$$N = \frac{\theta}{\pi}\ \text{where}\ \theta = \sin^{-1}(h\sigma_k/2\sigma_m) \tag{2.178}$$

instead of

$$N = \frac{1}{\pi} \tan^{-1} \left( \frac{(3 - 4\rho_1 + \rho_2)^{1/2}}{(1 - \rho_2)} \right)$$

where $\sigma_k$ and $\sigma_m$ are the standard deviations of curvature and slope respectively; $\theta$ becomes one of the characterizing parameters (actually dependent on $h$).

Similar expressions corresponding to equations (2.165)–(2.171) are found. Thus the mean peak curvature (expressed as a ratio to $\sigma_k$) is

$$\sqrt{\tfrac{1}{2}\pi} \, \sin\theta / \theta \tag{2.179}$$

and the variance of peak curvature (again expressed to $\sigma_k^2$, the variance of curvature) is

$$1 + (\sin 2\theta / 2\theta) - \frac{\pi}{2} (\sin \theta / \theta)^2. \tag{2.180}$$

The other parameter Greenwood uses to characterize the surface he calls $r$, where $r$ is given by

$$r = (\sigma_m^2 / \sigma) \sigma_k \tag{2.181}$$

where

$$r^{-2} = \alpha = \frac{m_0 m_4}{m_2^2} \,. \tag{2.182}$$

This $r$ represents the 'surface roughness character'. Between the $r$ and $\theta$, Greenwood goes on to establish summit properties using a third parameter

$$\tau = 2(1 - 2\rho_1 + \rho_{\sqrt{2}})/(3 - 4\rho_1 + \rho_2). \tag{2.183}$$

For a good comparison of the Nayak, Whitehouse and Greenwood models the reader should consult the paper by Greenwood [30]. The two methods of Whitehouse and Greenwood are merely different ways of tackling the discrete problem.

The Nayak model for the surface is continuous, the Whitehouse model is discrete and the Greenwood model attempts to provide a link between the two. It should be remembered, however, that although Greenwood's characterization provides simpler formulae, the parameters $\sigma_k$ and $\sigma_m$ are still determined by discrete methods which use the Lagrangian numerical analysis values. Whitehouse accepts from the start the overall dependence on discrete values and establishes the method of characterization completely on it; the characterization formulae and the numerical analysis formulae are made to coincide.

In these methods of characterizing surfaces there is almost universal acceptance of the fact that the ratio of zero crossings (or mean line intersections) and the peak density is a useful surface characterization parameter involving, in its various guises, $\alpha$, $r$ and $k$.

*The important result is that there are no intrinsic parameters for the surface*. No surface has a mean summit height of $\overline{z}$; it only has a mean summit height $\overline{z}$ for a given numerical model $f(z)$ and sample interval $h$. Providing that $f(z)$ and $h$ accompany $\overline{z}$ when results are compared (and assuming that there are no instrumental problems) this is acceptable.

To conclude this section, it has been shown that in any practical situation it is necessary to characterize the surface and the procedure if results are to be compared. Another point which this section on surface

roughness shows is the enormous difference in complexity which the manufacturing engineer requires to specify the surface as compared with that of the tribologist. The former requirement is straightforward—simple averages like $R_a$ can be used—whereas in the latter much more complicated parameters are required.

Finally, these methods of specifying the characteristics of surfaces have resulted in methods involving average values, which allow the probability of extreme values such as peaks to be estimated. This does not solve the problem of specifying freak behaviour (even if the freak can be found). Neither does it get much closer to solving the question of what surfaces are really like. At best the properties within very closely defined limits can be specified, although even with fractal-type surfaces (discussed later in section 2.1.8 and chapter 7) this is difficult.

### 2.1.6    Assessment of isotropy and lay

The characterization of the lay of surfaces has been sadly neglected, mainly because of its complexity. An elementary classification is shown in figure 2.64.

Isotropy is usually taken to mean the lay of random-type surfaces, thereby differentiating between grinding on the one hand and, say, shotblasting on the other. General lay takes all surfaces into account. It cannot, however, be ignored in very many functional situations, especially relating to bearings, mechanical seals and cylinder walls; it can be more important than the value of the roughness. The reason, as Wirtz [52] pointed out, is due to the large variety of motions of the machine tool in generating the surface. Simply using the lay as a means of typology or characterization is not enough.

One method suggested by Peklenik and Kubo depends on characterizing the plots of the variations in the correlation functions on a linear or polar scale. This method is convincing but difficult to do because of the need to map the surface in some detail. An alternative method is to use the 'long crestedness' criterion of Longuet–Higgins. This type of surface approximates to a series of parallel grooves. The definition that relates to that of Peklenik and Kubo requires taking profile graphs in different directions and estimating some feature of the plot (figure 2.56). In Longuet–Higgins' case the following long crestedness ($LC$) criterion has been defined [26] as

$$LC = \frac{m_{20} + m_{02} + [(m_{20} - m_{02})^2 + 4m_{11}^2]^{1/2}}{m_{20} + m_{02} - [(m_{20} + m_{02})^2 - 4m_{11}^2]^{1/2}}. \tag{2.184}$$

For a completely isotropic surface $LC = 1$; for a completely long crested wave $(m_2)_{min} = 0$ and $LC \to \infty$. If $(m_2)_{min}$ is made equal to 0, $m_{11}$ is equal to $(m_{02} m_{20})^{1/2}$, remembering that the $m_{02}$ and $m_{20}$ moments would have been taken at right angles.

An excellent account of this type of approach is given by Pandit *et al* [60]. This will be referred to in section 2.1.7.6. Although this method of quantifying isotropy is sensible for truly random surfaces, it runs into severe problems for other kinds of surface. Take for example a face-turned part. Such a criterion would produce an isotropic result only if the measured tracks started at the centre of the face; anywhere else as a starting point would give a different result. From the measurement point of view the surface is neither random nor stationary. One could argue that this technique should only be used for random surfaces. The problem here, however, is that it is never possible to be sure whether deterministic elements are present. This criticism applies equally to the correlation plotting methods proposed by Peklenik and Kubo.

A technique for describing isotropy was also tried by Ledocq [54] who made a succession of traces in radial directions on various machined surfaces (similar to the work of Peklenik and Kubo). He then plotted the variations of well-known parameters such as $R_a$ and $R_q$, the average wavelength and the 0.5 correlation length. Unfortunately, although results were obtained they proved to be inconclusive, not least because of the instrumental difficulty of measuring texture along the lay—the problem is that of knowing what, if anything, to do with the sampling length for different directions!

What is needed is a technique which deals with any surface and preferably gives a numerical value to the isotropy and lay.

One sensible approach has been proposed by Wirtz [51] who suggested that it should be possible to relate the lay to a number of basic geometrical features such as circles, lines, etc (figure 2.63). In principle, this technique would derive from the appearance of the surface the movements of the tool which generated it. This after all is the lay. There is, however, still the problem of quantifying it. It seems that this type of method as it stands has potential but not yet the capability because of the difficulty of measurement. Probably the best method up to now is to use a two-dimensional spectrum method of a scaled-down model of the surface using 'liquid gates'. This enables an areal diffraction pattern to be produced which reveals the lay. Perhaps the only method is that used at present following drawing practice.



Radius fixed,                    Generator fixed,
moving generator                 moving radius

**Figure 2.63**    Typical lay patterns of machined parts produced with single tool.

In this method the pattern is identified by reference to standard patterns in a way which is similar to the use of metallurgical atlases. In these, photographs of metals under different conditions are compared by eye. There is, however, no reason why the comparison cannot now now be done by a pattern recognition method by computer. There are many techniques for recognizing patterns. Many such patterns have been presented in early work in the French standards for surfaces and in British standards (figure 2.64)

A much more comprehensive approach has been suggested along these lines [63] which characterizes the lay along set theory lines into patterns comprising points, straight lines, circles (including secants), polygons and curves. Up to 17 different groupings have been observed, each group having two subset groupings. In all, hundreds of different observed patterns can be classified in this way. However, the complexity of the system probably makes the application of this method prohibitive. This technique, as in the Wirtz methods, is basically a pattern recognition technique and so, in principle, with the massive advances in this subject since the appearance of vision systems for robots, could be tried again.

Also, since the work on fractals is essentially a description of form and structure [56], it is likely that there is scope for a proper visual classification system.

The work on areal assessment is basically on two levels: one is visualization and the other is quantification. The point to be made here is that the use of simple visualization is not to be ignored because, as pointed out by Stout [53], a great deal of information, particularly about structure, is available from a surface map. Use of suitable software can enable the eye to see many different aspects of the surface which would normally be lost.

Stout lists a few suitable operations which can be carried out on a 3D map to help the operator [53]. These are:

(1)  inversion by turning the picture upside down;
(2)  magnification in any direction;

| Symbol | Interpretation | |
|--------|----------------|---|
| = | Parallel to the plane of projection of the view in which the symbol is used |  Direction of lay |
| ⊥ | Perpendicular to the plane of projection of the view in which the symbol is used |  Direction of lay |
| X | Crossed in two slant directions relative to the plane of projection of the view in which the symbol is used |  Direction of lay |
| M | Multi-directional |  |
| C | Approximately circular relative to the centre of the surface to which the symbol is applied |  |
| R | Approximately radial relative to the centre of the surface to which the symbol is applied |  |

**Figure 2.64** A classification of lay.

(3) truncation of the signal at any level of height;
(4) contouring.

One of the very few pieces of work which directly addresses lay measurement of non-random surfaces is by Boudrean and Raja [117]. The technique is quite simple: it makes use of the fact that two closely spaced parallel profiles will appear to be shifted whenever the dominant lay direction is not perpendicular to the profile direction. Boudrean and Raja utilize the cross-correlation function to measure these shifts and then simple linear models are applied to the shifts to quantify the lay characteristics. This contribution was interesting in the sense that it could cater for lay which was curved as well as straight. The only problem in the case of curved lay is that the shift between profiles depends on the position of the traces within the data map (figure 2.65).

From the information obtained by this method useful data concerning the manufacturing process could be obtained after considerable interpretation. This work certainly reinforced the view that all data on the surface of a workpiece is useful for something.



**Figure 2.65** Raja lay topology.

So far this chapter has virtually followed a chronological path. This development has been intentional because it keeps firmly in mind the problems that arose in the early days and the methods used to solve them. Although the techniques are much more powerful today new problems keep emerging.

The interesting point is that the old problems keep re-emerging as new processes are introduced or new functions demanded.

Many of the newer methods such as discrete analysis and random process analysis are now incorporated into instruments. However, there is a continual search for newer, perhaps better, parameters. In what follows a few of the alternative parameters will be revealed. Although these have been demonstrated in various research papers they cannot be said to have gained universal acceptance. This may or may not come. Nevertheless, it is useful to mention them because they point the way to possible developments in the future.

### 2.1.7 Potential methods of characterization

#### 2.1.7.1 Amplitude and hybrid parameters

The route taken so far—that of considering the profile parameters which are in use today, progressing to random process analysis and from this to the areal evaluation of surfaces proposed by Longuett–Higgins, Nayak and Whitehouse—is taken as the general direction of progress in surface characterization. There have, however, been many other attempts to classify the surface roughness which, although not major advances, can still be considered to be interesting and even valuable in some applications. Some of these are briefly considered to show the general way in which the subject is developing.

Many simple classification systems for a typology of surfaces have been based either on the amplitude characteristics alone or on the spacings of profiles. A notable exception in the early days was Myres [14] who used four parameters, $R_q$, $\Delta_q$, RMS curvature and the directionality described in equation (2.11). His interest was mainly in the influence of surfaces on friction. Incidentally, he found that over a wide range of roughnesses $\Delta_q$ had a correlation of 0.85 with $\mu$ the frictional coefficient—a result which would be received with some scepticism today.

There have been a number of attempts to classify the amplitude probability density function [8]. The first were probably Pesante [9] and Ehrenreich [10] who used the slope of the bearing (material) ratio curve even earlier but did not point out explicitly the statistical significance. Other people tried various ratios of the $R_a$ and $R_q$ values taken from the material ratio curve, and Oonishi appears to have been the first to consider the combination of different processes on the same material ratio curve [48]. He also tried to incorporate some ideas about including a peak count in addition to the height information, as did Reason [7].

These attempts showed two things: one was that the existing parameters were insufficient to characterize the surface; the other was that it was not very clear which way to go forward!

#### 2.1.7.2 Skew and kurtosis

More recently a more statistical method was suggested by Al-Salihi [24] who included the skew and kurtosis values together with the $R_a$ or $R_q$ values as a complete set of characters. The two extra moments were an attempt to quantify the shape of the amplitude density curve. For many surfaces this technique was an acceptable step forward. However, problems can occur because the two parameters are not independent. An alternative method was proposed by Whitehouse [55]. His idea was to try to devise parameters which weighted peaks differently to valleys, thereby introducing a variable which could be useful in predicting function. The proposal was to use the beta function as a means of characterizing the amplitude density curve.

### 2.1.7.3 Beta function [55]

The beta function is one that is defined within a given range of $0(r)1$. It is expressed in terms of two arguments, $a$ and $b$. Thus $\beta(a, b)$ is given by

$$\beta(a,b) = \int_0^1 z^{a-1}(1-z)^{b-1}\,\mathrm{d}z. \tag{2.185}$$

It can be expressed in terms of a probability density $p_\beta(a, b, z)$:

$$p_\beta(a,b,z) = \frac{1}{\beta(a,b)} z^{a-1}(1-z)^{b-} \tag{2.186}$$

where $\beta(a, b)$ is the normalizing factor to make $p_\beta$ a true density. The usual parameters of the distribution can be determined in terms of $a$ and $b$. Thus, using gamma function identities,

$$\beta(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad \text{and} \quad \Gamma(a+1) = a\Gamma(a) \tag{2.187}$$

$$\bar{z} = \frac{1}{\beta(a,b)} \int_0^1 z z^{a-1}(1-y)^{b-1}\,\mathrm{d}z = \frac{\beta(a+1,b)}{\beta(a,b)} = \frac{a}{a+b} \tag{2.188}$$

the variance is

$$\sigma^2 = \frac{ab}{(a+b+1)(a+b)^2} \tag{2.189}$$

the skew is

$$S_{K\beta} = \frac{2(b-a)}{a+b-2}\left(\frac{a+b+1}{ab}\right)^{1/} \tag{2.190}$$

and the kurtosis is

$$K_\beta = \frac{6[(a-b)^2(a+b-1) - ab(a+b+2)}{ab(a+b+3)(a+b+2)}. \tag{2.191}$$

The basic philosophy was that any practical amplitude distribution could be approximated by a beta function (figure 2.66). Note that the function has two parameters, each independent of each other, and so could be used as a basis for a typology or characterization; $a$ is the weighting allocated to ordinates of the profile measured from the lowest valley upwards and $b$ is the weighting of the same profile as seen from the top of the profile and measured downwards. Peaks and valleys are therefore differently weighted. This could (a point evaluated later) provide useful functional discrimination. The problem arises of how to determine $a$ and $b$ from the profile. By changing the range in equation (2.185) from 0 to 1 to $R_p + R_v$ or $R_t$ and replacing $\sigma$, the standard deviation of the distribution, by $R_q$, the $a$ and $b$ parameters become

$$a = \frac{R_v(R_v R_p - R_q^2)}{R_t R_q^2} \quad b = \frac{R_p(R_v R_p - R_q^2)}{R_t R_q^2}. \tag{2.192}$$

**Figure 2.66** The beta function (left) symmetrical and (right) asymmetrical case.

The fact that odd peaks or valleys are only raised to a unit power suggests extra stability over the skew and kurtosis characterization method. Nevertheless, the weakness of the technique is that good estimates of $R_v$ and $R_p$ have to be made, which is difficult, indicating once again the basic problem of measuring peaks and anything derived from them! As an example of how the technique works from manufactured surfaces, compare this typology with one based on the skew and kurtosis. This is shown in figures 2.67($a$) and 2.67($b$).

The other problem associated with the beta function approach is that it cannot easily cope with multimode distributions as can be seen by reference to figure 2.66. So, for simple periodic profiles, there is sometimes a problem.



**Figure 2.67** Process identification: ($a$) central moments, ($b$) beta function.

It has correctly been pointed out by Spedding [56, 62] that the beta function is only one example of the class of Pearson distributions and that delving further into these classes points to a better and still more comprehensive classification system. It may also be said that other descriptions of distributions exist which utilize three parameters and so could, in principle, be used to improve this type of characterization. One such is the hypergeometric function. The problem then arises of having too many parameters. If three are needed for the shape of the distribution and one for the size or scale factor then already the situation exists where four parameters are needed to specify amplitude information alone. This would be difficult to implement or justify practically. Alternative methods of characterization of the amplitude density function have been tried by using different types of function other than the beta function. Examples of this have been the Chebyshev polynomial approximation and the Fourier characteristic function

An obvious choice is the characteristic function of the distribution. This is the Fourier transform of the amplitude distribution function where $\omega$ is usually replaced by $\zeta$. Thus

$$C(\zeta) = \int_{-\infty}^{\infty} \exp(j\zeta z) p(z) \mathrm{d}z. \tag{2.193}$$

Because this is a continuous function an arbitrary choice of $\zeta$ would be needed upon which to build a typology. However, a method is possible based upon the first few coefficients of the Fourier series of $p(z)$ using the range $R = R_t$ or $R_z$ as the fundamental wavelength.

Here a typology could be based on $F(n)$, $n = 0, 1$, where $F(n)$ is given by

$$F(n) = \int_0^R \exp(jn2\pi z/R_p)(z)\mathrm{d}z \tag{2.194}$$

which means matching to the distribution an expression of the form $f(A, B, C)$ where

$$f(A, B, C) = A + B\cos(2\pi z/R) + C\sin(2\pi z/R). \tag{2.195}$$

By a similar technique to that used for the beta function $A$, $B$ and $C$ can be found in terms of $R$, $R_q$ and $R_v$ or $R_p$. Thus

$$
\begin{aligned}
A &= 1/R \\
B &= \frac{2\pi^2}{R^3}\left(\frac{R^2}{6} + R_v^2 + R_q^2 - RR_v\right) \\
C &= \frac{\pi}{R}\left(1 - 2\frac{R_v}{R}\right).
\end{aligned}
\tag{2.196}
$$

At first sight this appears to have some advantages. One of these certainly is that there is a strong connection between this and power spectral analysis, which is often used for random process analysis and currently used in surface metrology analysis. Unfortunately, despite having three parameters available (two, if the first is regarded as trivial), the possible shapes that can be matched are much more limited than even the beta function. Also, the simple, deterministic types of profile have awkward values. For example, the sine wave has values of $C = 0$ and $B = 0.822$. Symmetry considerations of the distribution especially are not well dealt with by using the odd symmetry of the sinusoidal component. Another disadvantage is that the shape of the fitted distribution does not change as the coefficients $B$ and $C$ are changed in magnitude. This is not true for the beta function: as $a$ and $b$ get large the distribution becomes more 'spiky', a property which allows a little extra flexibility. Obviously the Fourier transform approach would be more suitable if more coefficients

were used, but this is not allowable because of the practical difficulty of adopting any system which employs more than one or two parameters.

### 2.1.7.4 Chebychev function and log normal function

Another possible type of function to use is the Chebychev function. This is not an arbitrary choice of function because serious attempts to analyse the profile itself into Chebychev rather than power spectral coefficients have been made recently. The method would involve identifying a curve $p(T_0, T_1, T_2)$ with the amplitude distribution of the profile $p(z)$, where $T_0$, $T_1$ and $T_2$ are Chebychev polynomial coefficients of the first kind and where $p(T_0, T_1, T_2)$ is given by

$$p(T_0, T_1, T_2) = T_0 + T_1 z + T_2 (2z^2 - 1) \qquad (2.197)$$

which for a range corresponding with $\pm R/2$ gives the coefficients in terms of surface profile parameters as

$$
\begin{aligned}
T_0 &= \frac{13}{46} + \frac{60}{46}\left[\left(\frac{R_q}{2R}\right)^2 + \left(1 - \frac{R_p}{2R}\right)^2\right] \\
T_1 &= \frac{3}{2}\left(1 - \frac{R_p}{2R}\right) \\
T_2 &= \frac{45}{46}\left[\left(\frac{R_q}{2R}\right)^2 + \left(1 - \frac{R_p}{2R}\right)^2 - \frac{1}{3}\right].
\end{aligned}
\qquad (2.198)
$$

An interesting alternative to the beta function which appears certainly to be as good is due to Murthy *et al* [57] who proposed a log normal distribution. Thus, if the profile ordinates are $z(x)$, then $v$ is $\ln(z)$. If the transformed variable $v$ is distributed according to a Gaussian (normal) model then $z$ has a log normal distribution which is often found in distributions of extreme values such as peak distributions.

This distribution is

$$p(z, \mu, \sigma) = \frac{1}{z\sigma\sqrt{2\pi}}\exp\left[-\frac{1}{2}\left(\frac{z - \mu}{\sigma}\right)^2\right]. \qquad (2.199)$$

Classical theory gives

$$
\begin{aligned}
\bar{z} &= \exp(\mu + \sigma^2/2) \\
\bar{\sigma}^2 &= \exp(2\mu + 2\sigma^2) - \bar{z}^2.
\end{aligned}
\qquad (2.200)
$$

When the expression is referred to surface characteristics $R_q, R_v, R_p$ as for the beta function, the arguments $\mu$ and $\bar{\sigma}$ become

$$
\mu = \tfrac{1}{2}\ln\left(\frac{(2R_p)^4}{(2R_p)^2 + (R_q)^2}\right)
$$

$$
\bar{\sigma} = \ln\left(\frac{R_q^2 + (2R_p)^2}{(2R_p)^2}\right).
\qquad (2.201)
$$

Plotting $\mu$ against $\sigma$ for a number of processes gives much the same amount of discrimination as the beta function. Unfortunately omitted from the investigation is the Gaussian curve which is discussed in more detail in the section on multiparameters. Murthy *et al* also consider a two-dimensional mapping of the material ratio curve [57]. Some current thinking, not entirely new, concludes that three basic parameters are needed for the height, one for scale and two for shape. It can, of course, be arranged that three parameters all have the dimensions of length and the ratios provide shape information. One such method separates the bearing curve into three parts: the peaks, the basic surface and the valleys (figure 2.13). This method, not yet standardized, seems to give good results for multiprocesses such as plateau honing but the construction necessary to identify the three regions of the curve is not suitable for periodic surfaces. It should be pointed out that this type of surface, which is split according to height regions as above, is often called a 'stratified' surface. This word is used to indicate that each stratum of the surface is targeted at a given function, that is load capacity, oil retention, etc.

### 2.1.7.5 Variations on material ratio curve

One method for characterizing stratified surfaces uses the material ratio curve. This note is purely practical and should not be used as an alternative to the cumulative peak variants mentioned earlier. In particular this method is used in plateau honing in the car industry. As is well known there are two important aspect of the profile which have to be independently assessed: one is the plateau comprising usually fine grinding or honing marks, and the other is the deep valleys. The former carries the load, the latter acts as an oil channel and debris trap.

This curve (figure 2.13) is split into three by means of drawing a secant to the region at the point of inflection corresponding to a 40% material ratio which is then drawn to intercept the axes. The height of the distribution is then split into three, $R_{pk}$, $R_{sk}$ and $R_{vk}$; $R_{pk}$ is a measure of the peaks, $R_{sk}$ measures the basic 'core' or 'kernel' of the surface and $R_{vk}$ estimates the depth of the valleys.

#### (a) Evaluation process for determining the parameters from the material ratio curve

*(i) Calculating the parameters $R_k$, $M_{r1}$ and $M_{r2}$ [57]*
The equivalent straight line, calculated according to the procedure, intersects the abscissae $M_r = 0\%$ and $M_r = 100\%$. These points define the intersection lines which determine the roughness core profile by dividing off the protruding peaks and valleys.

The vertical distance between these intersection lines is the core roughness depth $R_k$. Their intersections with the material ratio curve define the material portions $M_{r1}$ and $M_{r2}$. If these intersection points lie above the highest point or below the lowest point on the material ratio curve then they are set to give $M_{r1} = 0\%$ and/or $M_{r2} = 100\%$. In this case the parameters $R_{pk}$ and/or $R_{vk}$ are then zero.

*(ii) Calculating the equivalent straight line*
The equivalent straight line is calculated for the centre section of the material ratio curve which includes 40% of all the measured profile points. This 'central region' lies where the secant of the material ratio curve over 40% of the material portion shows the smallest gradient (see the figure).

A straight line is calculated for this 'central region' which gives the minimum mean square deviation in the direction of the profile ordinates. This method was developed for the German car industry and this and its variants are proving to be useful. It should be mentioned, however, that the straight part of the curve results more from the way the curve is plotted than because it is a reality. If the bearing or material ratio curve were plotted on probability paper rather than linearly there would be no 's' shape for a random surface, only a straight line.

To determine the region for the equivalent straight line calculation move a secant ($\Delta M_r = 40\%$) along the material ratio curve. The secant with the smallest gradient establishes the 'central region' of the curve for the equivalence calculation.

Note:

To ascertain the material ratio curve, the class widths of the ordinates of the roughness profile should be selected to be small enough to allow at least 10 classes to fall within the 'central region'.

For surfaces with very small roughness or surfaces with an almost ideal geometrical plateau such a fine classification may no longer be meaningful because of the limited resolution of the measurement system. In this case the number of classes used in the calculation of the equivalent straight line should be stated in the test results.

The 'central region' is ascertained as follows. Commencing at the top, the associated secants with the least material portion difference of $\Delta M_r \geqslant 40\%$ are determined. The least steep secant determines the 'central region' (if there are equal gradients then the one for the highest region is used).

*(iii) Calculation of $R_{pk}$ and $R_{vk}$*

The areas above and below the region of the material ratio curve that delimit the core roughness depth $R_k$ are shown hatched in figure 2.68. These correspond to the cross-sectional area of the profile peaks and valleys which protrude out of the roughness core profile.



**Figure 2.68** (*a*) Conventional plot, conventional surface; (*b*) probability plot, conventional surface and two-process striated surface.

The parameters $R_{pk}$ and $R_{vk}$ are each calculated as the side of a triangle which has the same area as the 'peak area' or 'valley area' (see figure 2.68(*a*)). The triangle corresponding to the 'peak area $A_1$' has $M_{r1}$ as the base and that corresponding to the 'valley area $A_2$' has $100\% - M_{r2}$ as the base.

Note:

The parameters according to this standard shall only be calculated if the material ratio curve is 's' shaped as shown in the figures and in practice only shows a single point of inflection. Experience has shown that this is always the case with lapped, ground and honed surfaces.

The conversion of 'peak area' and 'valley area' into equivalent area triangles is really a way of establishing the top limit of the $R_{pk}$ (and $R_{vk}$) parameter value.

From figure 2.68

$$A_2 = \frac{R_{vk}}{2}(100 - MR2)\mu m \tag{2.202}$$

$$A_1 = \frac{R_{pk}}{2} MR2 \mu m. \tag{2.203}$$

For the conversion to equivalent area triangles instead of the area, the curve itself allows $R_{pk}$ and $R_{vk}$ to be slightly smaller. One represents the fine process and the other the coarse process. In this plotting the kernel disappears altogether. Nevertheless it is proving useful and this is what counts.

Because of the rather artificial form of the 's' shape needed to define $R_{pk}$, $R_{vk}$ and $R_k$ a lot of research is under way to find a more rigorous breakdown. One way utilizes the fact that many surfaces are Gaussian or at least reasonably close in nature. This has led to the use of the material probability curve [58]. In this the percentage of material is plotted on a cumulative probability plot. In this way a straightforward finishing process shows itself as a straight line whereas a surface such as plateau honing shows itself as two distinct lines (figure 2.68(*b*)) which intersect. How to determine the point of intersection accurately and how to identify the three surface striations $R_{pk}$, $R_k$ and $R_{vk}$ from this type of presentation have caused problems. One rather complicated way [58] is briefly outlined below. It is not simple because a number of non-linear effects can confuse the situation. These are:

(1) the presence of debris or rogue peaks in the data,
(2) non-statistical valleys,
(3) the curvature of the material probability curve in the region of the transition from one process to another.

One technique for characterizing the curve, due to Scott [120,128], is as follows:

1. Fit a conic (assumed non-elliptical) to the material probability curve $z = Ax^2 + Bxz \times Cz^2 + Dx + E$ where $x$ is the material probability curve expressed in standard deviations and $z$ is profile height. From this conic the asymptotes are determined.
2. The asymptotes are bisected with a line which is drawn to intersect the conic. This intersection can be used as the material ratio point at which the processes change. In any event it is a first approximation to it.
3. The second derivative of the material ratio curve is computed using a data window of (0 .025 standard deviations. This computation is carried out upwards above the point of intersection and downwards below the point of intersection separately. When the magnitude of the second differential exceeds six standard deviations in either region this determines the points of the upper plateau limit (UPL) and lower plateau limit (LPL) respectively.
4. The $z$ axis of the material probability curve is normalized.
5. Better asymptotes can be constructed at this point specifically by fitting linear regression lines through the regions from the intersection to UPL and LVL (lower valley limit) respectively. The bisector can be redrawn more accurately.
6. Further boundary points of the lower plateau limit (LPL) and upper valley limit (UVL) are found by bisecting the angle between the asymptotes and the principal bisector.
Thus the material probability curve now has three regions.
7. A linear regression is then performed within each region of the original non-normalized material probability curve. The slope of the line in the region UPL to LPL can be called $R_{qp}$ while the slope of the line in the region UVL to LVL can be called $R_{qv}$. The intersection point $M_{pt}$, is the material ratio at the plateau-to-valley intersection.

It has to be made clear that these alternatives to $R_{pk}$, $R_k$ and $R_{vk}$ (i.e. $R_{qp}$, $M_{tp}$, $R_{qv}$) are only tentative and are included to illustrate the diversity of ways of characterizing multiprocess surfaces. Both methods have their problems. The former involves characterizing a rather arbitrary 's' shape into three regions. The latter assumes as a basis the Gaussian nature of surfaces, and utilizes somewhat arbitrary curvature criteria to identify the plateau and valley regions. Both methods are useful; which, if either, will ultimately be used remains to be seen.

The discussions above have concentrated on the height discrimination of striated surfaces. It should not be forgotten that any surface generated by a number of processes in which the surface is not completely eradicated by the following one has complex spatial characteristics. This means, for example, that deep valleys can be very wide and can cause considerable distortion of any filtered profile. Hence if a conventional instrument is to be used to measure striated surfaces the longest filter cut-off possible should be used otherwise errors will be introduced. For a detailed discussion of these problems see the analysis by Whitehouse [119].

Figures 2.70 to 2.75 show the stages involved in getting the parameters $R_k$ etc. What is missing is the method of filtering. This is discussed in the next section. Figure 2.69 shows that the method is dependent on maximum peak dependence.



**Figure 2.69** Problem of identifying limiting maximum and minimum.



**Figure 2.70** Layering of profile—core.



**Figure 2.71** Identification of 40%.



**Figure 2.72** Identification of $M_{r1}$ and $M_{r2}$.

**Figure 2.73** Layering of profile $R_{k1}$, $R_{pk}$, $R_{vk}$.



**Figure 2.74** Identification of areas *A1, A2*.



**Figure 2.75** Material ratio curve.

**Figures 2.70–2.75** Choice of filter and sampling length.

The effect of not using a valley cut-off can be quantified. Basically the cut-off should be as long as possible, e.g. 2.5 mm, to prevent the mean line dropping into the valleys (figure 2.76).

*$R_k$ filtering*
This filtering technique is according to ISO 13565 pt. 1 and DIN 4776.

(*a*) 0.8 mm



(*b*) 2.5 mm



**Figure 2.76** Choice of sampling lengths.

(*a*) Profile

(*b*) Modified profile

**Figure 2.77** Standard 2CR filter.

*Filtering process to determine the roughness profile*

The filtering process is carried out in several stages given the modified profiles.

The first mean line is determined by a preliminary filtering of the primary profile with the phase correct filter in accordance with ISO 11562 using a cut-off wavelength $\lambda$c in accordance with clause 7 and corresponding measuring conditions in accordance with Table 1 of ISO 3274:1996. All valley portions which lie below this mean line are removed. In these places the primary profile is replaced by the curve of the mean line.

The same filter is used again on this profile with the valleys suppressed. The second mean line thus obtained is the reference line relative to which the assessment of profile parameters is performed. This reference line is transferred to the original primary profile and the roughness profile according to this part of ISO 13565 is obtained from the difference between the primary profile and the reference line.

Selection of the cut-off wavelength $\lambda$c = 0.8 mm. In justified exceptional cases, $\lambda$c = 2.5 mm may be selected and this should be stated in the specification and test results.

**Table 2.9** Relationship between the cut-off wavelength $\lambda$c and the evaluation length *ln*

| $\lambda$c | *ln* |
|---|---|
| 0.8 | 4 |
| 2.5 | 12.5 |

Both the $R_k$ filter and the long cut-off are compromises. This is inevitable because there is more than one process present in the profile. Attempting to embody both in one procedure is bound to produce errors.

*Scope*

ISO 13565 describes a filtering method for use with surfaces that have deep valleys below a more finely finished plateau, with a relatively small amount of waviness. The reference line resulting from filtering according to ISO 11562 for such surfaces is undesirably influenced by the presence of the valleys. The filtering

(a) Unfiltered primary profile (valleys shown hatched)



(b) Unfiltered primary profile after suppression of valleys



(c) Position of the reference line in the primary profile



(d) Roughness profile in accordance with this standard

**Figure 2.78** $R_k$ filtering.

approach described in figure 3.26 suppresses the valley influence on the reference line such that a more satisfactory reference line is generated. However, a longer cut-off should always be tried [119].

*Normative references*

The following standards contain part of ISO 13565. At the time of publication, the editions indicated are valid.

ISO 3274:1996, Geometrical Product Specifications (GPS)—Surface texture: Profile method — nominal characteristics of contact (stylus) instruments.

ISO 4287:1997, Geometrical Product Specifications (GPS)—Surface texture: Profile method — terms, definitions and surface texture parameters.

ISO 11562:1996, Geometrical Product Specifications (GPS)—Surface texture: Profile method

*Definitions*

For the purposes of this part of ISO 13565, the definitions given in ISO 3274 and ISO 4287 apply.

**ISO 13565-1:1996 (E)**

*Reference Guide*

To measure profiles in accordance with this part of ISO 13565, a measuring system which incorporates an external reference is recommended. In case of arbitration the use of such a system is obligatory. Instruments with skids should not be used.

*Traversing direction*

The traversing direction should be perpendicular to the direction of lay unless otherwise indicated.

Another point which causes confusion is concerned with the material ratio obtained from the profile and that obtained over an area. In fact these are equal. The profile *MR* squared is not the areal value. See figure 3.27(a) which is a series of blocks of profiles. This can be moved about to look like figure 3.37(b) which is still the same *MR* value yet it is the complete areal view from whichever direction is taken.

If a large number of tracks are made on figure 3.27 the material ratio is always the same. It is not valid to take just one reading on the lower surfaces. This equivalence is one reason why it is valid to base the use of the various functional parameters on the material ratio.

A completely different approach which is of considerable importance is the modelling of the surface profile by means of time series analysis. Because of its importance some space here will be devoted to the subject. It will also be examined again in the section on surface generation in chapter 3.

*2.1.7.6  Time series analysis methods of characterization—the characterization of spatial information*

The basic idea behind this technique is to consider the surface profile as being part of a system [56, 59–63, 65–67]. In fact it is considered to be the output from a system (as might be seen as a force or temperature or whatever) when a completely random input (white noise) is put into it (figure 2.79).

By assuming the input and knowing the output, the surface can be considered to be an operation on white noise; in other words it modifies the white noise in a way that can be viewed as a classification of the surface. This has some appealing features because surfaces are produced by a manufacturing system, each block of



**Figure 2.79**   Time series method characterization.

which has transfer characteristics (transfer function). Why not consider the surface as part of the overall system and deal with it accordingly (see chapter 7)? This not only makes sense, it also fits into the basic philosophy of chapter 6 where the whole manufacturing system is considered to centre on the component produced. Unfortunately, this time series subject can fit into processing methods of chapter 3 as well as here because it relates to digital or discrete data. However, as it is a form of characterization it will be dealt with briefly in this section.

Linear physical systems can be categorized in terms of with respect to their dynamic behaviour by systems of simultaneous linear differential equations. The simultaneous systems can be transformed into a single differential equation which is usually denoted in discrete form. (For simplicity the time domain will be used to follow existing terminology but it should be remembered that, in practice, the time axis is spatial. It can usually be related via the instrument to time.)

$$a_n \frac{\mathrm{d}^n z(t)}{\mathrm{d}t^n} + a_{n-1} \frac{\mathrm{d}^{n-1} z(t)}{\mathrm{d}t^{n-1}} + \ldots + a_0 z(t) = b_m \frac{\mathrm{d}^m u(t)}{\mathrm{d}t^m} + \ldots + b_1 \frac{\mathrm{d}u(t)}{\mathrm{d}t} + u( \tag{2.204}$$

with $m < n - 1$; $u(t)$ is white noise with zero mean and $\sigma^2$ variance. The discrete form of this equation is denoted by an autoregressive moving average equation (ARMA) model of order $N, M$:

$$(A_N D^N + A_{N-1} D^{N-1} + A_{N-2} D^{N-2} + \ldots + A_1 D + A_0) z(t) = (B_M D^M + B_{M-1} D^{M-1} + \ldots + B_1 D + B_0)i \tag{2.205}$$

where $z(t)$ is the profile, D is a difference operator ($D^K z_t, = z_{t-k}$) or shift operator, $A_k$ are autoregressive coefficients, $B_k$ are moving average coefficients and $u(t)$ is the supposed input to the system (assumed to be white noise or a variant).

The problem is that of determining the values of the $A$ and $B$. Once found these can be thought of as describing the profile.

For simplicity, to see how this is done, given a measured profile signal in digital form $z_t$, a signal $z'_t$, which is generated by an ARMA $(N, M)$ equation, shall be used to approximate this measured signal.

To determine the $A$ and $B$ values the least-squares method can be used. So

$$\sum_{i-1}^{n} (z_i - z_i')^2 = S \tag{2.206}$$

where $S$ is a minimum, and $z_i - z'_i$ is the error, or residue.

In principle, once the $A$ and $B$ are known the spectrum of the profile can be obtained and from this the spectral moments derived. Thus

$$P(\omega) = \frac{\sigma_u^2}{2\pi} \frac{[1 + \sum_{i=1}^{n} (B_i/B_0) \exp(-j i\omega)]^2 (A_0/B_0)}{[1 + \sum_{k=1}^{n} (A_k/A_0) \exp(-j k\omega)]^2} \tag{2.207}$$

where $\sigma_u$ is the standard deviation of $u$. If the model is autoregressive only then the numerator in equation (2.207) disappears.

It has been pointed out that the estimation of the coefficients in the full ARMA model is difficult [63] because it is a non-linear estimation and has to be achieved by iteration. However, for simplicity, if the model is reduced to an autoregressive one then clearly a linear system exists and the parameters can be found.

Take for example the case of an AR (2) model: where

$$E[(z_i - (A z_{i-1} + B z_{i-2} + u_i))^2 \tag{2.208}$$

is a minimum.

Taking expected values and differentiating with respect to $A$ and $B$ respectively gives the Yule–Walker equations. These are solved to get $A$ and $B$.

In this case if the variance of the profile is unity and has zero mean and $E[z_i z_{i-1}] = \rho_1$ and $E[z_i z_{i-2}] = \rho_2$, then

$$\begin{pmatrix} 1 & \rho_1 \\ \rho_1 & 1 \end{pmatrix} \begin{pmatrix} A \\ B \end{pmatrix} = \begin{pmatrix} \rho_1 \\ \rho_2 \end{pmatrix} \tag{2.209}$$

from which

$$A = \frac{\rho_1(1 - \rho_2)}{(1 - \rho_1^2)} \quad B = \frac{\rho_2 - \rho_1^2}{(1 - \rho_1^2)} \tag{2.210}$$

or

$$z_i = \frac{\rho_1(1 - \rho_2)}{(1 - \rho_1^2)} z_{i-1} + \frac{(\rho_2 - \rho_1^2)}{(1 - \rho_1^2)} z_{i-} \tag{2.211}$$

is the autoregressive equation representing the profile, where $\rho_1$ and $\rho_2$ are points taken from the autocovariance function (autocorrelation function of the surface at spacings equal to one and two sample intervals respectively).

It is interesting to note that this result could have been obtained directly from the multinormal distribution in equation (2.65), which describes the probability density functions of $z_i$, $z_{i-1}$ and $z_{i-2}$. Thus the probability density of $z_i$, given $z_{i-1}$ and $z_{i-2}$ is given by

$$p(z_i | z_{i-1}, z_{i-2}) = \frac{\sqrt{z - \rho_1^2}}{\sqrt{2\pi(1 - \rho_2)(1 + \rho_2 - 2\rho_1^2)}} \exp\left(-\frac{\{z_i - \rho_1(1 - \rho_2)z_{i-1} + [(\rho_1^2 - \rho_2)/(1 - \rho_1^2)]z_{i-2}\}}{[2(1 - \rho_2)(1 + \rho_2 - 2\rho_1^2)/(1 - \rho_1^2)]}\right) \tag{2.212}$$

The denominator of the exponential represents twice the variance of the constrained $zi$, which means that this is the variance of the white noise random distribution from which the $z_i$ can be picked, given that the mean value of the white noise distribution is determined by the square root of the numerator — which depends on the specific values of $z_{i-1}$ and $z_{i-2}$ as shown in equation (2.212). A further check can be used by the method outlined by Watson and Spedding [62].

In general for an AR($N$) model $\sigma_u^2 = (1 - A\rho_1 - B\rho_2 - \ldots - N\rho_n)\sigma_z^2$, where the $A$, $B$, etc, are taken from equation (2.205).

The above variances and means follow for a simple AR(2) surface:

$$\rho_1 = \frac{A}{1 - B} \quad \text{and} \quad \rho_2 = \frac{A^2}{1 - B} + B. \tag{2.213}$$

The time series approach has been criticized because it nominally applies to random waveforms (i.e. Gaussian). However, variations on the basic approach have been made which allow asymmetrical distributions to be modelled.

For non-normal distributions using the same method as equation (2.213) the relationships have been worked out [56] for the skew and kurtosis of the profile in terms of those of the generator. The main point is that skew and kurtosis can be forecast, within certain limits, for the non-Gaussian case, so the former criticism

is not really a problem. The issue is not so much how to characterize surfaces, as is addressed here, but how to simulate surfaces in a computer for functional tests. The surface generation problem will be considered in chapter 3 section 3.11.

DeVries [63] discusses how AR(1) and AR(2) models of surfaces with non-Gaussian distributions can be achieved. The general case for ARMA models was developed by Davies *et al* [56] and subsequently applied to surfaces [62].

The philosophy is that any ARMA model can be expressed in an infinite series of moving averages (MA). Thus

$$z_x = a_x + c_1 a_{x-1} + c_1 a_{x-2} + . \tag{2.214}$$

where $c_1$, $c_2$ are constants (e.g. for the AR(1) process) and

$$z_i + \varphi_1 z_{i-1} + a_i \qquad \text{but} \qquad z_{i-1} = \varphi_2 z_{i-2} + \tag{2.215}$$

so $z_i = ai + \varphi_1 a_{i-1} + \varphi_1^2 a_{i-2} + \ldots$ and so on. $c_1$, $c_2$, etc, can be found from repeated use of (2.215) from which the skew of the profile is

$$Sk_p = \frac{\sum_{i=0}^{q} c_i^3}{(\sum_{i=0}^{q} c_i^2)^{3/2}} Sk_{rn} \tag{2.216}$$

and the kurtosis of the profile is

$$K_p = \frac{\sum_{i=0}^{q} c_i^4 K_{rn} + 6 \sum_{i=0}^{q-1} \sum_{j=i+1}^{q} c_i^2 c_j^2}{(\sum_{i=0}^{q} c_i^2)^2} \tag{2.217}$$

where $q$ is the number of terms in the series of equation (2.214), $Sk_{rn}$ is the skew of the random number generator and $K_{rn}$ is the kurtosis of the random number generator.

These formulae show that it is therefore possible to specify the skew and kurtosis of the generated profile as characterization parameters (for the shape of the amplitude distribution) together with the parameters of the ARMA model for the spatial parameters. So, for an AR(2) model, there would be five parameters in total (including $R_a$ or $R_q$) to specify the surface profile.

Autoregressive methods are really ways of modelling the spectrum of the surface, either through the spectra or in terms of the correlation coefficients.

Another of the problems associated with ARMA methods is that of trying to incorporate some measure of periodicity. One way suggested has been to separate out the deterministic component from the random component and to characterize them separately. Wold's discrimination [64] allows the breakdown of waveforms into random and deterministic components, at least in principle. As mentioned before, it is highly suspect to allow these theorems to mask the intrinsic difficulty of measuring the surface.

DeVries [63] and others [43] characterize the surface profile by first determining the order of the AR model and the associated parameters and then transforming to the spectrum as in equation (2.207). From the spectrum the moments, say $m_0$ and $m_2$, can be found and from these a parameter describing the RMS slope of the profile or the average wavelength $2\pi(m_0/m_2)^{1/2}$ is obtained. The advantage of such an approach is that, whereas the parameters $A$ and $B$ rely heavily on the sample, a smoothed version of the power spectrum does not.

A typical order required for an AR model of a turned surface worked out to be over 10 and for other complex surfaces over five, which obviously means that too many parameters are being involved—hence the need to reduce them by taking moments of the spectrum instead. The benefits of this approach over that

of a direct estimation of the power spectrum are that the time series model does allow a systems approach in the sampled data domain and so allows surface simulation of the required spectral characteristics.

A more recent targetted classification of surfaces has been made by Lukyanov [47] who classifies the autocorrelation function into basically only two types; these have been shown earlier, one representing random surfaces and the other narrow band surfaces. Each type has a subset, for example

$$\begin{aligned} \text{type I:} \quad &\text{either correlation} \frac{1}{1+\alpha\tau^2} \text{ or } \exp(-\alpha\tau^2) \\ \text{type II:} \quad &\text{either } \frac{\cos\omega_1\beta}{1+\alpha\tau^2} \text{ or } \exp(-\alpha\tau^2)\cos\omega_2\tau. \end{aligned}$$

(2.218)

### 2.1.7.7 Possible methods of classification based on a Fourier approach

Transforms have been mentioned earlier in this chapter, in particular with respect to height information. The basic idea of breaking down the surface waveform into a series of numbers by means of a Fourier analysis is not new. In this method the idea is that, when applied to a surface, the coefficients of the operator can be used as a basis for characterization. It is tempting to apply this technique to other than the Fourier series or Fourier transform. Quite a lot of effort has gone into investigating alternatives. Transforms such as those of Walsh, Hadamard and Wigner [65–67] have been and still are being tried. These other transforms are related to the Fourier transform. How they are applied to surfaces is considered in chapter 3 on data processing but a brief mention of them will be given here.

#### (a) Hartley transform
This is the difference between the real and imaginary parts of the Fourier transform and is defined as

$$\begin{aligned} F(u,v) &= F_{\text{real}}(u,v) + jF_{\text{imag}}(u,v) \quad \text{the Fourier transform} \\ H(u,v) &= F_{\text{real}}(u,v) - F_{imag}(u,v) \quad \text{the Hartley transform} \\ jF_{\text{imag}}(u,v) &= \tfrac{1}{2}(F(u,v) - F(-u,-v)). \end{aligned}$$

(2.219a)

There are a few interesting differences between these two transforms. One is that the Hartley transform does not use the complex plane as such, unlike the Fourier transform. For this reason, it is not really sensitive to phase and there can be some ambiguity in sign.

The transform has a similar equation relating its two argument functions as the power spectrum and the correlation function. Thus

$$H(u) = \int_0^L f(x)\text{cas}(kux)\text{d}x.$$

(2.219b)

Here the operator cas means cos and sin.

The original intention of using the Hartley transform was for speed. However, its computation is not as simple or straightforward as the Fourier transform, so it has not been found to be a serious competitor to the Fourier transform. This is given in chapter 3 together with the transforms below in rather more detail.

#### (b) Square wave function

##### (i) Walsh functions
Instead of correlating the signal with a sine wave as in Fourier analysis a square wave is used. This has some advantages because the clipped signal is very easy to achieve digitally and has already been used in the evaluation of the pseudo-autocorrelation function using the Steltjes integral.

The square wave signal is also faster to use than the sine wave. However, not surprisingly, the very fact that it is a square gives it extra usefulness in examining surfaces having sharp changes or discontinuities. It does not, however, break down into very meaningful coefficients if the surface is continuous and undulating, which is where the Fourier transform comes into its own.

*(ii) Hadamard function*

There are other transforms, which use a binary approach, such as the Hadamard transform. This again suffers from the fact that if the surface wave is anything like sinusoidal instead of giving a simple single coefficient, it breaks the sine wave down into a rather unconvincing set of coefficients.

Consequently, although these and similar transforms like the wavelet transform have not been taken up seriously in surface metrology as yet, they may well be in the future. This is because the nature of surface texture is changing.

It may be that as the presence of sharp features on surfaces such as ridges, crevices or scratches increases, the square function transforms will be used more often.

### 2.1.7.8   A general note about space-frequency functions

Most engineers accept the Fourier transform and Fourier analysis as a meaningful way to break up a signal into an alternative form, i.e. sinusoids (equation (2.220)). Many signals have an oscillating base as in vibration. This representation has been very useful.

$$F(w) = \int_{-\infty}^{\infty} z(x)\exp(-jwx)dx \qquad (2.220)$$

Equation (2.220) gives the amplitude and phase of the coefficient $F(w)$ when the function $z(x)$ is multiplied by the exponential function and integrated over all $x$. The value $F(w)$ owes much of its stability to the fact that the integral range is very wide. Unfortunately this attribute has drawbacks. If there is a change in the statistical nature of $z(x)$ over a small range of $x$ it is completely swallowed by the integral. Local variations cannot be picked out. Often the global average such as in equation (2.220) is a measure of good performance. For example, in roundness measurement if the average value of the third harmonic is low, then it can be deduced that there is no workpiece clamping problem. On the other hand if one of the tool used in face milling is broken, e.g. the first tooth on a twelve tooth tool, measuring the twelfth harmonic will not show up the presence of the broken tooth. What is needed is a function which has the attributes of a differentiator to highlight small changes as well as an integrator to have stable results. There have been many attempts to solve this dual problem. The space-frequency (or time-frequency) functions are a general class of function designed to have the dual role.

One function—the structure function (first used in surface metrology in 1971 by Whitehouse[126]) has the elements of a dual function. It is defined as

$$E\left(z\left(x - \frac{\tau}{2}\right).z\left(x + \frac{\tau}{2}\right)\right) = 2\left(\sigma^2 - \sigma^2\rho\right) \qquad (2.221)$$

where $\sigma$ is $R_q$ the root mean square roughness and $\rho$ is the correlation.

The term

$$z(x) - z(x + \tau) \sim \tau z'(x) \qquad (2.222)$$

so

$$S(\tau) = E\left(z(x) - z(x + \tau)^2\right) = \tau \frac{2_1}{L}\int_0^L z'(x)^2 \, dx \qquad (2.223)$$

or
$$S(\tau) = E\left(z\left(x - \frac{\tau}{2}\right).z\left(x + \frac{\tau}{2}\right)\right) \tag{2.224}$$

which is very closely related to the Wigner function kernel.

Unfortunately this function has only one argument spatially and none in frequency. Therefore the search has been for alternatives to the simple random process power spectrum and autocorrelation.

In equation (2.224) there is the one operational parameter $\tau$. To be useful another parameter is needed. It seems that an obvious extension is to take the Fourier transform of the kernel.

So

$$Y(\tau, \beta) = \int_{-\infty}^{\infty} z\left(x - \frac{\tau}{2}\right).z\left(x + \frac{\tau}{2}\right)\exp\left(-j\beta x\right)dx \tag{2.225}$$

In equation (2.225) there are two variables $\tau$ and $\beta$ which satisfies the criterion for isolating positions in space and time. The equation (2.225) is obviously contrived but it is very close to the Wigner distribution which has a formal physics pedigree and does develop from the Fourier transform and autocorrelation. The kernel in fact is made up of an autocorrelation term, equation (2.225), and a Fourier term, $\exp(-j(x))$. It fits in well with an engineering approach.

The ambiguity function has a formula similar to equation (2.22) so it also falls under the Fourier umbrella. An example of their use is found later in the book in chapter 5.

It will be seen in chapter 6 and elsewhere that the autocorrelation function and the power spectrum are very useful. Their use derives mainly from the fact that because they are insensitive to phase change they are inherently stable. In the case of the autocorrelation function this revolves around the kernel

$$\langle z(x)z(x + \tau)\rangle. \tag{2.226}$$

It is fairly obvious that this sort of average, although useful for stationary signals (in the spatial sense), is not very revealing with non-stationary signals, that is where the statistics of the surface change as a function of $x$. This is because for a given $\tau$ the dependence on $x$ is lost because the kernel is integrated with respect to $x$. To be useful for characterizing non-stationary signals $x$ has to be conserved. One very powerful possibility is to redefine the kernel in such a way that the value of $x$ is retained. Such a possibility is given by

$$C(\tau, x) = \left\langle z(x + \tau/2)z^*(x - \tau/2)\right\rangle. \tag{2.227}$$

The ensemble average is centred on $x$. This effectively means that the expression for the average is centred on $x$ rather than being a distance $x$ from an arbitrary origin.

Equivalently a short time spectrogram can be postulated, $P(\omega, x)$. This is

$$P(\omega, x) = \left|\int_{-\infty}^{\infty} z(\beta)h(\beta - x)\exp(-j\omega\beta)d\beta\right|^2 \tag{2.228}$$

where $h$ is a window function located at $x$ and which weights $z(\cdot)$. This results in an average $\omega$ centred on $x$ and extending over the width of the window. It does not give the instantaneous frequency at $x$. $\beta$ is a dummy space variable.

Using this thinking directs attention to two types of function, the Wigner function and the ambiguity function. Some of their principal properties are given below.

*(a) Wigner distribution functions and ambiguity functions*

The Wigner function is

$$W(x,\omega) = \int_{-\infty}^{\infty} z\left(x - \frac{\chi}{2}\right) z^*\left(x + \frac{\chi}{2}\right) \exp(-j\omega\chi)\mathrm{d}\chi \tag{2.229}$$

and the ambiguity function is

$$A(\chi,\overline{\omega}) = \int_{-\infty}^{\infty} z\left(x - \frac{\chi}{2}\right) z^*\left(x + \frac{\chi}{2}\right) \exp(-j\overline{\omega}x)\mathrm{d}x \tag{2.230}$$

The very unusual feature of these is that the functions $W(x, \omega)$ and $A(\chi, \overline{\omega})$ can be obtained equally easily from the frequency domain. Thus

$$W(x,\omega) = \frac{1}{2\pi}\int_{-\infty}^{\infty} F\left(\omega - \frac{\overline{\omega}}{2}\right) F^*\left(\omega + \frac{\overline{\omega}}{2}\right) \exp(+jx\overline{\omega})\mathrm{d}\overline{\omega} \tag{2.231}$$

$$A(\chi,\overline{\omega}) = \frac{1}{2\pi}\int F\left(\omega - \frac{\overline{\omega}}{2}\right) F^*\left(\omega + \frac{\overline{\omega}}{2}\right) \exp(-j\chi\overline{\omega})\mathrm{d}\omega. \tag{2.232}$$

It is obvious from this that both functions can be arrived at equally easily from the space or frequency direction. They can be said to exist in between the two — hence the dual term space-frequency. They can also be expressed in terms of each other:

$$W(x,\omega) = \frac{1}{2\pi}\iint_{-\infty}^{\infty} A\left(\chi,\ \overline{\omega}\right) \exp[-j(\omega\chi - \overline{\omega}x)]\mathrm{d}\overline{\omega}\mathrm{d}\chi$$

$$A(\chi,\overline{\omega}) = \frac{1}{2\pi}\iint_{-\infty}^{\infty} W\left(x,\ \omega\right) \exp[-j(\overline{\omega}x - \omega\chi)]\mathrm{d}x\mathrm{d}\omega. \tag{2.233}$$

By looking at the equations it can be seen that both functions utilize what are effectively the Fourier kernel $\exp(\cdot)$ and the correlation kernel $\langle z(x)z(x + \tau)\rangle$ thereby making use of the benefits of both. This means that both functions $W(\cdot)$ and $A(\cdot)$ are more flexible than either the autocorrelation or the power spectrum.

If it comes to a choice between $W(x, \omega)$ and $A(\chi, \overline{\omega})$, the winner is marginally the Wigner function (as far as can be seen at present) because it alone retains the actual values of $x$ and $\omega$ and so is most suitable for non-stationary characterization.

In principle, therefore, the possibility of being able to characterize the average statistics of the surface, and also the non-typical characteristics of the surface such as defects or flaws, using one function is now a possibility.

The fact that the values of $x$ and $\omega$ are retained makes the Wigner function behave like a two-dimensional convolution. Making $x$ or $\omega$ zero simulates filtering in the space or frequency domain. On the other hand, the ambiguity function acts more as a correlator, $\chi$ and $\overline{\omega}$ being effectively lags. For $\overline{\omega} = 0$ the ambiguity function is the autocorrelation function.

In effect, in the space-frequency domain the ambiguity domain moves, whereas the Wigner domain expands or shrinks. In both cases there is a fixed element. In the Wigner function it is the position of the centre of the domain; in the ambiguity function it is the size of domain (i.e. the bandwidth-extent product is constant).

It is clear that the Wigner and ambiguity functions are likely to work on a surface profile. It is obviously impossible to visualize the areal functions, which are of four dimensions in space-frequency.

Together, the ambiguity function and the Wigner distribution function have much potential. The problem now arises of formulating a way of characterization based on these functions. There is now far too much information available. To make a sensible and practical attempt to reduce the data, it seems plausible to resort to taking the moments of the functions as for previous functions. Because of the double arguments, the

moments can be in frequency, space or both. Thus the Wigner function has moments which are given below and which are surprisingly simple:

$$\text{first frequency moment} = \text{Im}\left(\frac{z'(x)}{z(x)}\right) = \frac{\int \omega F(\omega)\mathrm{d}\omega}{\int F(\omega)\mathrm{d}\omega} \tag{2.234}$$

$$\text{second frequency moment} = -\frac{1}{2}\text{Re}\frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{z'(x)}{z(x)}\right)$$

$$= \frac{\int \omega^2 F(\omega)\mathrm{d}\omega}{\int F(\omega)\mathrm{d}\omega} - \left(\frac{\int \omega F(\omega)\mathrm{d}\omega}{\int F(\omega)\mathrm{d}\omega}\right)^2 \tag{2.235}$$

$$\text{first spatial moment} = \text{Im}\left(\frac{F'(\omega)}{F(\omega)}\right) = \frac{\int xz(x)\mathrm{d}x}{\int z(x)\mathrm{d}x} \tag{2.236}$$

$$\text{second spatial moment} = -\frac{1}{2}\text{Re}\frac{\mathrm{d}}{\mathrm{d}\omega}\left(\frac{F'(\omega)}{F(\omega)}\right)$$

$$= \frac{\int x^2 z(x)\mathrm{d}x}{\int z(x)\mathrm{d}x} - \left(\frac{\int xz(x)\mathrm{d}x}{\int z(x)\mathrm{d}x}\right)^2. \tag{2.237}$$

This form is obvious when it is realized that for the spatial moments $\omega = 0$ and for the frequency moments $x = 0$. The moments are in fact the set obtained from the profile and spectrum directly—only they are linked together by the Wigner function.

Similar results are obtained for the ambiguity function. The global first moment for frequency is given by

$$\frac{1}{2\pi}\int\int \omega W(x,\omega)\mathrm{d}x\mathrm{d}\omega = \frac{1}{2\pi}\int \omega |F(\omega)|^2 \mathrm{d}\omega / \sigma_F^2 \tag{2.238}$$

$$= -\int x|z(x)|^2 \mathrm{d}x / z^2 \quad \text{for space} \tag{2.239}$$

where $\sigma_F^2$ is the variance of $F(\omega)$ and $\sigma^2$ is the variance of $z$.

As an example of how these moments can be used, consider the dynamics of machine tools. For this it is necessary to examine the types of signal which could be imprinted on the surface by vibration of the tool. It turns out that amplitude modulation characteristics are revealed by the second local moment in frequency whereas non-stationarity tends to be revealed by the first-order local moments in frequency [72]. The spatial moments can isolate scratches and faults in the same way.

Thus for $z(x) = a(x)\exp(\mathrm{j}\varphi(x))$ (pitch vibration of tool)

$$\text{zeroth moment (frequency)} = p(x) = a^2(x)$$
$$\text{first moment (frequency)} = \Omega(x) = \rho(x)$$
$$\text{second moment (frequency)} = m(x) = -\frac{1}{2}\frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{a'(x)}{a(x)}\right). \tag{2.240}$$

The instantaneous power and frequency are found at $x$.

If the surface wave is due to yaw vibration of the tool

$$z(x) = a\exp(\mathrm{j}(\alpha/2)x^2)$$
$$p(x) = \alpha^2 \tag{2.241}$$
$$\Omega(x) = \alpha x.$$

For frequency modulation

$$z(x) = a\exp[\omega_0 x + \varphi_0 + b\sin(\omega_m + \varphi_m)]$$
$$p(x) = \alpha^2 \tag{2.241a}$$
$$\Omega(x) = \omega_0 + b\omega_m\cos(\omega_m x + \varphi_m).$$

From equation (2.241) it is clear that the first moment clearly identifies the type of signal.

The question arises as to whether the space-frequency functions can be useful for characterizing other surface features. The answer is probably positive in the case of surface defects and flaws.

How does a scratch of width $2x_1$ and position $x_0$ on the surface fare with this moment treatment? This has a spectral response given by

$$F(\omega) = \left(\frac{1}{j\omega}\right)^2 \left\{-\exp(-j\omega x_0) + 2\exp[-j\omega(x_0 + x_1)] - \exp[-j\omega(x_0 + 2x_1)]\right\}. \tag{2.242}$$

This gives the first moment value of $x_0$ and the second moment equal to the dispersion $x_1^2\big/24$ which corresponds to the position and width of the defect. From this data it seems plausible to suggest that useful information on defect classification can result by using the Wigner function. The ambiguity function would be useful in classifying the shape of the defect rather than its size and position.

### (b) Gabor transform [131]

Another function of the space-frequency type is called the Gabor transform. This is a function which makes use of the transform symmetry of the Gaussian-shaped pulse; that is, the shape of a Gaussian pulse in the time (or space) domain is the same in the frequency domain. So, the input signal is broken down into a linear combination of shifted Gaussian pulses in the time (space) domain. Obviously these same pulses can be thought of as also existing in the frequency domain. Because the responses to Gaussian pulses are well known they can be put in a look-up table and used instead of a calculation of transformation. This considerably speeds up the calculation (some details are given in section 3.8.5); it has not yet been used in surface metrology but may be in the future.

### Wavelet Transforms

The Wigner and ambiguity functions are time (space) functions. Both time (space) and frequency are to be selected so that different parts of the signal can be investigated by basically shifting the original analysis and also by changing the scale. The effective kernel in both contains an exponential term and a correlation term. In effect the kernel provides a short term function which is a mixture of a correlation term and a Fourier term. The slope of this depends on the actual values of the kernel arguments.

An alternative approach is provided by wavelet theory which has as its kernel an arbitrarily shaped function with two arguments: one position (or time) and the other a scale factor. The fact that the waveform is arbitrary gives tremendous flexibility. It can be represented in a number of forms, for example short-time Fourier transform decomposition of signal into basis functions, sub-band signal decompositions etc.

There is only one constant in the above list and this is the equation (i.e. the form of the wavelet transform).

$$W(ab) = \frac{1}{\sqrt{a}}\int h^*\left(\frac{t-b}{a}\right)z(t)dt \tag{2.243}$$

In the equation for the transform the wavelet $h(t)$ operates on the signal $f(t)$ under the integral. $W(a, b)$ is the wavelet transform.

Historically the first glimmering of the wavelet theory was due to Gabor [131] in 1946. What he did was to localize the signal by a 'window' of finite width before Fourier analysis was used. This is often called the Gabor representation. The penalty of this localization is loss of frequency resolution. However, this approach can isolate transients in signals.

One representation often used is to divide the frequency into cells of constant bandwidths e.g. octaves or decades rather than constant absolute size (i.e. giving the representation in terms of function of constant shape). In this guise the wavelet transform is characterized by position and scale rather than position and frequency. Wigner and ambiguity functions are genuinely time (space) frequency functions so the wavelet transform can be considered to be complementary to them both. The other point is that Wigner and ambiguity functions have fixed form, whereas wavelets can have any form providing that the basic equation is equation (2.243).

Wigner and ambiguity functions have components of autocorrelation and Fourier or power spectrum; the wavelet does not. It is not possible to evaluate moments in a general form as in Wigner and ambiguity functions. The freedom of the wavelet transform allows position and bandwidth to be simultaneously changed, whereas Wigner and ambiguity functions have fixed position and variable bandwidth or fixed bandwidth with flexible position.

The main problem with the use of wavelets is their variety. It started out with Gabor in 1946 who tried to localize features in the waveform under test by pre-processing with a finite duration window before the Fourier analysis, thereby ensuring some degree of localization albeit at expense of frequency resolution.

The key to the approach was to divide the frequency domain into cells of constant (usually octave) bandwidth thus giving a representation in terms of functions of constant shape (the original ones were Gaussian).

The main difference between the Fourier and wavelet approaches is as follows.

The Fourier representation of analysing the signal with a view to localization is twofold.

(i) splitting up the real axis into units of given length
(ii) representing the function is *each interval* by its Fourier series (i.e. a block transform).

Such a representation has position and frequency parameters. The wavelet on the other hand describes signals in terms of *scale* rather than *frequency*.

Thus the wavelet of $f(x)$ is the set of coefficients as equation (c) with two parameters—position index $n$ and scale index $m$.

The important thing to note is that the wavelet method is very versatile, unlike the Wigner, but its versatility means that the shape of the function has to match the application if good results are to be obtained. For example, it is no good trying to represent a sine with Gaussian shaped packets, it is a complete mismatch. Wigner treats every waveform equally.

The point here is that it is sensible to have some idea of the function before applying wavelet theory.

For general applications the wavelets should be orthonormal in the same way as in a Fourier series. Wavelets which have this property can be useful in the processing of surface data. Raja *et al* are probably the first investigators to use wavelets [132].

One of the problems with wavelet application is that to get best results the shape of the wavelet should match the feature of interest on the surface and yet at the same time have the orthonormal property indicted above. These two requirements are not necessarily easy to satisfy simultaneously. Take for example the hat transform shown earlier. This shape—or more precisely the negative of it $f(x) = -f(x)$ —is very suitable for identifying the unit grain impression left on the surface during grinding but it does not have the required orthogonal properties.

The wavelet technique does allow fractal analysis [78]. The multiscale property of possible fractal surfaces can be explored by having the wavelet scale variable step down in octaves. It is unlikely that many surfaces would scale down by more than two decades so there is a real problem of knowing through how many orders of magnitude self similarity holds true.

*Space frequency function wavelet theory*
The main point is that space frequency functions are used for dispersive functions, basically when the characteristics are changing in time (space) and/or frequency.

The wavelet transform method decomposes a space function (in surface metrology) into a set of basic functions or wavelets in the same way that a Fourier analysis breaks down a signal into sinusoids of different frequency

Thus
$$W(a,b) = \int_{-\infty}^{\infty} f(x)\left[\psi^*(a,b)(x)\right]dx \qquad (2.244)$$

Where $\psi(ab)(x)$ is given by equation (2.245) (see equation (2.243) for time version).

$$\psi(ab) = \frac{1}{\sqrt{a}}\psi\left(\frac{x-b}{a}\right) \qquad (2.245)$$

The basic wavelet is $\psi(x)$. It can be seen from equation (2.245) that $\psi(ab)(x)$ is a member of the wavelet set obtained from the base by a spatial shift $(x - b)$ and a scale reduction of $a$. These operations are shift and dilation.

It is easy to see that by changing $a$ and $b$ the wavelet can be positioned anywhere in the space domain $x$ by $b$ and adjusted in size by $a$.

In effect $a$ and $b$ are used to adjust the space and frequency location. Small $a$ produces a large wavelet which corresponds to a high spatial resolution and vice versa, large $a$ shrinks the wavelet function in the frequency domain to concentrate in a small region in frequency. Obviously poor spatial resolution goes with this in exactly the same way as in a Fourier transform pair.

This localization property in space and frequency is exactly the same as with the Wigner function. Small variations can be targeted and not lost in the large averaging of the Fourier transform method.

It is usual to arrange the various wavelets in such a way as to be orthogonal and hence independent of each other.

In discrete terms the input signal $f(x)$ is multiplied by discrete wavelet functions $\psi_{jk}(x)$ equivalent to equation (2.245)

Where
$$\psi_{jk}(x) = \frac{1}{\sqrt{2^j}}\psi\left(\frac{-2^j k}{2^j}\right) \qquad (2.246)$$

In equation (2.246) $2^j$ represents the frequency of each of the wavelet basis functions in equation (a) and $2^j k$ represents the location in space corresponding to the $b$ term in equation (2.245).

The wavelet coefficients $C_j k$ may be regarded as a space frequency map of the original signal.

$$C_j k = \int_{-\infty}^{\infty} f\left(x\right)\psi_{jk}^*\left(x\right)dx \qquad (2.247)$$

As already noted there are many variations of the wavelet theory. One popular variant, multi-resolution is given here.

This equation (2.246) is used to sample the input signal and its approximations at different resolutions. At each resolution the scaling coefficients and the wavelet coefficients are

$$C_{j+1,k} = \sum_{i=-\infty}^{\infty} g\left(i - 2k\right)d_{jk} \qquad (2.248)$$

$$d_{j+1,k} = \sum_{i=-\infty}^{\infty} h\left(i - 2k\right)djk \qquad (2.249)$$

Equations (2.248) and (2.249) represent a decomposition of the $(j-1)$ scaling coefficients into low frequency and high frequency terms. The terms $g$ and $h$ are high pass and low pass filters derived from the analysis wavelet $\psi$ and the scaling function.

One application uses orthogonal filters of length 4 units.

Thus

$$h(0) = \left(1 + \sqrt{3}\right)/4\sqrt{2}, h(1) = \left(3 + \sqrt{3}\right)/4\sqrt{2}$$
$$h(2) = \left(3 - \sqrt{3}\right)/4\sqrt{2}, h(3) = \left(1 - \sqrt{3}\right)/4\sqrt{2}$$
$$g(n) = (-1)^n h(3 - n)$$

(2.250)

The actual functions to multiply the raw data are very simple, as can be seen in the example above. Computationally they are fast and efficient especially if the functions are orthogonal (i.e. independent). Very often the wavelets are Gaussian in shape which allows for smooth, well behaved frequency characteristics. Also the Gaussian shape is a good default characteristic because it is a reasonable match for most functions. Incidentally, in this form the wavelet transform is very much like the original Gabor transform.

A summary of the essential points is shown below in figure 2.80(*a*).

It is undeniable that the space frequency functions are useful in spotting changes in statistics which result from machine tool and process problems. Can these functions be used in a functional way? The answer to this is yes with some provisos. Functions which involve flow can utilize space frequency function *S*.

| Data | Data modification | Multiplier | Result |
|---|---|---|---|
| $f(x)$ | | $\exp(-jwx)$ (continuous) | Fourier - F($w$) |
| $f(x)$ | $f\left(x + \dfrac{\tau}{2}\right) \cdot f\left(x - \dfrac{\tau}{2}\right)$ | $\exp(-j\bar{w}\tau)$ (continuous) | Wigner - W($\bar{w}\tau$) |
| $f(x)$ | | $\psi\left(\dfrac{x - b}{a}\right)$ (localized) | Wavelet - W($ab$) |

**Figure 2.80(*a*)** Space frequency functions.



**Figure 2.80(*b*)** Space frequency and the function map.

However, the space variable has to be replaced by a time variable. Where space/time frequency functions are not really useful is in functions which involve the 'normal' contact of surfaces. Here the contact mechanism is parallel rather than serial and is concerned with peak behaviour in a somewhat non-linear way. Using the classification given in chapter 7.

Can wavelets help with fractals [133]? Should a Markov process with exponential correlation function be described as a fractal surface or be left as it is? There is a danger of seeing fractal behaviour in all waveforms.

The basic idea behind wavelet use is to analyse the surface into different scales. Here 'scale' implies 'resolution'. In effect, wavelet analysis could be described as a 'mathematical zoom lens'.

According to Raja [132] the dyadically dilated wavelets constitute a bank of octave band pass filters and the dilated scaling functions form a low pass filter bank.



**Figure 2.81**

Depending on how many octaves are covered by the wavelet (e.g. three in figure 2 .82(a)) the surface spectrum can be successively approximated.

Usually the zero frequency response is taken to be zero because the wavelet has zero area. The use of 'raised' wavelets does allow a dc component. Also placing a Dirac delta function central and opposite to the wavelet can simulate the high pass version.

There is a real possibility of using wavelets to unveil specific features of manufacture or function. For example, take one of the most popular wavelets which is in the form of a hat of form

$$g(x) = (1 - x^2)\exp\left(-\frac{x^2}{2}\right)$$
(2.251)

This takes the form shown in figure 2.82(b)i.

It seems difficult to spot any manufacturing identity with this shape, so it would appear to be inappropriate to use this wavelet shape. However, rotating the hat through 180° about the $x$ axis produces the shape shown in figure 2.82(b)ii. This is precisely the shape of a single grain impression left on the surface by grinding and would have a very high correlation with this shape. It is true that autocorrelation would identify such a waveform, but in its correlation form, which is not as obvious the spatial form given in figure 2.82(b)ii.

*(a)*

Transmission

Band pass wavelet

$w$

Surface spectrum

$w$

$w/8$  $w/4$  $\dfrac{ws}{z}$

Low pass wavelet

**Figure 2.82(*a*)**

*(b)*

(i) Hat

(ii) Grain impression

$x$

$x$

**Figure 2.82(*b*)**  Hat transform.

The obvious path, therefore, would be to use the autocorrelation function of the surface profile to reveal the presence and horizontal scale of the grain impression. Knowing this, the hat shape could be used as a wavelet to find out exactly where the impression is strongest within the profile.

Without the pilot examination of the surface using autocorrelation the chances of picking a suitable wavelet shape and scale would be remote.

In fact this hat shape would not be used directly because it does not have the orthonormal properties required for the wavelet, so the wavelet method would probably be replaced with a cross correlation of the hat with the profile.

It would be necessary to find a wavelet set with the nearest shape and scale using one of the techniques used by Daubechies [134].

### 2.1.8  Fractals

Recently, it has been observed that some surfaces have a particular type of behaviour called the fractal property. Mandelbrot [52] described structures in nature such as snowflakes, coastlines, clouds, etc, as having the fractal property. Simply, this is that the geometry, regular or random, exhibits 'self-similarity' over all ranges

of scale. This means that these structures are characterized by some dimension which is not Euclidean. For example, a snowflake is not two dimensional; in fractal terms it is just over unity.

The fractal property of some engineering surfaces was first investigated by Sayles and Thomas [70].

There are two terms that are relevant in terms of scale. One is 'self-similarity' which refers to the property that it has the same statistics regardless of scale. In terms of roughness parameters therefore self-similarity implies that any parameter should be independent of the scale of sampling. The second, a 'self-affine' fractal, is only self-similar when scaled in one direction. Single-valued fractal functions must be self-affine because a small feature can only occur on a larger feature if its slope is larger. On this basis a surface profile should preferably be considered to be single-valued because surface slopes are steeper for finer sampling.

The difference between the self-similar property and the self-affine property is that self-similar fractals require only one parameter to define them — $D$, the fractal dimension — whereas self-affine surfaces require an additional one that relates to the scale being viewed.

In roughness terms it is necessary to specify, in addition to $D$, the way in which the ratio of vertical to horizontal magnification has to be changed in order to preserve self-similarity. The name given to this other parameter is called 'topothesy' after Berry [71].

Another way of classifying fractal surfaces is by the power spectral density. If the power spectrum is

$$P(\omega) = k/\omega^v \tag{2.252}$$

the dimension $D$ is given by

$$D = (5 - v)/2. \tag{2.253}$$

Yet another way is using the structure function $S(\sigma)$ where $\delta$ is the sample interval

$$S(\sigma) \propto \delta^{v-1}. \tag{2.254}$$

In these terms the topothesy A is the horizontal distance over which the chord joining the ends of the sample interval between measurements has an RMS slope of 1 radian:

$$S(\Lambda)/\Lambda^2 = 1. \tag{2.255}$$

The constant of proportionality in equation (2.196) is $\Lambda^{2D-2}$, so

$$S(\delta) = \Lambda^{2D-2}\delta^{2(2-D)} \tag{2.256}$$

This formula enables the dimension of a profile of a surface to be found. By plotting the logarithm of structure function against the logarithm of the sample interval a straight line should result, whose slope enables $D$ to be found. The intercept at unity gives $\Lambda$ (figure 2.83).

There is another way to evaluate the fractal dimension of a profile. This is using the high spot and local peak counts $m$ and $n$. Thus, given the structure function formula,

$$S(\delta) = \Lambda^{2D-2}\delta^{2(2-D)}, \tag{2.257}$$

the structure function can be evaluated in terms of the derivatives of the autocorrelation function.

$$S(\delta) = -A''(0)\delta^2. \tag{2.258}$$

**Figure 2.83** Determination of fractal properties.

Also the curvistructure function C($\delta$) is given

$$C\left(\delta\right) = A^{iv}\left(0\right)\delta^4 \tag{2.259}$$

which can be related to $S(\delta)$ by differentiating $S(S(\delta))$.

So equating these enables $D$ the fractal dimension to be obtained. For example, in terms of $m$, the local peak spacing

$$m = \frac{1}{2\pi}\sqrt{\frac{A^{iv}\left(0\right)}{-A''\left(0\right)}} = \frac{1}{2\pi\delta}\sqrt{4D^2 - 6D + 2} \tag{2.260}$$

Hence $D$ is found in terms of a simple count of peaks at the scale of size (which is usually determined by the instrument).

$\Lambda$, the topothesy, is less important but it is obtained from

$$\Lambda = \left(\pi R_q \delta^{D-1}\right)^{2/(2D-1)} \tag{2.261}$$

The topothesy has the dimension of length and can have any value, but the fractal dimension $D$ can only have values between 1 and 2; 1 corresponds to a straight line while 2 is an infinitely rough profile that fills in the whole plane of the profile. Hence the dimension $D$ is not very sensitive to different processes as is seen below in table 2.10.

**Table 2.10**

| Process | $D$ | Topothesy ($\mu$m) |
|---|---|---|
| Ground | 1.17 | $3.38 \times 10^{-4}$ |
| Turned | 1.18 | $2.74 \times 10^{-4}$ |
| Bead blast | 1.14 | $4.27 \times 10^{-4}$ |
| Spark eroded | 1.39 | $3.89 \times 10^{-1}$ |

In practice no real surface has fractal properties over all scales of size, only a very limited range, which means that there is no unique $D$ value for the surface any more than there is a unique value for other parameters.

The reason why there is interest in fractal behaviour is fundamentally to get out of the problem which has beset all tribologists and surface investigators for the past twenty years. This concerns the dependence of

the surface parameters on instrument characteristics which include evaluation length or area covered, resolution and sampling interval. It will be shown in the next chapter how non-intrinsic surface parameters like peak curvature and slopes depend heavily on the short and long wavelengths of the surface-measuring instrument. Obviously if a parameter can be found which is scale invariant then its value when measured within the bandwidth of the instrument would be perfectly valid for scales above and below the instrument's capability. Results could genuinely be extrapolated from one scale to another. The fractal dimension $D$ and topothesy $\Lambda$ are such parameters. It is therefore very tempting to find fractal dimensions in surfaces. These would be true intrinsic parameters.

The important question, therefore, is whether surfaces, and engineering surfaces in particular, are fractal. There seems to be a desire for them to be so, for the reasons given above, but is this justified? Obviously the question is mostly related to the method of manufacture of the surface. Fractal comes from fracture mechanics, so it would seem plausible to say that any process based on fracture mechanics would produce fractal surfaces. This may well be true for the machining of brittle components such as ceramics; it is not true for many processes in which plastic flow occurs.

Also, fracture mechanics is based on the laws of crack propagation, which follow a Brownian pattern. It would seem plausible again to suggest that any surface made up of a growth mechanism such as in the deposition of materials [73], solidification of a liquid [72] and fracture [75] (i.e. growth of cracks) would have fractal properties. Some fractal behaviour in stainless steel [138] has been observed.

Some of the basic properties of fractals in an engineering concept are attempted by Majundar and Bhusan [76].

In terms of characterization, if $z(x)$ is a profile of a fractal surface (everywhere continuous but not differentiable everywhere) it can be characterized as

$$z(x) = G^{(D-1)} \sum_{n=n_1}^{\infty} \frac{\cos(2\pi\gamma''x)}{\gamma^{(2-D)n}} \quad 1 < D < 2 \qquad \gamma > 1. \tag{2.262}$$

Here $G$ is a scaling constant related to topothesy by

$$\Lambda = G\big/(2D - 2)\sqrt{2\ln\gamma} \tag{2.263}$$

and $\gamma^n$ corresponds to a frequency term, that is reciprocal wavelength $\gamma^n = 1\big/\lambda^n$. This is very similar to expanding $z(x)$ in a Fourier series except that, instead of the harmonics increasing linearly with $n$ (i.e. $n\gamma$), they increase as a power law $\gamma^n$.

This is a case of the Weierstrasse — Mandelbrot function (W — M). It has interesting properties

$$z(\gamma x) = \gamma^{(2-D)}z(x) \tag{2.264}$$

which make $z(x)$ a self-affine function because the scaling of $z$ and $x$ is unequal; $n$ corresponds to the low cut-off frequency.

*(a) Fractal relation to power spectrum and its moments*

Since $z(x)$ comprises a superposition of infinite frequency modes it is a multiscale function. Although the series for $z(x)$ is convergent that of $\mathrm{d}z/\mathrm{d}x$ is divergent.

The relation to the power spectrum $P(\omega)$ is

$$P(\omega) = \frac{G^{2(D-1)}}{2\ln\gamma}\frac{1}{\omega^{(5-2D)}} \tag{2.265}$$

The structure function is

$$\langle [z(x_1) - z(x_2)]^2 \rangle = G^{2(D-1)} \Gamma(2D-3) \, \sin\!\left(\frac{(2D-3)\pi}{2}\right)(x_1 - x_2)^{(4-2D)}. \tag{2.266}$$

Also $A(\tau) = m_0(1 - \frac{1}{2}s(\tau))$

$$\langle z^2 \rangle = \sigma^2 = m_0 = \frac{G^{2(D-1)}}{2\ln\gamma} \frac{1}{(4-2D)} \left( \frac{1}{\omega_L^{(4-2D)}} - \frac{1}{\omega_h^{(4-2D)}} \right) \tag{2.267}$$

$$\left\langle \frac{dz(x)}{dx} \right\rangle = \sigma_m^2 = m_2 = \frac{G^{2(D-1)}}{2\ln\gamma} \frac{1}{(2D-2)} \left( \omega_h^{(2D-1)} - \omega_L^{(2D-2)} \right) \tag{2.268}$$

$$\left\langle \frac{d^2z(x)}{dx^2} \right\rangle = \sigma_c^2 = m_4 = \frac{G^{2(D-1)}}{2\ln\gamma} \frac{1}{2D} \left( \omega_h^{2D} - \omega_L^{2D} \right). \tag{2.269}$$

All these terms for curvature, slope and RMS are in terms of $D$ and are restricted by the instrument restraints $\omega_L$ and $\omega_h$ — the lower and higher bandwidth limits.

If the length of the sample and the resolution $\omega_h$ are known ($\omega_L$ is related to the evaluation length) then by measuring $m_0$, $m_2$ and $m_4$ given above, $G$ and $D$ can be determined and so the fractal dimension $D$ and the parameter $G$ which characterizes the roughness at all scales.

The equations above [76] show that the averages of the profile, slope and curvature are functions of the two length scales and their dependence involves the $D$ dimension of the surface. The variances of $z$, $dz/dx$ and $d^2z/dx^2$ in themselves give little information about $D$ and the multiscale structure of the surface.

The asperity density depends on $\alpha$:

$$\alpha = \frac{m_0 m_4}{m_2^2} = \frac{(D-1)^2}{D(2-D)} \left( \frac{\omega_h}{\omega_L} \right)^{(4-2D)} \tag{2.270}$$

It can therefore be seen that for all $D$ less than 2 the parameter $\alpha$ depends on $\omega_h$ and $\omega_L$ and is therefore instrument dependent and thus non-unique — the same old problem!

The scale dependence of surface parameters like the roughness RMS $R_q$ etc, are shown in table 2.11.

**Table 2.11**

| Parameter | Dependence on $\omega_h$ and $\omega_L$ |
|---|---|
| $R_q$ | $\omega_L^{(D-2)}$ |
| Mean peak height | $\omega_h^{(D-2)}$ |
| $\Delta_q$ | $\omega_h^{(D-1)}$ |
| Mean peak curvature | $\omega_h^{D}$ |

Majundar and Bhusan [76] suggest that the only way to remove this restriction of instrument dependence is to use the multiscale structure of self-affine asperities. This could be simulated by using the relationship connecting the profile ordinate $z(x)$ by

$$z = G^{(D-1)l(2-D)} \cos\!\left(\frac{2\pi x}{l}\right) \quad \text{where } -\frac{1}{2} < x < \frac{1}{2}. \tag{2.271}$$

If the dimension of a surface is taken to be $D = 2$ then the power spectrum behaves as $1/\omega^2$. If $\omega_h > \omega_n$ then the height varies as $\ln(\omega_h/\omega_L)$. If now the surface profile is multiplied by a length scaling factor $\beta > 1$, it increases the spatial resolution and decreases the sampling length and $\omega_n$ becomes $\beta\omega_n$ and $\omega_L$ becomes $\beta\omega_L$. The prediction is that the heights remain constant, the slopes increase as $\beta$ and the curvatures by $\beta^2$. The relationships follow this trend as can be seen in figure 2.84. How fractal dimensions occur in contact will be seen in chapter 7.



**Figure 2.84** Dependence of surface properties on scale factor $\beta$.

*In conclusion, the tendency to look for fractal behaviour is justified in the cases of the finer surfaces which tend to be manufactured by non-plastic processes and more by growth or deposition mechanisms.*

If it can be assumed that a $D$ value is common to all scales of size then this truly is an intrinsic parameter to the surface so that, in principle, measuring $D$ for any scale of size suffices for all. Hence if $D$ (and A) is estimated from one scale of size it may be common to all scales, but there is no guarantee as different process mechanisms may come into play. Measuring the standard parameters of height, slope and curvature definitely only relates to the scale of size in which they are measured, as will be seen in chapter 3.

By far the safest path is to find what scale of size is most important and measure the relevant parameter corresponding to that size. This method has to be the best whether fractal characteristics are assumed or not.

If fractal properties are assumed then the dimension parameter can be measured at the most convenient scale, usually corresponding to the most available instrument.

The properties of other scales can then be inferred directly and models of contact etc worked out for those of a scale of size that cannot be readily measured, perhaps because they are too small.

All the evidence is that for the finer surfaces there is more of a chance of having fractal properties than for the rougher. However, it could be that with the growing use of ductile grinding in which there is little or no fracture behaviour, fractal behaviour will be restricted to a few processes in which case there is no real benefit in pursuing the fractal path for general characterization. Also, at the longer wavelengths where the mechanisms are less violent and more continuous, in the sense that they are usually produced by variations in the tool path rather than the machining process, fractal behaviour is less likely.

The interesting point to note is that the fractal approach is so closely related to the order of the power spectrum that these two approaches could be considered to be different views of the same thing. In fact the value of $D = 1.5$ corresponds to the exponential autocorrelation function which is a natural consequence of the process having Poissonian statistics and would be the most likely shape of function for any truly random method of generation. It has already been noted by Mulvaney [77] that surfaces have spectra which are almost exponential anyway. He suggests

$$P(\omega) = \frac{K}{1 + \left(\omega/\omega_c\right)^v} \qquad (2.272(a))$$

rather than

$$P(\omega) = K/\omega^v. \tag{2.272(b)}$$

This fits in with the earlier observations on spectra by Whitehouse and Archard [28].

In conclusion, it seems that the fractal approach is in fact the spectral one with a different emphasis being put on the order of the power. The only real danger is that investigators try to impose fractal characteristics when they do not exist. The danger is that scale insensitive parameters (fractals) are used to try to characterize scale sensitive behaviour, e.g., dynamics [135].

It has been stated earlier that engineering surfaces such as milling, turning, etc. are highly unlikely to have fractal properties. One criterion which can be used to tell whether or not a process is likely to be fractal is simply to determine whether or not it is, or has, a growth mechanism (or decay) as part of its characteristics. Single point cutting obviously has not and nor does grinding, except when in the brittle cutting mode.

Any process involving sputtering, coating or even painting could well be fractal. The surface of thin films or doped semiconductor materials also have the deposited element present and so can be fractal.

In other words, conventional engineering processes are unlikely to be fractal despite some investigation [136]. However, with miniaturization, more specialized processes are being evolved often with bombardment methods; the result is that in the nanotechnology regime fractal processes are more likely.

Can friction and wear be characterized by fractal methods? How can fractal methods be applied to chaos? How sensitive are fractal methods to their definition? The last of these questions [139] has been discussed as a function of scale and the sensitivity of the fractal dimension with orientation of the surface profile. Brown [139] uses patchwork analysis to investigate roughness characterization and seems to go further with his straightforward method than many using more complicated techniques.

### 2.1.9    *Surface texture and non-linear dynamics in machines*

One of the ways in which surface texture is useful is in acting as a fingerprint of the manufacturing process and machine tool behaviour. This will be seen in chapter 6 where it is demonstrated that tool wear, built-up edge and tool vibration can be controlled by use of the correlation function and power spectrum. The surface reacts to the way in which the process is changing or vibration is built up in the spindle and traverse table.

One emerging way in which the link from the surface to the machine tool is being carried out is briefly mentioned below and in more detail in chapter 6. This concerns the use of non-linear dynamic system theory. Using this technique it should be possible to improve knowledge of the relationship between the machine tool and the surface (e.g. for turning or milling operations).

What makes the system non-linear is that for some part of the cutting cycle, especially in milling [79], the tool is not in contact with the workpiece. For turning this is not so evident. In milling the stability boundary may be two or three times wrongly calculated by conventional methods. This makes the understanding of such effects desirable — especially in acceptance testing of milling machines.

The idea is to build up a mathematical model of the machine tool dynamics from the observed data — that is, the surface roughness which has been produced [80]. This means that a state space has to be constructed from what is in effect a time series (i.e. the surface profile). This method has the potential of getting actual values of machine tool parameters.

The cutting model consists of a single non-linear oscillator coupled to a mass edge with the cutting attached to it. This gives

$$M\ddot{z} + T\dot{z} + \lambda z = F \sin c \tag{2.273}$$

The cutting force excites vibration $z$ in the oscillator of the form

$$\text{force} = bc_s c_t \tag{2.274}$$

where $b$ is the chip width, $c_s$ is the cutting stiffness and $c_t$ is the chip thickness

According to Scott [80] the equation of motion is

$$\ddot{z} + C\dot{z} + Dz = Bc_t \tag{2.275}$$

Normalizing the time domain for one rotation of the spindle is conventional and convenient. In equation (2.215)

$$C = \lambda/M, \; B = \sin(\alpha)bc_s/M, \; C = T/M. \tag{2.276}$$

Because the long-term behaviour best illustrates the non-linearity rather than the transient response — the system can then be steadied to reveal 'strange attractor' behaviour — a state-space representation of the system consists of $z$, $\dot{z}$ and the last profile (this means the profile produced in the last revolution).

The basic idea is that from this state-space representation using the model on the one hand and the surface on the other, a convergence can be achieved in which the parameters of the model and, thus, the cutting conditions, can be optimized.

Obviously the use of such a technique is in its infancy, but providing that the model is reasonable in the first place it should enable a much greater understanding of the relationship between the texture and the process.

The fundamental difference between this approach and the typology arguments given in much of this chapter is this: conventionally, the characterization of the surface profile or areal data has resulted in a set of numbers which preserve as much of the real information in the data as possible, consistent with the use to which it will be put; these values are then used to help pin down the control of manufacture or function. In the state-space approach all the profile data are used — the convergence to a pattern which reduces to numbers, allowing the actual parameters of the process itself to be estimated in value. This latter method must be better in that it is a more direct link but it does have recourse, at least in the first instance, to a realistic model of the process and the machine tool.

## 2.2 Waviness

This is an integral part of the surface texture. There is a considerable controversy as to whether waviness should be included with the roughness evaluation or not. Probably more has been written on methods of separating the two than on any other subject in surface metrology. Some investigators think that it has all been a waste of time and that, functionally, the geometry as a whole is important. That this is undoubtedly true in a number of function applications, such as light scattering, is a valid point. However, there are a number of reasons why waviness should be measured apart from the roughness; the principal one being that it represents an important symptom of machine tool behaviour. From the point of view of manufacturing control, its measurement is most valuable. On the other hand, there are functional uses where roughness is much more important and should not be corrupted with irrelevant information, static contact being a typical example. Anyway, the fact that they may be assessed separately should not preclude their integration when required. What is more questionable is the wavelength at which the separation between roughness and waviness is supposed to occur. This obviously changes from process to process and even within a process. Problems such as this have ensured that there has been a considerable metrological gap in the measurement and understanding of waviness.

Here the two basic metrological problems, that of measurement and that of providing a numerical index, are retarded, mainly because it is only recently that any real functional significance has been attached to waviness. Because of this inheritance of neglect any attempt to measure it has been with either form-measuring

instruments or surface-roughness-measuring instruments. The former have generally been too insensitive while the latter have had insufficient range. Fortunately this state of affairs has changed rapidly, most instrument manufacturers now providing suitable instrumentation. Typical lengths of traverse are of the order of 10 mm or more and with sensitivities of the order of a tenth of a micrometre.

It is in the horizontal properties rather than in the vertical where the main differences lie. The heights of waves need not be and indeed often are less than that of the peak-to-valley roughness but the horizontal spacings are usually much larger than those of the roughness, often an order of magnitude or more. Therefore, methods of separation are usually based on some kind of wavelength discrimination. At the other end of the scale it is often stated [81] that unless at least three successive waves exist on the surface then no waviness can be measured.

It is difficult to isolate the measurement of waviness from the reference lines used in roughness because, in some instances, they are one and the same thing. Whereas waviness is often regarded as no more than a demarcation between roughness and errors of form, the roughness reference line also has an additional job of segregating roughness from the set-up effect of the instrument. In surface metrology there is difficulty in isolating intrinsic references (those derived from the surface itself) from the references used for convenience which are derived from a profile and which themselves may be tainted with errors due to the set-up of the part relative to the measuring instrument.

This discrepancy between the waviness profile and the roughness reference is most severe when using filters, because these operate not only on the profile (the boundary of metal to air itself) but also on its position and angle in space. It is only when the filter has had time (or space) to find this steady-state position and angle that the true geometry can be assessed. So it is better to talk about waviness after this situation has been reached. It has been common practice to regard waviness as the reference line for roughness, which only works if the filtering has been done properly using a phase converted filter.

Because of the inherent difficulties involved, waviness will be defined more by familiarity rather than formally. A number of aspects will be considered, some historical and some theoretical. In this way, the basic problems will be brought into the open even if only indirectly.

A good starting point is to consider some of the features of waviness (sometimes called secondary texture) and the question of what is meant by waviness, or if in fact it exists or needs to be measured at all. These questions depend crudely on three things:

(1) what is common usage;
(2) whether purely causative arguments can be used; or
(3) whether functional effects are the real criterion?

Unfortunately the answers to these problems not only are unknown, they are largely mixed up.

For this reason what will be discussed will be to some extent disjointed. However, the historical definition will be introduced here that 'waviness is some sort of amplitude effect induced by the machine tool, normally to the detriment of the performance of the workpiece'. It is usually felt that waviness is a result of poor manufacture, that is, chatter due to a machine tool which is too elastic. Also, it is not always detrimental, as will be seen in chapter 7. The problem is that if it is produced by poor manufacture it cannot be properly controlled anyway! If it is also causative (i.e. produced by the machine tool) then it imposes not only a distinctive geometry on the surface but also subsurface properties. These in turn have different functional effects additional to the subsurface effects produced by roughness. On balance, therefore, it seems prudent to try to evaluate them separately.

There is another reason for separating the components of the surface out rather than leaving them together. The reason for this is that the process marks on the surface are not important just for their geometry. They are important in the subsurface sense.

The reasoning is thus. Process marks — the machined marks — are produced at very high speeds, producing very hot spots on the surface which seriously stress the subsurface. Residual stress is severe and is

concentrated close to the surface, e.g. $1\mu m$. Plastic movement also introduces stresses. On the other hand waviness due to the machine tool is produced slowly and elastically at quite large Hertzian depths (fig. 2.85) — ten times the surface roughness.



**Figure 2.85** Subsurface stresses.

A graph showing the relative averages of roughness and waviness is shown in figure 2.86 remembering that the power spectrum $P(w)$ is not the power.



**Figure 2.86** Distribution of energy.

From figure 2.86 it is seen that far more energy goes into the surface skin via the roughness than via waviness and form simply because the energy is proportional to the square of the frequency.

It can be said that the geometrical profile of a surface not only gives roughness and waviness geometry but also gives an indirect stress picture of the subsurface. Very sharp curves on the surface are produced with high energies hence high temperatures and produce abusive conditions in the subsurface.

The important point is that if there are two surfaces involved, such as in contact, it is not possible to give equal weight to all components of the geometric spectrum. Each one is a reflection of the energies which produced it and each one represents different mechanical properties. So roughness and waviness should be separated and dealt with in a way compatible with their thermal pedigree. If the function does not involve two surfaces, as in optical scatter, by all means treat them as one geometry but not otherwise.

One extra problem sometimes encountered is whether the waviness which is visible on the workpiece actually exists! Sometimes, if a number of summits spaced widely apart in the $x$ and $y$ directions are slightly burnished the eye can make up very convincing patterns. The low resolution of the eye tends to pull detail together.

In any event the patterns formed on the surface by waviness can be varied. Some, like pronounced chatter marks and coarse feed marks of a badly trued grinding wheel, can be identified at a glance; others may need

an instrument to reveal their presence; and some cannot be measured. In the case of surfaces of revolution, those extending along the lay of the roughness often become the circumferential departures from roundness.

Waviness as seen in a profile graph can often be appraised both as an undulation of the mean line and as an undulation of a line drawn through the more prominent crests. It is thus possible to speak both of mean line waviness and of crest line waviness. Generally the two are not the same, although they can sometimes be much alike.

One of the few papers dealing exclusively with waviness devotes much space to this problem [81].

Whichever of these two methods of definition is functionally the more significant is open to question. Another problem is finding the real nature of waviness as opposed to form error and roughness. It is often asserted that the constituents of the geometry making up the surface are additive. Certainly in the case of waviness this is open to doubt; it depends on whether waviness is defined as a band of undulations within the total geometry. Typically the presence of waviness is detectable visually as a pattern of marks spread more or less periodically along the surface. This sort of effect on the eye can be produced by different types of wave-form. Three idealized examples are shown in figure 2.87.



Figure 2.87 Different modulated waviness.

This figure shows the single case often encountered in which the waviness is superimposed on the roughness. It also shows a waveform having the form of an amplitude-modulated signal, which means that there must exist multiplicative components in the geometry, and it shows the case where the phase or frequency of the tool mark is changing with passage along the surface. Each of these geometries can cause a visual effect, and certainly (*a*), (*b*) and possibly (*c*) could also be caused by faulty machining, particularly when the machining has been done by numerically controlled machines. At present little or no provision is made for differentiating between these cases but it can be done using the Wigner function. It may well be that more than one of these effects are present on any one surface. The shape of the waviness component depends to some extent upon the type of reference used to isolate it from the roughness. For example, the crest line drawn through figure (*a*) follows almost exactly the mean line, as it does in figure 2.87(*c*). However, the mean line is straight for figure (*b*) which is certainly not true for the crest or valley line.

More practical cases are shown in figure 2.88 which shows both geometric and visual pictures.

In one case (face turning with a loose cross-slide) there is an undulation which is many times higher than the tool scratch marks, although it is hardly detectable by eye, while in the other (straight flat grinding), what might be taken for pronounced chatter marks are found to have a height hardly greater than that of the normal scratch marks. The illusion of deep chatter marks here is caused by a periodic change in the direction of the scratch marks without much change in their general level, this phenomenon being quite common.

Obvious differences in waviness between just these two ways of representing the surface highlight one difficulty of measurement.

It seems more than likely that the selected definition of waviness will depend upon the function of the surface. This is because all measurements of waviness take place on the surface before it is used. What is really essential in some instances, such as wear, is that the waviness of the run-in profile is measured, even if it is by prediction from the unworn profile — under these circumstances the waviness determined from a mean line may well be the most significant. On the other hand it seems equally likely that in the problems of limits and



**Figure 2.88** Graphs of real surfaces showing different errors in combination.

fits a crest line may be the best reference to use simply because, under these conditions, the initial contact points are not removed because there is little load or movement.

Work based on these different methods of defining the reference has been carried out by a number of workers and various definitions have already been introduced into the national standards of some countries.

The first task before putting a number on waviness is to get a waviness profile. In surface roughness the instrument itself did this. The profile represented an air-metal boundary. Fitting parameters to this boundary, assuming a reference line to have been suitably defined, was justifiable. Problems arising from a less than perfect reference could be compensated for afterwards. This is not so easy in waviness. It is not as a rule a boundary — it is a profile derived from a profile. There seem to be more degrees of freedom to be constrained than in roughness.

The way in which the waviness profile behaves on average for random waveforms has been evaluated [84]. This will be discussed in the chapter on instrumentation. However, it does show that again both the height $h$ and the spacings of asperities $\lambda$ are involved;

$$\text{average waviness} \sim \left[ 2\ln\left(\frac{8Rh}{\lambda^2\sqrt{2\pi}}\right) \right]^{1/2}. \tag{2.277}$$

In each of the above cases the standard roughness filter is compared in figure 2.31 with the so-called phase-corrected filter. Notice how the phase-corrected filter follows the line that the eye would probably judge to be the waviness.

Figure 2.31 also shows how in some severe practical cases the phase-corrected mean line of roughness looks like a convincing waviness profile. It is still true to say, however, that such distortion due to the standard wave filter does not always occur. It is only when microcontact and macrocontact are being investigated that these effects are better minimized.

Yet other methods exist, one even more open ended than the others. This is based simply on experience of the process. If an inspector can detect poor machining by eye from the surface or the profile graph it should be possible to measure it by taking advantage of the versatility of a computer. Once the geometrical outcome of a fault has been identified on a graph an algorithm can be written to simulate it. This method suffers from the disadvantage that, for each machine, process and material part, the algorithm might have to be changed. There is evidence, however, that this new approach is being used more often in industry. Unfortunately computer versatility sometimes creates more problems than it can resolve.

The ways in which the waviness profile can be obtained in practice will be briefly discussed in chapter 4. In many cases the fitting of the reference line to roughness has many common features with the waviness profile so that much overlap is possible. In those situations where both waviness and roughness are required at the same time, it is becoming apparent that the use of a computer coupled to the instrument considerably reduces the measurement problems.

Another case in which waviness can be best isolated by a computer is when random process theory is being used. In this technique it is the generally periodic nature of the waviness which is relied upon as a basis for separation rather than the different frequency bandwidth.

As has been explained in section 2.1, the autocorrelation function can identify periodic from random components providing that certain conditions are met; the principle being that there must be a clear distinction as to whether or not phase-modulated periodicity can or cannot be tolerated.

Leroy [85] has attempted to isolate the waves on the surface by an iterative process. He identifies the dominant periodicity first by a novel level-crossing procedure similar to the motif method. Then he assumes first estimates of the values of the amplitude and phase from that of the nearest subharmonic of the length of the chart. This enables some sinusoid to be drawn on the profile. From the differences between it and the profile small adjustments are made until the maximum deviation is a minimum, or the least-squares deviation is a minimum. Once this sinusoid is optimized to a predetermined degree it is removed from the profile and the procedure repeated on the next sinusoid. The whole operation is cycled until the profile has been broken

down into a sufficient number of sinusoids. These are not necessarily orthogonal sinusoids. A similar result can be obtained by breaking down the profile length into its harmonics. Providing that the profile length is large compared with the longest surface waveform, the spectrum is generally satisfactory.

Methods of assessing waviness have sometimes been carried out using a twofold attack. This involves determining the presence of waviness and measuring the amount of waviness, preferably at the same time. Some investigators have merely concentrated on developing a waviness profile, others on the assessment. In general, though, the numerical assessment of waviness is still being argued about.

There is also the problem of what to measure once the profile of the waviness has been decided, whether based on a mean line method or an envelope.

One such approach based on the theory of crossovers developed by Rice and others has been advocated by Peklenik [86]. He assumed that the waviness takes the form of a second-order random process (which it quite often does). He then determined the superior and inferior envelopes on a bandwidth criterion. Typical graphs are shown in figure 2.89.



Surface profile $X(t)$     Envelope $B(t)$ Mean line $m_x$

(a)

(b)

**Figure 2.89** Concept of surface envelope: (a) ground surface, (b) turned surface.

Peklenik [86] uses the mean height of the envelope as a starting point for assessing the usefulness of the waviness concept in envelope terms which, when worked out using the theory of narrow-band processes, gives a value of $a$ above the mean line of the profile, where $a$ is the RMS value of the surface. This established a level above which it is considered that functional effects may well be concentrated.

From these two parameters are worked out the average width of the profile above this level and the time $T$ of one period of the envelope. Thus the concept of the dominant wavelength of the waviness emerges. This has also been advocated by Spragg and Whitehouse [13] as an extension of their average wavelength concept from that of surface roughness to waviness.

Peklenik further suggests that the rate of change of the envelope and phase angle are important, a point which emerges later with the Wigner distribution function.

Thus, the probability density of the envelope, slope and phase change have, subject to the assumptions above for an amplitude-modulated profile, the following Gaussian form

$$p(\dot{\varepsilon}) = \frac{1}{\sigma\sqrt{\omega_2^2 - \omega_1^2}.\sqrt{2\pi}} \exp\left(\frac{-\dot{\varepsilon}^2}{2\sigma^2(\omega_2^2 - \omega_1^2)}\right) \tag{2.278}$$

$$p(\dot{\varphi}) = \frac{\omega_2^2 - \omega_1^2}{2\left[(\dot{\varphi}-\omega_1)^2 + (\omega_2^2 - \omega_1^2)\right]^{3/2}} \tag{2.279}$$

Further possibilities include using the correlation function of the envelope and the correlation function of the phase as the parameters of waviness characterization. From these investigations the most probable useful parameter for the narrow-band Gaussian model appears to be

$$\alpha_\varepsilon = 2M_{\lambda\varepsilon}/T_1. \tag{2.280}$$

This expresses the estimated average peak length of solid material $M_{\lambda\varepsilon}$ related to half the envelope wavelength $T_{1/2}$ [86]. But a problem arises with this suggestion when the waviness has no amplitude-modulating factor: frequency modulation can result from tool vibration. Under these circumstances the waviness characterization must lie only in the phase terms.

Usually when waviness amplitude has been specified in the past it has been in terms of a measure of the separation of the envelope maxima and minima, in some cases this measured value being an average of a number of measurements. The measurement of waviness according to Perthen (reference [10] in chapter 4) and Von Weingraber (reference [8] in chapter 4) has been the average drop between the straight line representing the form error and the envelope line established by a rolling circle having a 3.2mm radius (figure 2.90). One advantage of this method is its consistency, the roughness and waviness numerical values being additive. This could usefully be called mechanical waviness. However, the original choice by Perthen and Von Weingraber of 3.2 mm is seen to be too small in figure 2.90 because the transmission characteristic for waviness never approaches the value of unity.



**Figure 2.90** Problem of roughness, waviness and form — transmission characteristics.

The mechanical methods of detecting waviness by means of a stylus, which is a blunt foot that rides the roughness, is well known as the basis of the E system. It is very straightforward for profiles in which roughness and waviness signals are additive (figure 2.89). Also, even-amplitude modulation as in figure 2.89 can easily be picked up, whereas other forms of waviness, such as a 'chirp' signal or frequency modulators, hardly make any impression on a mechanical system. The Wigner method finds these anyway, as will be seen next. Note that this disadvantage should not have too high a weighting because some waviness effects are additive to the roughness signal. This could be caused, for example, by the slideway ball screw having an error in pitch.

Modulation effects in waviness are usually caused by the tool vibrating radially relative to the component or axially due to tool column stiffness problems (figure 2.89) or self-generated chatter between the tool and the workpiece. It should be added here that these causes of waviness are not likely to introduce fractal properties to the envelope because by definition they are much slower moving than the effects that cause the roughness.

**Figure 2.91** Mechanical waviness and form.

Obviously in the M system the RMS values of the waviness (defined as the mean line of the roughness) and the roughness are also additive whether or not the filters used are phase corrected. However, the advantage of using a phase-corrected method is that average values also become additive, that is the transmission characteristics of the waviness line and the roughness always add up to unity, whereas if phase shift is ignored this is not so. For simple waveforms this is illustrated in figure 2.92.



**Figure 2.92** Typical waviness on surfaces.

A better method is to use a phase-corrected filter whose characteristic intersects the roughness wave filter at 50% (figure 2.93).



**Figure 2.93** Relationship between waviness and roughness amplitude for phase-corrected filter.

In such cases, the two components are exactly complementary in phase and in amplitude. The fact that the waviness reference includes the tilt, misalignment of the surface and form error cannot be avoided.

In practice, the filter characteristics of waviness (i.e. those characteristics excluding roughness) are $W_e(\lambda) = 1 - R(\lambda)$.

One of the issues is the slope and shape of the curve. Linear characteristics are to be preferred even if only to be consistent with the fact that the phase characteristics are linear in frequency. However, current thinking favours the Gaussian shape.

So, one criterion for waviness has been that it is deterministic, for example has a known geometrical form. The roughness on the other hand is deemed to be random. Obviously this idea is simplistic because of the large number of cases in which either the process is also deterministic or the waviness is random albeit of long wavelength.

Hence it can be seen that use of the local frequency moments can isolate the various forms that envelopes can take. This is most likely to be important in functional cases.

It seems, therefore, that more discriminating functions like the Wigner distribution function and the ambiguity function can help in quantifying the waviness.

## 2.3   Errors of form

### 2.3.1   Introduction

One of the biggest problems with waviness is that it is difficult to deal with from a purely geometrical standpoint. It is not a measure of a metal-air boundary as is roughness. Neither is it a deviation from a perfect Euclidean shape as are the measurements of the deviations from straightness, flatness and roundness. At least in these latter cases there is some defined perfection. (The definition can be written down formally.) This makes the measurement problem easier. On the other hand the wide range of different engineering shapes that have to be contended with is considerable. Some general classification will be discussed in what follows. Figure 2.94 shows a breakdown of relevance.

To be consistent with what has already been said deviations from, or concerned with, a linear causative variable will be considered. Form can be considered to extend into straightness and then flatness, after which different goemetric variables such as circles, epi-trochoids, etc, have to be examined.

The only practical way is to relate the geometry to cause and use, as in figure 2.94.

There are, as in most metrological problems, three separate difficulties which need to be considered. The first is the nature of the problem, the second is the method of assessment and the third is the method of display.



**Figure 2.94** Quality regimes for surface geometry.

Table 2.12 gives display symbols for the various components of form error.

**Table 2.12** Error of form symbols.

| Symbol | Name | Symbol | Name |
|---|---|---|---|
| ◯ | Roundness | ▬ | Straightness |
| ◉ | Concentricity | \ / | Co-axiality |
| ‖ | Parallelism | ◖◯ | Cylindricity |
| ▱ | Flatness | ╱ | Runout |
| ⊥ | Squareness | | |

Errors of form suffer from the opposite of roughness. In roughness the short wavelengths are known subject to instrumental constraints such as stylus tip dimension. It is the long-wavelength boundary of waviness which is difficult to define. In form errors the long wavelengths are determined by the ideal shape specified, for example a circle, and it is the short-wavelength boundary with waviness that has to be specified. Because of the difficulty of defining anything in absolute terms, in waviness it is usual to define the short wavelength of form in terms of the size of the piece. It is usually defined as a fraction of the relevant dimension of the workpiece. A factor of one-third or a quarter has been used. Wavelengths less than this are difficult to explain in terms of errors of form.

Errors of form are relatively easy to characterize: they are broken down into Euclidean shapes such as circles, planes, etc. This is easy when compared with the problem of characterizing roughness, especially now that the surface characterization using fractals has started. However, complex problems arise in the methods of assessing form error. In what follows emphasis will therefore be placed on the assessment problem.

The ideal form itself can be regarded as a skin in space. It needs a certain minimum number of points to describe its shape and position.

The ideal skin so defined is infinitesimally thin but practical surfaces are not. If zonal methods are being used in assessment, such as a minimum zone in which the thickness is a measure of the minimum peak-to-valley distance from the skin, one more data point is needed to fix the thickness of the skin. That is, the minimum zone sphere needs five points of constraints corresponding to three points for origin, one for size (a radius) and one for thickness of zone (see table 2.13).

Errors of straightness are often due to errors in machining and slideway error, but they can also be the result of sagging of the workpiece under its own weight, thermal effects produced during machining, stress relief after machining and many other reasons.

The types of component usually involved in this sort of assessment are shafts, slideways, etc.

There are, as in most metrological problems, three separate difficulties which really need to be considered. The first is the nature of the problem, the second is the method of assessment and the third is the method of display of the very small deviations.

A more comprehensive list is shown in table 2.14.

The points in the column do not refer to degrees of freedom: for example one degress of freedom is a translation in the 2 direction is a movement along a line which needs two points to establish it.

Also a rotation about one axis is one degree of freedom yet it requires 3 points to establish the circle. The points are geometric constraints imposed by the shape. Knowledge of these points enables unambigous paradigms to be developed.

**Table 2.13**

| Feature | No of points |
|---|---|
| Line | 2 |
| Plane | 3 |
| Circle | 3 |
| Sphere | 4 |
| Cylinder | 5 |
| Cone | 6 |

**Table 2.14**  Zonal points

| Function 2D | Figure | Points for definition |
|---|---|---|
| Line |  | 2 |
| Minimum deviation from line |  | 3 |
| Plane |  | 3 |
| Minimum Deviation from plane |  | 4 |
| Circle |  | 3 |
| Maximum deviation from minimum circumscribed circle (ring gauge) |  | 4 |
| Maximum Deviation from maximum inscribed circle (plug gauge) |  | 4 |
| Manimum zone |  | 4 |

| Function 3D | Figure | Points | Deviations from |
|---|---|---|---|
| Sphere |  | 4 | 5 |
| Cylinder |  | 4 | 5 |
| Cone |  | 5 | 6 |

In addition to these three obvious considerations, there are some practical points that are also relevant. These are basically concerned with the way in which the data is gathered. This step is vital because upon it often depends the way of classifying the data. Obviously the more data points the better, providing the errors are reduced rather than increased by getting them.

Ideally there should be as small a number of data points as possible, preferably taken at those points which are known to be functionally sensitive. Also the sample pattern for obtaining the data should be picked so that the points have equal uncertainty. The other issue is matching the data pattern to the shape of the workpiece. Rather than consider these practical points separately they will be dealt with as they occur. Also, as in the case of roughness, measurements to get the raw data and characterization of the data are difficult to separate: indeed they must often be carried out together.

### 2.3.2    Straightness and related topics

It should be noted that the terms straightness and roundness are somewhat misleading. They should really be 'the departures from true roundness' and for straightness 'departures from a straight line'. Common usage has reduced then to the single word.

In measuring straightness it is common practice to measure the error by methods that take account of the length of the surface to be measured. For example, if a surface plate is to be measured for flatness, the instrument chosen would be such that it would ignore scratches and scraping marks and detect only the longer wavelengths on the surface. However, if the workpiece is the spindle of a watch component only a few millimetres long, then a different approach is required, since the spacings of the undulations which could be classed as errors of form are now very short and it becomes much more difficult to separate form error from the surface texture.

Because the measurement of the straightness of small components to a high order of accuracy is fairly easily achieved by the use of instruments having accurate slideways, the graphical representation is usually satisfactorily displayed on a rectilinear chart.

In using such instruments the parallelism or taper of opposite sides of a bore or shaft can be made simply by taking multiple traces on the recorder. Note that this is correct providing that there is no relative movement between the workpiece and the datum when taking the traces.

For the assessment of independent bores there is an added complication. Repositioning of the part relative to the datum is usually necessary in order to fit the pick-up into the different holes. Accurate relocation methods need to be used in this situation. Figure 2.95 shows a typical example of the parallellism of surface generators.

As has been stated, it is difficult to isolate the method of assessment of a parameter from the method of measurement. Straightness measurement can be achieved in a number of ways, either directly or indirectly. The direct method involves comparing the surface with a straight line datum, which can be either a mechanical reference or the line of sight of a telescope, the reference being the line, the test piece being the surface. Alternatively the instantaneous slope of the surface or its curvature can be measured and the straightness



**Figure 2.95**  Worn cylinder bore generators.

obtained by integration. As a first step in characterization the data has to be provided. Then the characterization proper can take place usually by fitting a suitable function to the graph.

All of these methods ultimately must end up with a graph that shows the deviations of the surface from a straight line reference somewhere in space as in figure 2.96. In general the reference is not in the same direction as the surface, neither is it within the measured profile.



**Figure 2.96** Measurement of straightness using autocollimator of level.

The value of $a_i$ would be a direct measurement (figure 2.96), or it could be given by

$$a_i = \Delta x \sum_{j=1}^{i} \alpha_j \qquad (2.281)$$



**Figure 2.97** Measurement of curvature using multiprobes.

If the measurements are made with an autocollimator or a level which measures angle rather than deviation [98]. This method has the advantage that the arbitrary value $b$ is removed, but the problem of the arbitrary tilt shown in figure 2.96 is left. However, this can be removed if instead of angle the curvature or higher orders of deviation are measured. This can be achieved by the method shown in figure 2.97.

In the example of straightness there may be variations in the linear distance between the part and the datum, there may be variations in the tilt between the two and there may even be relative curvature variations. Such a situation involving all of these might arise in measuring the form of a moving flexible sheet of material, for example.

### 2.3.3    *Generalized probe configurations for variable errors*

Because of the diversity of problems of this nature met within metrology it is necessary to identify a more generalized configuration. Take the variations mentioned above, for instance: the distance $z$ between the datum and the test object can be expressed as a random variable in time. Thus

$$z = d(t) + m(t)x + c(t)x^2 \qquad (2.282)$$

where $d, m, c$ are independent random variables representing average separation, tilt and curvature respectively; $x$ is distance. Those coefficients above the second can be regarded as a characterization of errors of straightness.

To eliminate this number of random variables four probes are needed in the general case [88]. The probes would have sensitivities of 1, $a$, $b$, $c$ and they would be at distances $l_1$, $l_2$, $l_3$ and $l_4$ from the centre of the measuring system in figure 2.97. Three equations need to be satisfied:

$$1 + a + b + c = 0 \qquad (2.283)$$

$$-l_1 - al_2 + bl_3 + cl_4 = 0 \qquad (2.284)$$

$$l_1^2 + al_2^2 + bl_3^2 + cl_4^2 = 0. \qquad (2.285)$$

Equation (2.283) takes account of average separation, (2.284) the tilt and (2.285) the quadratic terms. The tilt term has to have odd symmetry about the mid-point of the carriage — hence the negative signs in the equation. Obviously just the same order of systematic error can be determined by having multishifts instead of multiprobes. Solving these equations reveals a very useful fact that formulae for numerical differentiation can be used in multiprobe technology. For example, consider the case where $l_1 = -2h$, $l_2 = -h$, $l_3 = +h$ and $l_4 = +2h$. The probe combination signal to satisfy equations (2.283–285) with these constraints imposed is given by $C$ where

$$C = V_1 - 2V_2 + 2V_3 - V_4 \qquad (2.286)$$

and $a = -2$, $n = 2$, $c = -1$.

Notice that equation (2.286) corresponds to the simple formula for the third numerical differential where the measured ordinates $f_1, f_2$, etc, have been replaced by probes. Thus,

$$h^3 f''' = \tfrac{1}{2}(f_{+2} - 2f_1 + 2f_{-1} - f_{-2}). \qquad (2.287)$$

In this case there is a gap of $2h$ between $V_2$ and $V_3$. $h$ is the unit of spacing between the probes. Making all the probes equidistant gives $a = -3$, $b = 3$, $c = -1$, which reduces the overall distance from $4h$ to $3h$.

Such a system has a harmonic weighting function W given by

$$W_{\mathrm{h}} = [\exp(-\mathrm{j}2\pi l_1/\lambda) + a\exp(-\mathrm{j}2\pi l_2/\lambda) + b\exp(\mathrm{j}2|l_3/\lambda) + c\exp(-\mathrm{j}2\pi l_4/\lambda)]. \qquad (2.288)$$

$W_{\mathrm{h}}$ will be explained with respect to similar methods in the measurement of roundness and refers to the harmonic distortion introduced by the method and which, because it is known, can be removed easily by computer [88]. Techniques like this are becoming more popular because computation allows errors and tedium to be removed from the calculation. This trend is a progression from the earlier methods of error reversal used to get rid of errors in the straight-edge used as a reference [100].

### 2.3.4  Assessments and classification

There is one serious problem with the indirect methods as opposed to the direct method which to some extent offsets their benefits. This is the enhancement of any noise. These methods are basically differentiation

methods which reduce the signal-to-noise ratio and, whilst this is not too serious a problem in the measurement of straightness, it can be very serious in measuring flatness. Methods have to be adopted to reduce noise. This is a good example of the balance which has to be achieved between ease of measurement and complexity in processing. Luckily the latter is becoming less expensive and faster so that clever spatial configurations for probes can be utilized more often.

Assuming that such a data set as seen in figure 2.96 has been obtained by the appropriate means, the problem of characterizing it remains, that is what it means in terms of departures from straightness.

A proposed way has been first the use of a best-fit least-squares line drawn through the profile to establish the reference (figure 2.98). Then from this line the sum of the maximum positive and negative errors is taken as a measure of the departure from true straightness.

Deriving such a line is equivalent to fitting a first-order regression line through the data set representing the profile of the workpiece, as in fitting reference lines to separate waviness from roughness, although here the subtle difference is that the drawn line is an end in itself. Deviations from it will be the deviations from the intended shape. Let the vertical deviation from an arbitrary level at distance $x$ be $z$. Also let the equation of the least-squares line to the data be $Z = C + mx$.

It is required to find the coefficients $C$ and $m$ in such a way that the sum of the differences squared between $Z$ and $z$ for all $x$ is a minimum (figure 2.98):

$$S = \int (z - Z)^2 (\cos^2 \alpha). \tag{2.289}$$



**Figure 2.98** Straightness — best-fit line.

Using exactly the same technique as for equation (2.16) but only using the digital equivalent

$$2\alpha = \tan^{-1} 2 \left( \frac{\sum z_i x_i - (1/N) \sum_{i=1}^{N} x_i \sum_{i=1}^{N} z_i}{\sum x_i^2 - (1/N)(\sum x_i)^2 - \sum_{i=1}^{N} z_i^2 + (1/N)(\sum z_i)^2} \right); \tag{2.290}$$

because the angles are small this reduces to

$$m = \frac{\sum x \sum z - N \sum xz}{(\sum x)^2 - N \sum x^2}. \tag{2.291}$$

However, these small angles will not be valid for microminiature parts, as they are evaluated with small-range coordinate-measuring machines, and so

$$C = \bar{z} - mx \quad \text{where } \bar{z} = \frac{1}{N} \sum_{i=1}^{N} x_i. \tag{2.292}$$

The process of estimating $m$ and $C$, the tangent and the intercept, can be regarded as an application of the principle of maximum likelihood assuming $N$ sets of random independent variables.

The calculation of the flatness and straightness deviations using formulae like (2.291) instead of (2.292) is getting less justifiable as instruments become more integrated. Precision and miniature coordinate-measuring machines will soon be available to measure straightness etc as well as size and there will be no guarantee that the general direction of measurement will be close to the plane or line direction being sought. This is a very important consideration for the future.

A very simple method of getting an estimate of the slope often used in practice is that of merely joining the first and last data points together with a line. Another method consists of defining a zone comprising two separated parallel lines which just contain the straightness data between them.

A few points need to be noted. One is that such a zonal method needs three points of contact. It is easy to show that if there are only two points of contact, one at a peak and one at a valley, then the separation of parallel lines through them is not a minimum. Three points of constraint are needed, which could be two peaks and one valley or two valleys and one peak. It can be shown that in situations involving the minimization of linear problems the maxima and minima will always be alternate [90]. This shows the principle of the Steifel exchange technique of linear programming. Three unknown parameters, $m$, $c$, to establish one line and $E_s$ to establish the separation of the reference line and the minimum case of a straightness minimum zone, and three contact points are needed (figure 2.99).



Straightness error

**Figure 2.99**  Straightness — Steifel exchange.

How this zonal method differs from those used in roundness will be discussed in the relevant section on roundness.

Usually, but not necessarily, the zonal method described here will give a smaller actual value of the straightness than the least-squares method. Obviously by definition it cannot ever be bigger because it is the minimum value. Differences between the two methods usually do not exceed 10% or thereabouts. Similar considerations apply to flatness. Also the obvious problem with zonal methods is their susceptibility to the odd very large peak or valley. This is the same problem as that which occurred in the evaluation of the surface texture peak parameter. The danger is that a large peak or valley may not be typical of the process.

Numerical techniques for evaluating the minimum zone will be given in section.

It could be argued that alignment should be encompassed in the subject of straightness because it is the measurement of discrete points from a line. The basic technique is similar to that of straightness measurement but the art is in its application to a multiplicity of engineering components such as large frames, bearings and plant. These details are covered adequately elsewhere [87].

Errors in gears, screw threads, etc, are specialized subjects dealt with in engineering metrology books and will not be covered here. Needless to say the same problems always exist: that of establishing a reference and assessing the errors from it.

*2.3.5    Flatness*

Flat surfaces of large size such as surface plates and tables are of considerable importance in precision engineering because they act as reference surfaces for the inspection of other workpieces.

Assessment of the flatness of a surface can be accomplished in a variety of ways and with good fidelity. One of the oldest ways of testing a surface for flatness is to sweep it in several places with a straight-edge and

observe where and how much light leaks through. This method suffers from the sag of the edge and the diffraction of the light through the gap if it is very small. Also the edge has to be rotated on the surface to ensure true planarity (figure 2.100).



**Figure 2.100** Straightness — straight-edge.

Other methods of measuring flatness such as the use of coordinating-measuring machines working relative to a fixed axis, or the use of autocollimators, levels and curvature devices as previously used in straightness measurement, have their attendant problems of noise generation. For this reason the data points have to be used with great care.

In addition to the techniques discussed in the measurement of straightness there are additional ones in flatness. The most important of these is interferometry, especially in the heterodyne mode, which enables the contour of the surface to be evaluated in absolute terms. This technique will be described in the section on optical methods. It must suffice here to say that the technique is viable and does produce, under certain conditions, absolute values.

Flatness is an extension of straightness errors to two dimensions (called areal here as in roughness). One equation of a plane is

$$z = c + m_1 x + m_2 y \tag{2.293}$$

where the $x$ and $y$ are here taken to be the causative variables and $z$ the observational variable.

There are a number of methods of assessing flatness, some similar to straightness. However, because the assessment is somewhat laborious there are often pilot measurements based on a 'Union Jack' pattern (figure 2.101) and limiting the measurements somewhat. Historically it appears that the first reference plane was simply taken to be that plane defined by passing it through three corners of the plate. It has since been suggested that the plane parallel to the diagonals is becoming more widely used.

These two methods are not analytically defined and yet are relatively easy to carry out. The 'Union Jack' pattern is often criticized because of the relatively poor coverage of the area. However, as was seen in straightness, errors are often in the form of a bow and the maximum errors are generally at the edges of the plate or mirror where coverage is most extensive, so the choice of pattern is not so haphazard as might appear.

Not all flatness problems involve surface plates. More and more investigations are involved in the measurement of the flatness of rough objects. In these cases autocollimators and inclinometers are not so much used. Probe methods can be used which measure relative to a fixed internal datum plane or movement.

Consider first the least-squares plane method. As an example the technique adopted in straightness will be extended. The normal equations are similar to those in straightness but have one more variable in number. The variables of minimization are derived from making $S$ — the sum of deviation squared — a minimum, that is

$$S = \sum_{i=1}^{N} [z_i - (c + m_1 x_i + m_2 y_i)]^2 \tag{2.294}$$

is a minimum.

**Figure 2.101** Testing of engineer's surface plate

Note:
The constraint on this equation is the same as for straightness.

Leaving off the limits of summation and subscripts the normal equations become

$$\sum z = cN + m_1 \sum x + m_2 \sum y$$
$$\sum xz = c \sum x + m_1 \sum x^2 + m_2 \sum xy \qquad (2.295)$$
$$\sum yz = c \sum y + m_1 \sum xy + m_2 \sum y^2.$$

From these $m_2$ can be found, yielding

$$m_2 = \frac{N\left(\sum x^2 \sum xy - \sum xz \sum xy\right) - \sum x\left(\sum x \sum yz - \sum y \sum xz\right) + \sum z\left(\sum x \sum xy - \sum y \sum xz\right)}{M}$$

where $M$ is a determinant given by

$$M = \begin{vmatrix} N & \sum x & \sum y \\ \sum x & \sum x^2 & \sum xy \\ \sum y & \sum xy & \sum y^2 \end{vmatrix} \qquad (2.296)$$

obtained directly from equations (2.236). After finding $m_2$ two equations remain. Thus

$$\sum z - m_2 \sum y = cN + m_1 \sum x$$
$$\sum xz - m_2 \sum xy = c \sum x + m_1 \sum x^2 \qquad (2.297)$$

where the right-hand side is similar to that for straightness. From equation (2.291) $c$ and $m$ can be found.

It can also be shown, as in straightness, that the best-fit plane always passes through the centroid of the observed points $x,\ y,\ z$. Referring all measurement to this as origin therefore eliminates the term $c$ in the equations given (also true in straightness).

If uniform increments are taken in orthogonal directions it is possible to simplify these equations considerably, yielding the new plane

$$z = m_1 x + m_2 y \qquad (2.298)$$

where

$$m_2 = \frac{\sum y^2 \sum xz - \sum xy \sum yz}{\sum x^2 \sum y^2 - (\sum xy)^2}$$

$$m_2 = \frac{\sum x^2 \sum yz - \sum xy \sum yz}{\sum x^2 \sum y^2 - (\sum xy)^2}.$$

(2.299)

The error from the mean plane can be expressed either as the maximum peak plus the deepest valley from it or as a standard deviation [91] (RMS value) $\sigma_F$:

$$\sigma_F = \left( \frac{\sum (m_1 x + m_2 y - z)^2}{(n-2)(m_1^2 + m_2^2 + 1)} \right)^{1/2}$$

(2.300)

where the summation is over all points in the plane.

There are much simpler ways to establish a plane, however, and it has to be said that usually there is not much difference between estimates of flatness from any method, as can be seen with reference to table 2.15 below. Perhaps the simplest is to join any three corners together with lines. The plane containing these lines is then used as the reference. This is called the three-corner plane and is often used in industry.

**Table 2.15**

| Diagonals | $E_0$ | $E_1$ | $E_1$ | $E_1$ | $E_1$ | Least Squares $E_2$ | Minimum Zone $E_3$ |
|---|---|---|---|---|---|---|---|
| 1 | 26.0 | 24.0 | 22.0 | 28.7 | 28.7 | 24.9 | 24.0 |
| 2 | 36.0 | 39.0 | 35.4 | 35.6 | 39.0 | 38.5 | 35.4 |
| 3 | 40.0 | 40.0 | 47.0 | 48.0 | 41.0 | 40.2 | 40.0 |
| 4 | 44.0 | 73.0 | 76.6 | 82.2 | 83.0 | 47.2 | 41.9 |
| 5 | 46.0 | 52.1 | 64.1 | 68.5 | 70.9 | 41.8 | 38.1 |
| 6 | 46.0 | 60.0 | 70.5 | 76.7 | 71.8 | 42.7 | 40.8 |
| 7 | 49.0 | 44.8 | 60.0 | 84.0 | 64.0 | 42.8 | 40.7 |
| 8 | 66.0 | 141.1 | 141.4 | 138.4 | 122.0 | 71.0 | 66.0 |
| 9 | 85.0 | 91.8 | 86.8 | 100.0 | 87.4 | 74.0 | 84.3 |
| 10 | 100.0 | 98.3 | 127.0 | 126.0 | 125.0 | 93.0 | 92.2 |
| 11 | 262 | 273 | 273 | 271 | 263 | 264 | 262 |
| 12 | 280 | 597 | 584 | 577 | 537 | 302 | 280 |
| 13 | 800 | 930 | 1210 | 1220 | 1020 | 823 | 752 |

From *Microtechnique* XXVI (7); 1 unit = 1 x $10^{-5}$in.

A least-squares plane is an example of multiple regression as opposed to simple linear regression for straightness.

Yet another approach has been proposed by Miyazaki [102]. He takes the analysis one step further and uses the best-fit least-squares criterion to estimate the *most probable values* of the measurements from those actually measured. To do this he applies the constraint that the total measurement around each of the unit cells making up the measurement area must be consistent. Thus in one cell unit shown in figure 2.102(*a*) the errors must be consistent:

$$\mu_{i,j-1} + v_{i,j} - \mu_{i,j} - v_{i-1,j} = 0$$

(2.301)

**Figure 2.102** Miyazaki method for surface plate.

where the $i$ and $j$ subscripts here refer to the coordinates of the measured points on the plane. For each lattice unit this condition must hold.

For the integrated matching of each of these constraint equations over the whole of the plane two well-known techniques are used. The first makes use of Taylor's theorem to first order. It represents any deviation from the intended or estimated value at a different point as being first order in Taylor's theorem, making the necessary assumption that changes between values are small compared with the values themselves. Thus, in two dimensions, where the deviations are $z(x, y)$,

$$z(x + \Delta x, y + \Delta y) = z(x, y) + \frac{\partial z \Delta x}{\partial x} + \frac{\partial z \Delta y}{\partial y} \qquad (2.302)$$

or

$$\Delta z \simeq \frac{\partial z \Delta x}{\partial x} + \frac{\partial z \Delta y}{\partial y}. \qquad (2.303)$$

Second, use is made of the Lagrangian multipliers to enable the conditional differential equations (each represented by an individual measurement cell) to be solved. Methods like this involving best estimates rather than raw data points are becoming increasingly used in metrology. A point to note is that two-dimensional unit mapping of the surface rather than three-dimensional is all that is required.

Miyazaki therefore obtained *weighting factors* with which to multiply each of the measured values in order to get the most consistent set of measurements overall.

In practice a balance has to be struck between the accuracy of the method and the labour involved. The most accurate method is based on the least-squares method, but it is laborious to measure.

Note here that the method relies upon uniform spacings in the two orthogonal directions.

As an example of how this works, consider a plate whose datum has been taken as that giving the minimum of the square deviation. The weighting factors are shown in figure 2.102(b) taking the simple case where there are only four contiguous cells of measured data.

In the figure, $\delta_{ij}$ represents the deviations from the ideal plane at position $x_{ij}$. Similar tables have been evaluated for different plates.

The whole method is equivalent to an areal weighting function applied to the ordinates obtained relative to the least-squares plane obtained by conventional methods.

The best-fit method of assessing straightness and flatness is also used in measuring circular parts of spheres as well as cylinders and cones. In fact it is the most generally used method, as will be seen later in chapter 3. It has refinements, as in the previous method, using Langrangian multipliers and allowing for the errors in each individual measurement. However, errors in the algorithm generally produce more errors than those in the individual point measurements, so a general well-conditioned algorithm should be used. In chapter 3 such a general method based on the work by Forbes at NPL is used [93].

Such treatments usually form the background to national standards documents, for example BS 7172, the British Standard Guide to the Assessment of Geometric Features, Position, Size and Departure from Nominal Form (1989), and follow from methods of determining the true shape of nominal flats, cylinders, spheres and cones from coordinate-measuring machine measurements.

Cardon *et al* [94] have proposed a different method similar to the minimum-zone method for determining straightness. The definition is a min-max algorithm taken with a coordinate system such that $C=0$. The error in flatness is called $E_f$

$$E_f = \min(\max(z - (m_1 x + m_2 y)) - \min(z - (m_1 x + m_2 y))). \tag{2.304}$$

Whereas for straightness the minimum zone always has two points of contact on the peaks and one on the valleys or vice versa, for flatness the situation is rather different: four constraints are needed and these may well be two peaks and two valleys (or, in the areal texture nomenclature, summits and troughs). This is one occasion where roundness and flatness have a degree of commonality.

The four constraints are needed because of the need to derive the four variables $c$, $m_1$, $m_2$ and the separation of the zones containing the surface.

Some comparisons have been made by the above authors who compared different results taking different datum surfaces [94].

By definition the minimum-zone method must give the smallest error for the plate. As in straightness the divergences between this method and least squares is about 10%.

However, another criterion has to be adopted which is concerned with the time taken to obtain the data, which itself can be related to accuracy as a result of any temperature drift effects. The authors had to use a computer optimization procedure to get the minimum zone.

It has been pointed out that many of these sophisticated methods are unnecessary and uneconomic. For example, using three plates [105] and with hand scraping to remove high spots can give plates flat to about $5\ \mu$m, so it may be only in cases where the surface plate suffers wear and needs to be checked that the more automatic inspection methods of testing are needed. How to achieve automatic inspection quickly is at present being investigated.

Another reason for using an optimized datum plane is that the manufacturers of plates could penalize themselves unnecessarily. The actual flatness error value reported depends on the method of assessment used.

It is obviously worthwhile trying to get better methods of evaluation even for this reason alone. The manufacturer then picks the one that shows the smallest deviation and presumably stands a better chance of being sold.

A great many algorithms have been devised [95] in which the grid system takes into account the shape of the surface to be measured. This will be seen for the methodology roughness measurement covered later in chapter 3.

Figure 2.103 shows an example of triangular cells as opposed to rectangular ones as for the 'Union Jack' methods.

**Figure 2.103** Triangular measurement cell.

Ideally the order in which the points are taken should be specified so as to reduce the errors — and to check on the build-up of errors (e.g. in figure 2.104).

Key positions for the data, particularly those on the diagonals, are taken with extra care, the verticals usually acting as checkpoints.

Another shape is shown in figure 2.104. Notice that there is always some gap left at the edge of the plate. If wear has occurred on the plate it is likely to be at the periphery and so is never used for high-accuracy work.

Triangular cells are being used more often because of the better coverage of the surface. Also the possibility of hexagonal sampling patterns is now being tried because of good coverage and high accuracy.



**Figure 2.104** Triangular cell to fit circular plate

### 2.3.6 Roundness

So far, deviations from ideal shapes that are linear or substantially so, can be related to a straight generator. There are, however, many other different shapes that are equally important. In particular the circle can be

regarded as one such unit geometric element from which more complex engineering forms can be derived. About 70% of all engineering components have an axis of rotational symmetry in them somewhere. The man who first realized its importance was RE Reason of Rank Taylor Hobson, who set out almost single-handed to build instruments to measure roundness [96]. Out-of-roundness or more simply roundness will now be considered as the next logical step in building up a more total geometry of a surface.

How the roundness error can be visualized with respect to the surface texture errors is shown in figure 2.105. Shown in the figure are roughness marks C and waviness marks A and B.



**Figure 2.105**   Roundness and texture.

### 2.3.6.1   *Nature of departures from roundness*

Roundness or strictly out-of-roundness is more important than is usually thought. It is not generally realized that the energy needed to turn or grind comes from rotation, either from the workpiece or the tool. This rotation is absolutely critical to the process and the efficiency of the energy transfer largely determined by the axis of rotation.

Control of the rotational axis is achieved by roundness measurement (Fig. 2.106(*a*)).

Very many of the common machining processes are rotational, e.g. grinding, turning, milling. Only planing and broaching are translational.

Also, in functional situations lateral movement between surfaces is usually of a rotational nature (Figure 2.106(*b*)). Roughness and roundness can both be closely associated with the function map (Fig. 2.106(*c*)).

The first problem is the definition. A multitude of alternatives are in use, some derived from the measuring instrument used and some from the use of the workpiece or its method of manufacture. Generally a workpiece is described as round in a given cross-section if all parts on the periphery are equidistant from a common centre (the Euclidean definition). This sounds obvious but in most cases the part has not got a diametral form described completely by the equation of a circle. Superposed are other types of form. Some are due to the method of manufacture, such as grinding marks, and some are due to errors in the particular machine tool producing the part. These are the phenomena described earlier in surface texture. Usually, however, they are excluded in the primary measurement by means of filters or a blunt stylus acting as a mechanical filter. Out-of-roundness is usually restricted to a consideration of lobing and to some of the higher frequencies. It is in the identifying of the different harmonic components that the characterization takes place. In roundness, however, there is a need for the characterization of the method of collecting the data, as

(a) Manufacture

Workpiece

Cutting tool

Energy to cut material–from rotation
*Roundness provides control of rotation*

(Rotational energy used in turning, milling, grinding,
translation energy only planing)

(b) Function

Lateral movement in function
usually rotation–roundness vital

Normal movement in contact
–roughness vital

Maximum energy transfer in rotation
*Roundness provides control of clearance*

(c) Process

Roughness
influence

Function map

$g$

$g=0$

Roundness influence
(machine tool)

Relative velocity
in lateral direction

**Figure 2.106** Consideration of roundness factors

will be seen. Because of the many different sizes of workpiece it is usually best to specify frequency charac-
teristics in terms of what is effectively a wavenumber, that is in terms of undulations per circumference of
revolution. Lobing is considered to be those deviations having a wavenumber from 2 to about 15, although
sometimes higher numbers are included. These are often of quite high magnitude (e.g. micrometres in height)
and most often of an odd wavenumber. These odd lobes are most frequently almost of constant-diameter form

as a result of being made by centreless grinding. Lobing may exhibit an even or odd number of lobes, and may be of more or less constant height and spacing as produced by grinding, lapping, etc, but can be more random. Lobing is regarded as important, especially in problems of fit. One basic consideration for fit is that the effective size variation for odd-lobed cylinders is positive for external profiles and *vice versa* for internal profiles. There is one big difference between odd and even lobing, and this is that even lobing is measurable by diametral assessment but odd lobing is not. For some applications this constant-diameter property is not important, for instance in ball bearings where only separation between components is the criterion. Because one criterion used for out-of-roundness is the variation in diameter, a comparison between the properties of figures 2.107 and 2.108 is informative. Figure 2.107 is constant radius and figure 2.108 is a Releaux triangle (spherical triangle).



**Figure 2.107** Constant-diameter figure.



**Figure 2.108** Constant-diameter figure — Relaux triangle.

The properties of a Releaux triangle are significant. The generator P′Q′ has two of the properties of diameters of circles. It has constant length and is normal at both ends to the tangent to the curve. The generator P′Q′ constrains the true centre at only three angles and the radius from the centre 0 is only normal to the curve at six points. (In general, for an *n*-lobed figure there are *n* and 2*n* points.) Notice also that the generator P′Q′ is normal to the curve because the constraining curves (figure 2.108) AB, BC and CA are arcs of circles. One further point which can be significant in instrumentation is that the instantaneous centre of rotation of P′Q′ is not the mid-point *R* but one of the three apexes of the triangle, for example point S shown in the figure.

Note:
The Releaux triangle has the smallest area and the greatest deviation from roundness of all constant-diameter figures.

Higher-order lobing up to nearly one thousand undulations per revolution is said to be important from a functional point of view, especially again in the bearing industry where acoustic noise may be important.

Another common form of out-of-roundness is ovality or the two-lobed effect. Multiplicity of this effect is sometimes called polygonation. In some circumstances, parts may need to be oval, for instance in piston rings, so that out-of-roundness errors may not necessarily be detrimental to performance. Oval or elliptical parts are unique in roundness because they can easily be made circular simply by the application of a force of compression on the long diameter, or *vice versa*. Other names have been given to different features of the spectrum of undulations. One used when there has been a chatter set up between the grinding wheel and workpiece is called 'humming' and owes its name to the acoustic noise that it produces and results in a characteristic peak in the spectrum of the workpiece. The typology of roundness is given in table 2.18 (p280).

Another classic example of the importance of odd-lobed roundness is in the conventional measurement of the pitch diameter of screw threads. The pitch diameter is measured indirectly over wires lying in the thread spaces above and below, as shown in figure 2.109 for a 60° screw-thread gauge. If the wire is, say, elliptical the diameter is suspect.



**Figure 2.109** Screw-thread measure.

The higher undulations are usually not as important from the functional point of view as the lower ones, but they may be quite large in amplitude. Instrumentation has to be used to isolate the two bands. Typical methods involve mechanical filtering by means of a blunt stylus and/or electrical filtering.

Roundness means different things to different people. This is why the characterization is important. There are three basic ways in which roundness can be measured. These will be outlined here but can be obtained in some detail from the works of Reason [96]. Some of the distortions of such methods are given here. Algorithms associated with the assessment are given in chapter 3 on processing.

The three methods are diametral, chordal and radial. The first is the measurement of the diameter of a component by means of a micrometer or a height gauge or any device which uses two parallel flat anvils as the sensor. The out-of-roundness is regarded as one-quarter of the difference between the maximum and the minimum readings of $D$ as the part is rotated (figure 2.110). Thus, to a good approximation

$$D(\theta) = r_1(\theta) + r_2(\theta - \pi). \tag{2.305}$$



**Figure 2.110** Radial representation.

To see how this technique responds to different equally spaced lobes the Fourier coefficients of equation (2.304) are used. From equation (2.305) $F_m(n)$ is the measured coefficient for the $n$th lobe and $F_{actual}(n)$ is the true coefficient which may be zero or may exist:

$$F_m(n) = F_{actual}(n)[1 + \exp(jn\pi)]. \tag{2.306}$$

Notice that when $n$, the number of undulations per revolution, is an odd number then $F_m(n)$ is zero — despite the fact that $F_{actual}(n)$ exists. Hence diametral methods cannot see bodies with an odd number of lobes. As mentioned before this might not be serious if the body is to be used as a spacer such as in a roller bearing, but it would be serious in, say, a gyroscope which spins about its centre.

### 2.3.6.2 Chordal methods

The second method represents the traditional vee-block method. Strictly, it is the variation in the chord joining the points of contact of the part and the two faces of the vee. This variation induces changes in the vertical position of the gauge when the part is rotated.

Note that chordal methods can only be converted into radial results from a common centre when undulation spacing is uniform and the vee angle is chosen correctly (figure 2.111).



**Figure 2.111**   Vee-block measurement of roundness.

As the part is rotated the indicator reading varies. The difference between the maximum and minimum readings is noted. This is called the total indicator reading (TIR).

It is interesting to note that this technique is a derivative of the intrinsic equation method of defining an arbitrary curve relative to itself by means of a target at P and the length of curve from an arbitrary axis. Here the target is where the vee contacts the part and the arc length is the circumferential distance between the two points of contact.

As in the case of the diametral method this technique introduces distortions into the assessment of roundness. These need to be understood because of the general use and misuse of the vee-block approach in industrial workshops. The distortion of the lobes in this case is given by.

$$F_{\text{measured}}(n) = F_{\text{actual}}(n)\left(1 + (-1)^n \frac{\cos\ n\gamma}{\cos\ \gamma}\right). \tag{2.307}$$

For example, if $n=5$ and $y=60$, then $F(5)$ is measured to be zero although it is obvious that the workpiece is not perfectly round!

This distortion factor obviously varies according to the angle of the vee. This should be remembered because recent instrumental methods discussed later include this basic principle. It is clear that it is pointless trying to characterize a signal if part of it is distorted by the method of measurement.

To be absolutely sure that a part is round using this method more than one vee has to be used. In order to use this technique effectively some idea of the probable lobing on the workpiece has to be known so that it is impossible to fall into the trap shown in equation (2.307). Even then this equation would not guarantee that problems could not arise, because it assumes uniform angular separation of the lobes. Methods based on this technique have been used for some time with limited success [97].

Other investigators [98] have worked out the technique and introduced an extra degree of freedom by offsetting the indicator dial (figure 2.112), the offset probe being in effect equivalent to using another vee to reduce aberrations in the measurement.

This investigation showed that it is possible to optimize the offset angle and the vee angle to minimize the spread of distortions obtained over the important range of undulations from 3 to 15 (figure 2.113).

Fortunately such methods of roundness measurement in which a degree of distortion has to be accepted [98,99] are now largely redundant with the advent of computers.

Similar techniques have been tried without recourse to just two points of contact. One is shown in figure 2.114.

Figure 2.112 Offset probe method.



Figure 2.113 Offset probe method.



Figure 2.114 Multiple floating skid method.

This method is not really a vee-block method because the reference is taken to be the integrated effect of a large number of point contacts which are freely mounted on the yoke. The indicator is able to move through the yoke and has a sharp point. Variations between the indicator stylus and the integrated yoke position constitute the out-of-roundness.

Although this is not an instrumentation chapter, it is worth mentioning here that such a method would clearly be an instrumentation and operation nightmare. However, recent work by M Svitkin of Technornach in St Petersburg has made this method viable in specification and cost [140].

### 2.3.6.3 Radial methods

The third and preferred method is to compare the workpiece shape directly with a true circle. This can be achieved by using a mechanical reference such as a disc or annulus or a spatial reference in the form of an arc in space generated by a very good bearing in the spindle in figure 2.115(*a*) and in the rotating table in figure 2.115(*b*). Again, as in flatness, it is necessary to refer to the measuring instrument before assessment can take place. Instrument distortion, as with the vee-block method, still needs to be understood before meaningful assessment can take place. In many cases the assessed value of a surface parameter should be prefixed with the instrument used to obtain the data.

The methods are shown in figure 2.115(*a-c*). In figure 2.115(*c*) there is a reference disc which is fixed to the workpiece to be measured. A gauge contacts the master disc and another contacts the workpiece. The difference between the two signals is processed, not the individual signals.



**Figure 2.115** Various methods of assessing roundness using radial techniques.

Sometimes the master disc itself is left off together with one probe and the accurate rotation left to the centres on which the component is held, but this method is fraught with problems due to sag of the workpiece, poor centres, etc.

Other methods using mechanical references have been devised, for example when the reference has been a ring into which the workpiece, assumed to be a shaft or cylinder, is placed. A gauge then measures the difference between the two; problems of concentricity between the reference part and the workpiece are sometimes serious and unavoidable. Also, other errors can occur because of misalignment of the axis of the workpiece and the reference piece.

Out-of-roundness is sometimes referred to in the USA as DFTC (departure from true circle) which best describes the nature of the measurement. Whilst on the subject of names, out-of-roundness is sometimes shortened to just roundness for convenience (a shortcut not used in most European countries) and sometimes it is referred to as OOR. In any case, as a reminder it represents the deviation of the workpiece periphery from a true circle drawn somewhere in space. The reference circle need not be the same size nor need it be concentric (although it usually is). Eccentricity errors are usually due to the set-up and not to real deviations, except in exceptional circumstances such as the egg shape. There is a trend towards integrated measurement in which the size and shape are measured at the same time by the same instrument. This will be justified when size tolerances are less than out-of-roundness errors. There is evidence that this situation is already here in some cases, such as in gyroscopes and compact-disc spindles. Under these circumstances great care has to be exercised with the algorithms to ensure that they are compatible.

*2.3.6.4  Nature of the signal produced by a radial departure instrument*

In order to see the out-of-roundness at all it is necessary to magnify the deviations from a perfect circle and to display them on a chart in some form or other. It is from this record that judgements are made about the nature of the deviations, for example the amount and type of lobing. Unfortunately the fact that high magnifications have to be used can cause a certain amount of distortion of the graph to occur under certain circumstances. These distortions can, and often do, interfere with the interpretation of the out-of-roundness. For this reason it is impossible completely to divorce any discussion of the nature of roundness from that of instrumentation. In what follows the nature of the relationship between the workpiece and the graph will be explored so that misinterpretations in practice will be minimized. The whole issue of interpretation will be brought out in the roundness section of processes.

Because, in principle, we are making a transformation between the real part deviations and some graph or storage medium, the visual picture of one need not resemble the other in all of its features. The choice of which coordinate system to use depends upon which detail needs to be enhanced. This choice, the problems, and some of the advantages of the type of display will now be considered.

This section highlights the fact that the roundness of the component itself and the signal received from the roundness instrument do not necessarily have a completely straightforward relationship. In order to get any satisfactory signal some distortions result. The nature of these distortions is sometimes misleading and it therefore has to be understood.

The departures from true roundness, as revealed by the instrument, are plotted using either polar or Cartesian coordinates. Both kinds have geometric properties which need to be understood if the graphs are to be interpreted correctly. This is important because a wrong assessment can easily be produced.

Polar coordinates provide a realistic display that can be immediately related to the workpiece. For example, if the workpiece is held over the graph in correct orientation, the directions in which the crests and valleys occur can be seen directly, and this is often convenient.

Polar graphs have the useful property that, even if there is a small amount of eccentricity between the workpiece and the axis of rotation, the resulting eccentric plot of the periphery of a truly round workpiece may still be substantially circular. Small residual eccentricity in setting up an instrument can therefore be accepted. Polar graphs as normally obtained have the further property that chords through the centre of the workpiece plot as chords through the centre of rotation of the chart regardless of the residual eccentricity. The extent to which the workpiece has a constant diameter can therefore be checked on the graph by direct measurement through the centre of the chart, as shown in figure 2.116.



**Figure 2.116**  Choice of coordinates

A point to be watched is that, generally, the peak-to-valley height of the plot must not be too great a proportion of the mean radius, otherwise the disproportionate compression of peaks and valleys may greatly falsify both the visual and numerical assessments. Cartesian coordinates permit closely spaced detail to be plotted on a more open circumferential scale than is generally practicable with polar graphs (where the scale is limited by the radius on which the plot is made) but do not permit the residual eccentricity to be so readily separated from the undulations to be measured. As an example, each of the Cartesian plots $A_r$, $B_r$, $C_r$, $D_r$ in

figure 2.117 was taken at the same time as the corresponding polar plots $A_p$, $B_p$, $C_p$, $D_p$ in figure 2.118. While the four polar plots are much alike and have substantially the same peak-to-valley height despite the different amounts and directions of eccentricity, the Cartesian plots are not alike either in appearance or in peak-to-valley heights.



**Figure 2.117** Cartesian display of roundness.



**Figure 2.118** Polar display of roundness.

This is because the proper reference line for the polar plot remains substantially circular even in the presence of a certain amount of eccentricity, whereas the reference line for the Cartesian plot becomes a sine wave, of which the amplitude cannot be determined by inspection. In consequence, the tolerance on centring for a polar graph can be several times the radial error, while for a Cartesian graph it must be negligible in comparison. For many years this was a dominant consideration. More recently it has become possible to compute and plot (or eliminate) the sinusoidal reference line electronically, so that Cartesian coordinates can now be used more often.

The practical point in favour of the Cartesian presentation, which will be clarified in the next section, is that area measurement is more valid; attempts to measure, say, wear scars on ball bearings using the polar graph can be fraught with danger. Simple areas measured off the chart will not be valid. This is due to the fact that the instrument is fundamentally measuring a 'skin' of the surface. This skin looks much more Cartesian in nature than polar.

### 2.3.6.5  Relation between the centred workpiece profile and the radius-suppressed polar plot

If the equation of the true workpiece as measured from its centre is $r(\theta)$, in order to be able to see and measure the out-of-roundness components on a polar graph (i.e. the departures of $r(\theta)$ from a true circle), it is necessary to magnify the radial values considerably. However, only a small component of the term $r(\theta)$ is the out-of-roundness, the rest being the physical radius of the workpiece. At the magnifications generally necessary to measure the small changes in $r(\theta)$ an excessive size of chart, perhaps ~ $10^3$ metres, would be necessary if the total radius term was magnified. The only practical solution is to magnify the differences between $r(\theta)$ and a true circle, say $L$, having a radius nominally equal to the physical mean size of $r(\theta)$ taken all round. Then, if the mean radius of $r(\theta)$ is $R$, what appears on the graph will be $\rho(\theta)$, where $\rho(\theta) = M(r(\theta) - L) + S$ with $S$ the inside radius of the chart (usually chosen to be of the order of 10mm or so).

The important point to note is that $\rho(\theta)$ does not have the same shape as $r(\theta)$ owing to the effect of radius suppression and the magnification.

This can be seen by looking at figure 2.119(a) which shows these effects. Simply removing a radial value of $X$ from a workpiece having a nominal radius $R$ and peak-to-valley out-of-roundness $H$ changes the appearance. What was a convex-sided shape is transformed into a concave-sided shape. The effect of magnification is similar. Magnifying the out-of-roundness by different amounts produces the same sort of effect (figure 2.119(b).



**Figure 2.119**  Effect of radius suppression.

The kind of distortion shown in these figures can be minimized by confining the graph to an outer zone of the chart, but even so, careful thought and some degree of calculation may be needed before the true shape of the part can be correctly deduced from the graphical representation of the undulations.

### 2.3.6.6  Effect of imperfect centring

Apart from the effects on the shape of the centred graph caused by the suppression of the radius, further distortion occurs due to imperfect centring of the workpiece on the reference axis, affecting its appearance and hence the correct comparison of radii, diameters and angles [100].

To see what this effect is, consider a truly circular workpiece of radius $R$.

The polar equation of such a workpiece, whose geometric centre O' is positioned eccentrically by an amount $e$ at an angle from the centre of rotation O of the instrument is given by a function $k(\theta)$ where:

$$k(\theta) = e \, \cos(\theta - \varphi) + [R^2 - e^2 \, \sin(\theta - \varphi)]^{1/2} \dots \tag{2.308}$$

or

$$k(\theta) = e\ \cos(\theta - \varphi) + R - \frac{e^2}{2R}\ \sin^2(\theta - \varphi)] + \ldots \ . \tag{2.309}$$

This function $k$, a function of $\theta$, is the same as $r(\theta)$ when $e=0$ and looks perfectly circular off axis ($e \neq 0$) as well as on axis ($e=0$).

Consider the situation where the same part is measured by a roundness instrument. What is actually measured by the transducer is given by the following function $\rho(\theta)$ where $\rho(\theta)=M(k(\theta) - L)$, where $M$ is the magnification and $L$ is the suppressed radius. When this is displayed on the polar chart of the instrument it becomes (figure 2.120).

$$\rho(\theta) = M(K(e) - L) + S \ldots \tag{2.310}$$

where $S$ is the inner radius of the chart



**Figure 2.120** Spatial situation $k = e\ cos\ (\theta - \varphi) + \sqrt{R^2 - e^2\ \sin^2(\theta - \varphi)}$ (circle).

A polar representation of $\rho(\theta)$ given by

$$\rho(\theta) = M[e\ \cos(\theta - \varphi) + M(R - L) + S - \frac{Me^2}{2R}\sin^2(\theta - \varphi) + \ldots \ . \tag{2.311}$$

This equation no longer represents the shape of a displaced circle because the relative sizes of the individual terms compared with those in equation (2.309) have been changed.

In particular this can be seen with reference to the first two terms in equation (2.308) $R$ is typically 25mm whereas $e$ is 0.02mm, the ratio being of the order of 1000:1. Comparison between the corresponding terms in equation (2.309) reveals that this ratio is much smaller. S is typically 12mm, $M$ is 1000, $R$-$L$=0.02 mm. Thus $M(R-L)+S$: $M\ e\ cos(\theta - (\alpha)$ is of the order of 1.1 for a typical roundness graph. The disproportioning of the relative magnitudes of the terms means that the function representing the transducer signal no longer looks circular when plotted. In fact it is not circular!

Because of the equivalence of the scale of size of the first two terms and the relative unimportance of the following terms, equation (2.311) may be written as

$$\rho(\theta) = M(r - L) + S + Me\ \cos(\theta - \varphi) \tag{2.312}$$

without significant loss of accuracy for most practical situations. Letting $M(R - L) + S = t$ and $M e = E$, equation (2.312) becomes

$$\rho(\theta) = t + E \ \cos(\theta - \varphi). \tag{2.313}$$



**Figure 2.121** Eccentric workpiece — limaçon effect, 8 mm eccentricity on chart (30% chart eccentricity).



**Figure 2.122** Eccentric workpiece — limaçon effect, 12mm eccentricity on chart (50% chart eccentricity).

Thus, whereas equation (2.308) represents the true equation of a circle displaced about a centre, the equivalent equation (2.313) displayed by a roundness-measuring instrument does not represent the equation of a circle but that of a limaçon as is shown in figures 2.121 and 2.122. Notice that the limaçon is the general form of a cardioid.

The limaçon form gives a clue as to the most suitable method of generating a true reference for a circular part as seen by a roundness instrument. It represents in effect the first Fourier coefficients of the transducer signal, that is the dc term and first harmonic.

### 2.3.6.7  Assessment of radial, diametral and angular variations

The eccentric polar graph has some interesting properties. It can be shown that, in general, diametral variation and angular relationships should be measured through O, the centre of the chart, and not through the centre of the plot on the chart. This is easily seen for the practical case where $e/R$ is small and the polar form of the mag-

nified workpiece as seen on the chart is a limaçon. Measuring diameters should always give a constant value if the workpiece is perfectly circular. Thus the diametral measurement through O is $\rho(\theta)+\rho(\theta + 180)=2t$ which is always constant. Measurement through the apparent centre of the workpiece will not yield a constant value but will depend upon O. It will be a maximum value at right angles to the direction of the eccentricity and a minimum in the direction of eccentricity of $2t$.

Obviously for measurements taken through the origin of the part on the graph rather than the chart centre to be valid then the magnitudes of the terms in equation (2.309) must compare with those of (2.313), that is $S \sim ML$, which corresponds with the situation where there is no radius suppression, that is the inner radius of the graph $S \geq LM$. Then the graph no longer takes on the shape of either a limaçon or a circle. However, this situation is very rare and only occurs when the magnification is small and the radius of the part is small, which automatically makes $L$ small and the apparent shape very different. Instead of the bulge at the centre at right angles to the origin it is a reduction!

Angular considerations are similar. In normal circumstances for eccentric parts the angular relationships of the component are only valid when measured through the centre of the chart; it is only in special cases where there is only a small amount of zero suppression that consideration should be given to measurement through a point in the region of the centre of the part. This is shown in figure 2.123.



**Figure 2.123** Angular relationship through chart centre.

It is possible to make the limaçon obtained when the workpiece is eccentric simply by adding further terms in equation 2.311. The displaced graph can look more and more circular despite being eccentric. However, this doesn't mean that angular relationships (in the eccentric 'circular' trace) are corrected. All angle measurements still have to go through the centre of the chart. Centring the corrected graph by removing the eccentricity term is the only way that the centre for roundness is at the circle centre — it is also at the centre of rotation so that there is no problem.

Radial variation can be measured from the centre of the profile itself rather than the centre of the chart, but the measurements are subject to the proviso that the decentring is small.

R E Reason has given maximum permissible eccentricities to allow reasonably accurate measurement of the diametral, radial and angular relationships, subject to the criteria that the differences are just measurable on the graph. These are listed in the following tables.

**Table 2.16**

| Eccentricity $E$ (mm) | Radial error $R'-R$ | Diametral error $D'-D$ |
|---|---|---|
| 1.35 | 0.0212 | 0.0425 |
| 2.5 | 0.0875 | 0.175 |
| 5.0 | 0.35 | 0.7 |
| 7.5 | 0.75 | 1.5 |

The errors in table 2.16 refer to the eccentric errors of a graph of mean radius 40 mm.

Table 2.16 is important because it shows the permissible tolerance on centring for different purposes as shown on the chart. For example, with 1.25 mm eccentricity the error is too small to be detected, while with 2.5 mm eccentricity it will only just be measurable. Above this the difference will only matter if it is a large enough proportion of the height of the irregularities to affect the accuracy of their assessment. Eccentricity up to 5 mm can generally be accepted for normal workshop testing, with 7.5 mm as an upper limit for a graph around 75 mm diameter.

Table 2.17 shows that the more perfect the workpiece the better it needs to be centred. This requirement is generally satisfied in the normal use of the instrument, for good parts tend naturally to be well centred, while in the case of poorer parts, for which the criterion of good centring is less evident, the slight increase in ovality error can reasonably be deemed of no consequence.

Some diametral comparisons with and without eccentricity are also shown in table 2.17.

**Table 2.17**

| Diameter of workpiece (mm) | Eccentricity of graph | | |
| | 2.5 mm | 5 mm | 7.5 mm |
| | Magnification must exceed: | | |
| --- | --- | --- | --- |
| 0.25 | 400 | 1600 | 3600 |
| 0.5 | 200 | 800 | 1800 |
| 1.25 | 80 | 320 | 720 |
| 2.5 | 40 | 160 | 360 |
| 5 | – | 80 | 180 |
| 12.5 | – | 40 | 72 |

A practical criterion for when diameters can no longer be compared through the centre of the chart can be based on the smallest difference between two diameters of the graph that could usefully be detected. Taking 0.125mm as the smallest significant change, table 2.17 shows the lowest permissible magnification for a range of workpiece diameters and eccentricities when the graph has a mean diameter of 3 in (75 mm).

Notice how in figure 2.123 that, even when the workpiece has been decentred, the valley still appears to point towards the centre of the chart and not to the centre of the graph of the component. This illustrates the common angular behaviour of all roundness instruments.

A criterion for when angles should no longer be measured from the centre of rotation can be based on the circumferential resolving power of the graph. Allowing for an error in the centring of the chart itself, this might be in the region of 0.5 mm circumferentially.

If lower values of magnification should be required then the definitive formulae should be consulted. This applies to both the diametral measurement and the angular relationships. These tables merely give the nominally accepted bounds. Summarizing the foregoing account of the properties of polar graphs, it will be seen that the following rules can generally be applied:

*Plotting*:

1. To avoid excessive polar distortion, the trace should generally be kept within a zone of which the radial width is not more than about one-third of its mean radius.
2. The eccentricity should be kept within about 15% of the mean radius for general testing, and within 7% for high precision.

*Reading*:

1. Points 180° apart on the workpiece are represented by points 180° apart through the centre of rotation of the chart.
2. Angular relationships are read from the centre of rotation of the chart.
3. Diametral variations are assessed through the centre of rotation of the chart.
4. Radial variations are assessed from the centre of the profile graph, but are subject to a small error that limits permissible decentring.
5. What appear as valleys on the chart often represent portions of the actual surface that are convex with respect to its centre.

Modern measuring instruments are making it less necessary to read or plot graphs directly, but the foregoing comments are intended to provide a suitable background from which all advances can be judged.

### 2.3.6.8 Roundness assessment

Clearly, there are many parameters of roundness that might be measured, for example diametral variations, radial variations, frequency per revolution (or undulations per revolution), rates of change (velocity and acceleration). Most, if not all, could be evaluated with respect both to the whole periphery and to selected frequency bands. Radial variations can be assessed in a number of ways, for example in terms of maximum peak-to-valley, averaged peak-to-valley, and integrated values like RMS and $R_a$. As far as can be seen, there is no single parameter that could fully describe the profile, still less its functional work. Each can convey only a certain amount of information about the profile. The requirement is therefore to find out which parameter, or parameters, will be most significant for a given application, remembering that in most cases performance will depend on the configuration of two or more engaging components, and that roundness itself is but one of the many topographic and physical aspects of the whole story of workpiece performance. Assessment is now widely made on a radial basis because the parameters so determined provide information about the quality of roundness regardless of the number of irregularities. It is with variations in radius that the present method is mainly concerned.

A basic point is that whatever numerical assessment is made, it will refer to that profile, known as the measured profile, which is revealed by the instrument and is in effect the first step in characterization.

The peak-to-valley height of the measured profile, expressed as the difference between the maximum and minimum radii of the profile measured from a chosen centre, and often represented by concentric circles having these radii and thus forming a containing zone, is widely used as an assessment of the departure of a workpiece from perfect roundness. This is called the 'roundness error', the 'out-of-roundness' error or sometimes DFTC (departure from true circle).

The centre can be determined in at least four different ways which lead to slightly different positions of the centre and slightly different radial zone widths in the general irregular case, but converge to a single centre and radial zone width when the undulations are repetitive. All four have their limitations and sources of error. These four ways of numerical assessment are referred to and described as follows (see figure 2.124):

(1) Ring gauge centre (RGC) and ring gauge zone (RGZ)
If the graph represents a shaft, one logical approach is to imagine the shaft to be surrounded with the smallest possible ring gauge that would just 'go' without interference. This would be represented on the chart by the smallest possible circumscribing circle from which circle the maximum inward departure (equal to the difference between the largest and smallest radii) can be measured. As mentioned in the introductory section this is a functionally questionable argument.

(2) Plug gauge centre (PGC) and plug gauge zone (PGZ)
If the graph represents a hole, the procedure is reversed, and the circle first drawn is the largest possible inscribing circle, representing the largest plug gauge that will just go. From this

**Figure 2.124** Methods of assessing roundness. $R_{max} - R_{min} = 0.88$ mm (*a*); 0.76 mm (*b*); 0.72 mm (*c*); 0.75 mm (*d*).

is measured the maximum outward departure, which can be denoted on the graph by a circum-scribing circle concentric with the first.

(3) Minimum zone centre (MZC) and minimum zone (MZ)
    Another approach is to find a centre from which can be drawn two concentric circles that will enclose the graph and have a minimum radial separation.

(4) Least-squares centre (LSC) and least-squares zone (LSZ)
    In this approach, the centre is that of the least-squares circle.

The obvious difference between the ring methods and the least-squares circle is that whereas in the former the highest peaks and/or valleys are used to locate the centre, in the least-squares circle all the radial measurements taken from the centre of the chart are used. Another point is that the centre of the least-squares circle is unique. This is not so for the maximum inscribing circle nor for the minimum zone. It can be shown, however, that the minimum circumscribing centre is unique, therefore joining the least-squares circle as most definitive. Methods of finding these centres will be discussed in chapter 3. Fig. 2.125 shows the reliability of these methods.

Of these four methods of assessment the least square method is easiest to determine.

Summarizing, the only way to get stable results for the peak-valley roundness parameters is by axial averaging i.e. taking more than one trace. The least squares method gets its stability by radial averaging i.e. from one trace. Put simply the least squares is basically an integral method whereas the other three methods are based on differentials and are therefore less reliable. As will be seen the calculation of the least square parameters is straightforward.

Although the minimum zone method is more difficult to find than the plug gauge and ring gauge methods it leads to a more stable determination of the common zone between a shaft and bearing as seen in Figure 2.126. Four points determine the clearance zone using the minimum zone method whereas only two determine the clearance zone using the plug/ring gauge method.

Note:

(1) the energy interaction takes place at the common zone.
(2) Common zone plug/ring — two point interaction.
(3) Common zone minimum zone — four point interaction.

Figure 2.125  Reliability of circular reference systems.



Figure 2.126  Shaft and journal bearing.

Conceptually, the ring gauge, plug gauge and minimum zone methods are graphical in origin, the inspector working on the chart with a pair of compasses and a rule. These methods therefore provide simple practical ways for assessing the out-of-roundness. They suffer to some extent from two disadvantages. First, the fact that the centres so derived are dependent on a few isolated and possibly freak peaks or valleys which can make the measurement of concentricity somewhat risky. Second, a certain error is bound to creep in if the graph is eccentric because the graphical method depends on circles being drawn on the chart whether or not the part is centred.

It has already been shown that the shape of a perfectly round part as revealed by an instrument will look like a limaçon when decentred, and therefore ideally the inspector should use compasses that draw limaçons [90].

The form of this error can easily be seen especially using some of the newer instruments, which use large graphs and small central regions.

For example, consider the ring gauge method (minimum circumscribed circle). This will be determined by the largest chord in the body, which is obviously normal to the vector angle of eccentricity. Thus from figure 2.127 it will be $c = 2 \rho \sin \theta$, where $\rho = t + E \cos \theta_1$ and $\theta_1$ is the angle at which the chord is a maximum, that is $\rho \sin \theta$ is maximum.



**Figure 2.127**   Use of large polar chart.

Thus $\sin \theta (t + E \cos \theta)$ is a maximum from which

$$\frac{dc}{d\theta} = 2E \cos^2\theta_1 + t \cos\theta_1 - E = 0. \tag{2.314}$$

Hence

$$\theta_1 = \cos^{-1}\left(\frac{-t + \sqrt{t^2 + 8E^2}}{4E}\right). \tag{2.315}$$

This will correspond to a chord through a point O″, that is $(t + E \cos \theta_1) \cos B \theta_1$.

Notice that this does not correspond to the chord through the apparent centre of the workpiece at O′, a distance of $E$ from O. The angle $\theta_2$ corresponding to the chord such that $(t + E \cos \theta_2) \cos \theta_2 = E$, from which $E \cos^R \theta_2 = t \cos \theta_2 - E = 0$,

$$\theta_2 = \cos^{-1}\left(\frac{-t + \sqrt{t^2 + 4E^2}}{2E}\right) \tag{2.316}$$

from which it can be seen that $\theta_2$ is always less than $\theta_1$. The maximum chord intercepts the diameter through O and O′ at a distance of less than $E$ (because the terms under the root are substantially the same).

The maximum chord value is given by $2(t + E\cos\theta_1)\sin\theta_1$ using the same nomenclature as before. The apparent error in the workpiece measured on the chart (due to using compasses on the badly centred part) will be seen from figure 2.131 to be $d$, where $d$ is given by

$$d = (t + E\cos\theta_l)(\sin\theta_l - \cos\theta_l) + E - t \tag{2.317}$$

which has to be divided by the magnification to get the apparent error on the workpiece itself. It should be emphasized that this effect is rarely of significance. It is not important at all when computing methods are used, but as long as graphical verification of results is carried out and as long as manufacturers use larger polar graph paper with smaller centres there is a danger that this problem will occur.



**Figure 2.128** Effect of centring on minimum zone.

This apparent out-of-roundness value of the graph on the chart would be measured even if the workpiece and the spindle of the instrument were absolutely perfect. Ways of computing these centres and zones without regard to this inherent error obtained when evaluating graphically will be discussed in chapter 3.

This distortion can also be troublesome when examining charts for lobing and in some circumstances can cause confusion in interpreting cylinders having a tilted axis relative to the instrument datum.

### 2.3.6.9 *Effect of imperfect centring on the minimum zone method*

On a perfectly circular eccentric part the form revealed by a roundness instrument is as shown in figure 2.128. It can be shown that the centre of the minimum zone coincides with that of the least-squares circle and plug gauge for this case. This is because moving the centre of the circle from the plug gauge centre at O′ inwards towards O minimizes the plug gauge radius faster than it minimizes the ring gauge radius. To find the apparent measurement of out-of-roundness it is necessary to find the maximum radius from O′, that is $d_{max}$ which is obtained when

$$d = (t^2 + E^2 \sin^2\theta)^{1/2} \tag{2.318}$$

is a maximum, that is

$$\frac{E^2\sin^2\theta}{(t^2 + E^2\sin^2\theta)^{1/2}} = 0$$

from which $\theta = \pi/2$ and $d_{max} = (t^2 + E^2)^{1/2}$. Hence the error is

$$\varepsilon = d_{max} - t = (t^2 + E^2)^{1/2} - t \tag{2.319}$$
$$\varepsilon \simeq E^2/2t.$$

This is what would be measured by the operator. For $E = 10$ mm and $t = 30$ mm the error could be 1.7 mm on the chart, a considerable percentage!

### 2.3.6.10  Effect of angular distortion

From figure 2.128 the relationship between $\alpha$ as measured from O and $\theta$ as measured from O′ is

$$\alpha = \tan^{-1}\left( \frac{(t + E\cos\theta)\sin\theta}{(t + E\cos\theta)\cos\theta - E} \right) \tag{2.320}$$

where for convenience the eccentricity has been taken to be in the $x$ direction.

As previously mentioned the angles should always be measured through the centre of the chart irrespective of the eccentricity for normal purposes.

This is easily seen by reference to figure 2.132. Merely magnifying the derivations from the centred part and transcribing them to a chart as in figure 2.133 does not significantly change the angular relationships between the arrows as marked off. However, if the angles between the arrows are measured relative to the



**Figure 2.129**  Actual extent of movement of workpiece relative to its size.



**Figure 2.130**  Effect of angular distortion.

apparent centre O′ (which an operator may think is the true centre), considerable distortion of the results occur. Figure 2.131 shows what actually happens to angles on the chart when the workpiece is decentred. Instead of the arrows being separated by 60° they are seemingly at $\alpha_1$ and $\alpha_2$ which are gross distortions.



**Figure 2.131** Effect of angular distortion.

Although this angular invariance of the centred and eccentric charts is somewhat astonishing, it is not usual for it to cause confusion except in special cases. One such case is measuring any type of rate of change or curvature of the profile from the chart.

A slope feature subtending an angle $\delta\theta$ in the centred case of figure 2.132(a) will still subtend it in the eccentric case of figure 2.132(b). The feature will appear to enlarge on the one side and shrink on the other. In fact, however, they still subtend the same angle $\delta\theta$ about the chart centre. But measuring any feature from the apparent centre at O′ will give considerable angular errors, which in turn give corresponding errors in slope because it is a function of $\theta$.



**Figure 2.132** Effect of angular distortion on slope measurement.

Differentiating the equation for angle shows that $\delta x_1$ is $(t + E)/t$ times as big as $\delta\theta$ and $\delta x_2$ is $(t - E)/t$ times $\delta\theta$.

For $E = 10$ mm and $t = 30$ mm the subtended angle of a feature as seen from the apparent centre O′ can differ by a factor of 3:1 depending on where it is relative to the direction of eccentricity!

Hence if the slope feature has a local change of radius $\delta r$ the value of $\delta r / \delta\alpha$ will vary by 3:1 depending on where it is. For $E = 10$ mm the variation is 2:1. The extent of this possible variation makes quality control very difficult. The answer will depend on the purely chance orientation of the slope feature relative to the direction of eccentricity.

Measuring such a parameter from O, the chart centre, can also be difficult in the highly eccentric case because $d\rho/d\theta = E \sin\theta$ which has a minimum value of $E$ length units per radian at a direction of $\theta = \pi/2$, that is perpendicular to the direction of eccentricity. More affected still are measurements of curvature because the $d\theta$ to $dx$ distortions are squared. The only safe way to measure such parameters is by removing the eccentricity by computation or by centring the workpiece accurately.

Amongst the most obvious difficulties associated with the application of zonal methods is the possible distortion of the centre position due to irregular peaks on the circumference. An example of this can be seen with reference to the measurement of an ellipse using the plug gauge method.

Apart from the effects of eccentricity, polar distortion can affect more especially the drawing of the inscribed circle.

Consider, for example, the representation of an ellipse at 200× magnification in figure 2.133(*a*). All possible centres coincide and there is no ambiguity. But if the magnification is increased to 1000×, the representation acquires the shape shown in figure 2.133(*b*). Although the MZ, RG (and LS) centres remain the centre of the figure, two centres can now be found for the circle representing the largest plug gauge. Thus, while the MZ, RG and LS evaluations are the same for both magnifications, the plug gauge value, if based literally on the maximum inscribed circle, is erroneously greater for the higher magnification. In the former, plotted on as large a radius as the paper permits, the ambiguity of centres is just avoided, but on a small radius at the same magnification the ambiguity reappears. It is therefore important, when seeking the plug gauge centre, to keep the zone width small compared with its radius on the graph, that is to plot as far out on the paper as possible and to use the lowest magnification that will provide sufficient reading accuracy.



**Figure 2.133** Problems in plug gauge assessment.

This is a practical example of how zonal methods based upon graphical assessment can give misleading results if applied literally. Again, they highlight the importance of knowing the nature of the signal. Because the best-fit circle is only a special case of that of a best-fit limited arc, the general derivation will be given from which the complete circle can be obtained. The only assumption made in the derivation is that for practical situations where what is required is a best-fit partial limaçon rather than a circle simply because of the very nature of the instrumentation. How this all fits in with the graphical approach will be seen in the next section.

*2.3.6.11   Equation of a reference line to fit a partial circular arc as seen by a roundness instrument [110]*

Keeping to existing convention, let the raw data from the transducer be $r(\theta)$, $r$ having different values as $\theta$ changes due to the out-of-roundness or roughness of the part.

Remembering that the reference line to the data from the transducer before display on the chart has the equation

$$M(R - L) + Me \, \cos(\theta - \varphi) \tag{2.321}$$

and letting

$$M(R - L) = S \qquad Me = E \tag{2.322}$$

and

$$E \cos \varphi = x \qquad E \sin \varphi = y \tag{2.323}$$

then the limaçon form for the reference line between $\theta_1$ and $\theta_2$ is

$$\rho(\theta) - S = R + x\cos\theta + y\sin\theta$$

and in order to get the best-fit limaçon having parameters $R$, $x$, $y$ to the raw data $r(\theta)$ the following equation (2.324) has to be minimized. (Here, for simplicity, the argument $\theta$ to the $r$ values will be omitted.) The criterion for best fit will be least squares. Thus the integral $I$, where

$$I = \int_{\theta_1}^{\theta_2} [r - (R + x \cos\theta + y \sin\theta)]^2 \, d\theta \tag{2.324}$$

has to be minimized with respect to $R$, $x$ and $y$ respectively. This implies that

$$\left(\frac{\partial I}{\partial R}\right)_{\bar{x},\bar{y}} = 0 \qquad \left(\frac{\partial I}{\partial \bar{x}}\right)_{\bar{R},\bar{y}} = 0 \qquad \left(\frac{\partial I}{\partial \bar{y}}\right)_{\bar{R},\bar{x}} = 0. \tag{2.325}$$

Solving these equations gives the desired values for $R$, $x$ and $y$ over a limited arc $\theta_1$ to $\theta_2$. Hence the general solution for a least-squares limaçon over a partial area is given by

$$\bar{x} = \left[ A\left( \int_{\theta_1}^{\theta_2} r\cos\theta \ d\theta - B \int_{\theta_1}^{\theta_2} r d\theta \right) + C\left( \int_{\theta_1}^{\theta_2} r\sin\theta \ d\theta - D \int_{\theta_1}^{\theta_2} r d\theta \right) \right] \Big/ E$$

$$\bar{y} = \left[ F\left( \int_{\theta_1}^{\theta_2} r\sin\theta \ d\theta - D \int_{\theta_1}^{\theta_2} r d\theta \right) + C\left( \int_{\theta_1}^{\theta_2} r\cos\theta \ d\theta - B \int_{\theta_1}^{\theta_2} r d\theta \right) \right] \Big/ E \tag{2.326}$$

$$\bar{R} = \frac{1}{\theta_2 - \theta_1} \int_{\theta_1}^{\theta_2} r d\theta - \frac{\bar{x}}{\theta_2 - \theta_1} \int_{\theta_1}^{\theta_2} \cos\theta \ d\theta - \frac{\bar{y}}{\theta_2 - \theta_1} \int_{\theta_1}^{\theta_2} \sin\theta \ d\theta.$$

In this equation the constants $A$, $B$, $C$, $D$, $E$ and $F$ are as follows:

$$A = \int_{\theta_1}^{\theta_2} \sin^2\theta d\theta - \frac{1}{\theta_2 - \theta_1}\left( \int_{\theta_1}^{\theta_2} \sin\theta d\theta \right)^2$$

$$B = \frac{1}{\theta_2 - \theta_1}(\sin\theta_2 - \sin\theta_1)$$

$$C = \frac{1}{\theta_2 - \theta_1} \int_{\theta_1}^{\theta_2} \cos\theta d\theta \int_{\theta_1}^{\theta_2} \sin\theta d\theta - \int_{\theta_1}^{\theta_2} \sin\theta\cos\theta d\theta$$

$$D = \frac{1}{\theta_2 - \theta_1}(\cos\theta_1 - \cos\theta_2) \tag{2.327}$$

$$E = AF - C^2 \qquad \text{where } A, F \text{ and } C \text{ are as defined here}$$

$$F = \int_{\theta_1}^{\theta_2} \cos^2\theta d\theta - \frac{1}{\theta_2 - \theta_1}\left( \int_{\theta_1}^{\theta_2} \cos\theta d\theta \right)^2.$$

In the special case of full arc equation (2.326) reduces to

$$R = \frac{1}{2\pi} \int_0^{2\pi} r \mathrm{d}\theta \qquad \bar{x} = \frac{1}{\pi} \int_0^{2\pi} r\cos\theta \, \mathrm{d}\theta \qquad \bar{y} = \frac{1}{\pi} \int_0^{2\pi} r\sin\theta \, \mathrm{d}\theta \qquad (2.328)$$

which are the first Fourier coefficients.

In equation (2.326) the only unknowns for a given $\theta_1$ and $\theta_2$ are the three integrals

$$\int_{\theta_1}^{\theta_2} r \mathrm{d}\theta \qquad \int_{\theta_1}^{\theta_2} r\cos\theta \, \mathrm{d}\theta \qquad \int_{\theta_1}^{\theta_2} r\sin\theta \, \mathrm{d}\theta \qquad (2.329)$$

which can be made immediately available from the instrument.

The best-fit circle was first obtained by R C Spragg (see BS 2370) [111], the partial case by Whitehouse [100].

For the instrument display the constant polar term S is added spatially about O to both the raw data and the computed reference line in exactly the same way so that the calculation is not affected, that is the term cancels out from both parts within the integral equation (2.326). Note that the extent to which the assumption made is valid depends on the ratio $e/R$ which, for most practical cases, is of the order of $10^{-3}$.

Two practical examples of how the lines derived from these equations look on difficult surfaces are shown in figure 2.134.



**Figure 2.134** Best-fit reference lines for partial circumference.

Because in practice the ratio $e/R$ is so small all angular relationships have to be measured relative to the origin of the chart and not the centre of the part as seen on the graph. Also, because the initial and final angles of a partial arc or circumference will be known from the instrument's angular reference, a considerable

reduction in computation can be achieved simply by ensuring that the limits of integration are symmetrical about the centre of the arc. This change can be effected by a polar transformation of the angular datum on the instrument by

$$(\theta_1 + \theta_2)/2. \tag{2.330}$$

Thus, if $\theta_3 = (\theta_2 - \theta_1)/2$, equation (2.264) becomes

$$\bar{x} = \left( \int_{-\theta_3}^{\theta_3} r \cos\theta \, d\theta - \frac{\sin\theta_3}{\theta_3} \int_{-\theta_3}^{\theta_3} r \, d\theta \right) \left( \theta_3 + \frac{\sin2\theta_3}{2} - \frac{1}{\theta_3} + \frac{\cos2\theta_3}{\theta_3} \right)^{-1}$$

$$\bar{y} = \left( \int_{-\theta_3}^{\theta_3} r \sin\theta \, d\theta \right) \left( \theta_3 - \frac{\sin2\theta_3}{2} \right)^{-1} \tag{2.331}$$

$$\bar{R} = \frac{1}{2\theta_3} \left( \int_{-\theta_3}^{\theta_3} r \, d\theta - 2\bar{x} \, \sin\theta_3 \right).$$

Use of these equations reduces the amount of computation considerably at the expense of only a small increase in the constraints of operation.

*(a) Graphical procedure for determining the best-fit circle and centre*
For measuring out-of-roundness around a complete circumference then equation (2.332a) is the important one. Although as will be shown later this is easy to instrument, it is more laborious than the other three methods of measurement to obtain graphically. The way proposed at present is based on the observation that any $r \cos\theta$ value is an $x$ measurement off the chart and any $r \sin\theta$ value is a $y$ measurement (see figure 2.135).



**Figure 2.135** Calculation of best-fit least-squares circle.

Thus replacing the $r \cos\theta$ values and $r \sin\theta$ values and taking discrete measurements around the profile graph rather than continuously, as will be the case using analogue instrumentation, then the parameters of the best-fit circle to fit the raw data $x$, $y$ of the coordinates of the centre and the mean radius $R$ will be given by

$$\bar{x} = \frac{2}{N} \sum_{i=1}^{N} x_i \qquad \bar{y} = \frac{2}{N} \sum_{i=1}^{N} y_i \qquad \bar{R} = \frac{1}{N} \sum_{i=1}^{N} r_i. \tag{2.332}$$

Equation (2.326) gives the best-fit conditions for a partial arc which can enclose any amount of the full circle. Often it is necessary to find the unique best-fit centre of a concentric pair of circular arcs. This involves minimizing the total sum of squares of the deviations. Arcs 1 and 2 have the sum of squares $S_1$ and $S_2$:

$$S_1 = \sum_{i=1}^{M}(r - R_1 - \bar{x}\cos\theta_i - \bar{y}\sin\theta_i)^2$$

$$S_2 = \sum_{j=1}^{N}(r - R_2 - \bar{x}\cos\theta_j - \bar{y}\sin\theta_j)^2. \tag{2.333}$$

Minimizing $S_1 + S_2$ and differentiating these polar equations with respect to $x$, $y$, $R_1$ and $R_2$ gives

$$\begin{pmatrix} \sum\cos^2\theta_i + \sum\cos^2\theta_j & \sum\sin\theta_i\cos\theta_i + \sum\sin\theta_j\cos\theta_j & \sum\cos\theta_i & \sum\cos\theta_j \\ \sum\sin\theta_i\cos\theta_i + \sum\sin\theta_j\cos\theta_j & \sum\sin^2\theta_i + \sum\sin^2\theta_j & \sum\sin\theta_i & \sum\sin\theta_j \\ \sum\cos\theta_i & \sum\sin\theta_i & M & 0 \\ \sum\cos\theta_j & \sum\sin\theta_j & 0 & N \end{pmatrix}$$

$$\times \begin{pmatrix} \bar{x} \\ \bar{y} \\ R_1 \\ R_2 \end{pmatrix} = \begin{pmatrix} \sum r_{ij}\cos\theta_i + \sum r_{ij}\cos\theta_j \\ \sum r_{ij}\sin\theta_i + \sum r_{ij}\sin\theta_j \\ \sum r_{ij} \\ \sum r_{ij} \end{pmatrix}. \tag{2.334}$$

These equations are useful when data is available in the polar form. But when data is available in the Cartesian form, the other criterion, namely minimizing the deviation from the property of the conic as an example, is useful as described below. In this case the equations of the arcs are written as

$$x^2 + y^2 - ux - vy - D_1 = 0$$

$$x^2 + y^2 - ux - vy - D_2 = 0 \tag{2.334a}$$

and the total sum of the squares of the deviation from the property of the arc/conic is defined as:

$$E_s = \sum(x^2 + y^2 - ux_i - vy_i - D_1)^2 + \sum(x^2 + y^2 - ux_j - vy_j - D_2)^2. \tag{2.335}$$

Differentiating partially with respect to $u$, $v$, $D_1$ and $D_2$, then the equation in matrix form for the solution of $u$, $v$, $D_1$ and $D_2$ is given by

$$\begin{pmatrix} \sum x_i^2 + \sum x_j^2 & \sum x_i y_i + \sum x_j y_j & \sum x_i & \sum x_j \\ \sum x_i y_i + \sum x_j y_j & \sum y_i^2 + \sum y_j^2 & \sum y_i & \sum y_j \\ \sum x_i & \sum y_i & M & 0 \\ \sum x_j & \sum y_j & 0 & N \end{pmatrix} \begin{pmatrix} u \\ v \\ D_1 \\ D_2 \end{pmatrix} = \begin{pmatrix} \sum(x_i^2 + y_i^2)x_i + \sum(x_j^2 + y_j^2)x_j \\ \sum(x_i^2 + y_i^2)y_i + \sum(x_j^2 + y_j^2)y_j \\ \sum x_1^2 + \sum y_1^2 \\ \sum x_j^2 + \sum y_j^2 \end{pmatrix}. \tag{2.336}$$

Then

$$\bar{x} = u/2 \qquad\qquad \bar{y} = v/2$$

$$R_1 = \sqrt{D_1 + (u^2 + v^2/4}} \qquad R_2 = \sqrt{D_2 + (u^2 + v^2|/4}. \tag{2.337}$$

Obviously the key to solving these sorts of problems is how to make the equations linear enough for simple solution. This is usually done automatically by the choice of instrument used to get the data. The fact that a roundness instrument has been used means that the centre *a, b* is not far from the axis of rotation. If a CMM (coordinate measuring machine) had been used this would not be the case unless the centre positions were carefully arranged.

### 2.3.6.12  Lobing coefficients

Because a signal taken from a roundness instrument is necessarily periodic it is straightforward to break down the profile of the part into a Fourier series whose fundamental component corresponds to one revolution of the instrument.

  This analysis has some useful features because, whereas all the methods of numerical assessment discussed so far have been in terms of amplitude, only the Fourier series gives an opportunity to introduce a frequency factor into the assessment.

  Thus using the same terminology as before, the raw data from the instrument $\rho(\theta)$ may be expressed as

$$\rho(\theta) = R + \sum_{n=1}^{\infty} C_n \cos\theta(n\theta - \varphi_n) \tag{2.338}$$

or

$$R + \sum_{n=1}^{\infty} (a_n \cos n\theta + b_n \sin\theta\, n\theta).$$

  In practice the series of harmonics is not taken to infinity but to some number *M* deemed sufficient to represent the profile of the workpiece adequately.

  In the equation (2.339) $C_n = (a_n^2 + b_n^2)^{1/2}$ and $\varphi_n = \tan^{-1}(b_n/a_n)$ represent the amplitude of the *n*th harmonic and $\varphi_n$ the phase, that is the orientation of the harmonic in the profile relative to the angle on the instrument taken as datum.

  The coefficients are obtained from the profile $(\theta)$ by the following expressions:

$$
\begin{aligned}
R &= \frac{1}{2\pi} \int_{-\pi}^{\pi} r(\theta)\, d\theta && \text{or in digital form} && \frac{1}{N}\sum r_i(\theta) \\
a_n &= \frac{1}{\pi} \int_{-\pi}^{\pi} r(\theta)\cos n\theta\, d\theta && \text{or} && \frac{2}{N}\sum_{i=1}^{N} r_i(\theta)\cos n\theta && (2.339) \\
b_n &= \frac{1}{\pi} \int_{-\pi}^{\pi} r(\theta)\sin n\theta\, d\theta && \text{or} && \frac{2}{N}\sum_{i=1}^{N} r_i(\theta)\sin n\theta.
\end{aligned}
$$

  The coefficients described in the equations above can, to some extent, be given a mechanical interpretation which is often useful in visualizing the significance of such an analysis.

  Breaking the roundness signal down into its Fourier components is useful because it enables a typology of the signal to be formulated. This is shown in table 2.18.

  From table 2.18 it can be seen that there are a number of influences that produce the signal which is assessed for the out-of-roundness of the workpiece. The point to note is that these influences are not related directly to each other: some are genuine, some are not. It is important to know the difference.

  Although useful there are problems associated with placing such mechanical interpretations on the coefficients because it can sometimes be misleading.

An analysis of these waveforms would yield that both have a term $C_1$ and yet an engineer would only regard figure 2.136(a) as being eccentric. In fact the analysis of the profile in figure 2.136(b) shows it to contain harmonics all through the spectral range of $n = 1$ to infinity so that it could also be thought of as being elliptical and trilobe etc. The general term $C_n$ is given by

$$C_n = \frac{2A}{n\pi} \sin\left(\frac{n\alpha}{2}\right)$$ (2.340)

**Table 2.18** Roundness signal typology.

| Fourier coefficient | Cause | Effect |
|---|---|---|
| 0 | (a) Dimension of part | Tolerance/fit |
| | (b) Instrument set-up | |
| 1 | Instrument set-up | Eccentricity on graph |
| 2 | (a) Ovality of part | |
| | (b) Instrument set-up | Component tilt |
| 3 | (a) Trilobe on part | Distortion of |
| | (b) Machine tool set-up | component due to jaws of chuck clamping |
| 4–5 | (a) Unequal angle — genuine | Distortion due to |
| | (b) Equal angle — machine tool | clamping |
| 5–20 | Machine tool stiffness | Out-of-roundness |
| 20–50 | (a) Machine tool stiffness | Out-of-roundness |
| | (b) Regenerative chatter | causes vibration |
| 50–1000 | Manufacturing process signal | Out-of-roundness causes noise |



**Figure 2.136** Problems with the harmonic approach.

so that

$$C_1 = \frac{2A}{\pi} \sin\left(\alpha / 2\right).$$ (2.341)

Fortunately examples like this are not very often encountered and so confusion is not likely to arise.

One way of utilizing the spectrum in the numerical assessment of roundness has been suggested by Spragg and Whitehouse [12] who worked out an average number of undulations per revolution ($N_a$) based upon finding the centre of gravity of the spectrum, in keeping with the random process methods of roughness.

$N_a$ is given by

$$N_a = \frac{\sum_{n=m_1}^{n=m_2} n(A_n \sin n\theta + B_n \cos n\theta)}{\sum_{n=m_1}^{n=m_2} a_n \cos n\theta + b_n \sin n\theta)} \tag{2.342}$$

where no account is taken of the sign of the coefficient. The limits $m_1$ and $m_2$ refer to the bandwidth.

The advantages of this method over simply counting the number of undulations around the circumference is that there are occasions where it is difficult, if not impossible, to make a count. In particular, this situation arises if the profile has a large random element or a large number of undulations.

Consider figure 2.137($a$). It is easy to count the undulations on a four-lobed figure but for figure 2.137($b$) it is much more difficult. The average wavelength values are shown for comparison.



($a$)                       ($b$)

**Figure 2.137** Average wavelength for roundness: ($a$) average height 0.4 $\mu$m, average UPR 4; ($b$) average height 1 $\mu$m, average UPR 20.

Another similar idea has been proposed by Chien [112] who defines effectively the root mean square energy. Thus $N$, the $N$th-order wavenumber equivalent in energy to all the rest of the spectrum, is defined as

$$N = \left( \sum_{n=1}^{m} (a_n^2 + b_n^2)n^2 \Big/ \sum_{n=1}^{m} (a_n^2 + b_n^2) \right)^{1/2}. \tag{2.343}$$

### 2.3.6.13 Roundness assessment without a formal datum

The vee method has been mentioned as a screening technique for measuring round workpieces more or less as an approximation to the radial method.

The question arises as to what is really needed as an instrument. Is it necessary to have recourse to the sophisticated radial methods described above? This question breaks down into whether it is possible to measure roundness radially without an accurate datum. This in turn poses the question as to whether both can be done without introducing serious *distortions* of the *signal* which need to be compensated for. The answer is that it is possible to arrange that a certain vee probe combination can remove the need for an accurate axis of rotation. It is also possible to use a multiplicity of probes as in the case of straightness. Consider a truly circular part and its measurement. Errors in the measurement of such a part would result from spindle uncertainty.

The probe configuration has to be such that any random movement of the part is not seen. This is equivalent to making sure that a movement of the workpiece as a whole during the measurement cycle is not detected and that it does not upset the signal. Suppose that there are two probes at angles of $-\alpha$ and $\beta$ to a datum angle. If the part, considered to be a simple circle, is moved an amount $e$ at an angle $\delta$ the two probes will suffer displacements of $e\cos(\delta + \alpha)$ and $e\cos(\delta - \beta)$ assuming that $e < 10^{-2}R$, where $R$ is the workpiece radius, and providing that the movement is confined to a region, within $10^{-2}R$, about the centre of action of the probes. For an instantaneous movement not to be registered the following equation has to be satisfied:

$$e\cos(\delta + \alpha) + e\cos(\delta + \beta) = 0. \tag{2.344}$$

This can only be true if $x = \pi + \beta$, the diametral probe configuration analysed in equations (2.305) and (2.306).

To get more out-of-roundness information another probe is needed. Let this other one be at the reference angle and the two existing probes be at angles of $-\alpha + \beta$ from it as in figure 2.138. The multiprobe signals corresponding with the multiorientation method become

$$
\begin{aligned}
V_1(\theta) &= s(\theta) + e\cos\delta \\
V_2(\theta) &= s(\theta + \alpha) + e\cos(\delta + \alpha) \\
V_3(\theta) &= s(\theta - \beta) + e\cos(\delta - \beta).
\end{aligned}
\tag{2.345}
$$

where $s$ is the workpiece signal and $e$ is the eccentricity.



**Figure 2.138** Variable error removal in roundness.

To eliminate workpiece movement some combination of the probe output voltages must be chosen so as to make the terms in $e$ and $\delta$ zero. Notice that $e$ represents an error movement at a particular moment in time; it is not dependent on $\theta$, unlike $e(\theta)$. The zero condition can be met if the probes are equispaced around the circumference at 120°, but this has a serious disadvantage, which is that they totally enclose the workpiece. However, use can be made of the fact that movement of a workpiece always affects at least two probes in the opposite sense. This is a result of their spatial configuration around the part. By moving one probe through 180° and changing its direction of sensitivity by 180° the same result can be accomplished, and the component is no longer totally enclosed!

This new feature is of great importance in measuring the roundness of large inaccessible objects to a high harmonic content, *in situ*. The advantage of having a probe assembly subtending a total angle of less than $\pi$ cannot be overemphasized; it means that the probes do not have to be fitted to a split collar. Consequently the problems of register and alignment of the probes are considerably reduced. Mathematically this mechanical repositioning of the probe can be simulated by changing the harmonic weighting function from $W_1 = \exp(jn\pi) + \exp(jn\pi/3) + \exp(-jn\pi/3)$, which is obtained from having to

ensure that $e \cos \delta + e \cos(\delta + \alpha) + e \cos(\delta - \beta) = 0$, to $W_2 = 1 - \exp(jn\pi/3) - \exp(-jn\pi/3)$, which is obtained from a similar angular imposition. In weighting function terms the removal of random eccentricity errors corresponds to the formula giving a zero value for the case where $n = 1$, that is $F_c(1) = 0$. This is true for $W_1$ and $W_2$. However, there is a considerable mechanical advantage to be gained by using the second configuration. Summarizing, a very important point has emerged from this analysis, which is that multiprobe methods are subject to one more constraint than multiorientation methods. In the weighting function there has to be complete suppression for $n = 1$, otherwise the variable errors would not be removed. (This is not necessary in the multiorientation case; eccentricity errors corresponding to $n = 1$ can be removed for each data set — the data are simply centred.) This constraint in the multiprobe method is in effect a fixing of the probe angles. To minimize the harmonic suppression of this technique it is possible to use probes having different sensitivities [103]. The general weighting function $W$ then becomes

$$W = 1 - a \exp(jn\alpha) - b \exp(-jn\beta) \tag{2.346}$$

which has the same amplitude and phase characteristics as before but $x$ and $\beta$ are constrained. A point to note is that because of these constraints imposed by the different error removal methods the weighting function for the multiprobe method need not be the same as the same-order multiorientation.

One important consequence of the first-order constraint $F_c(1) = 0$ is that the zeroth-order $F_c = 0$ is not zero. Thus

$$F_c(0) = F_s(0)(1 - a - b) \neq 0. \tag{2.347}$$

The case $a + b = 1$ is not allowed because this makes $x$ and $\beta$ both zero. Mechanically, equation (2.347) means that the multiprobe method is necessarily sensitive to radius differences between parts. This apparent shortcoming can be made use of. Measuring the average value of the combination signal over one revolution gives $F_c(0)$. This can be converted to $F_c(0) = 0$. The range of this radius measurement depends only on the range of the transducers and the overall accuracy depends on how accurately any one radius has been calibrated relative to a master workpiece. This new technique is inherently more accurate than conventional methods because it is substantially independent of random errors in the measuring system. In the multiorientation case the radius terms cannot necessarily be maintained from one run to the next in sequence because of drift etc, so for every orientation the data is preprocessed to remove the zero-order radius term (and the eccentricity term).

How the probe angles are fixed in the general case can be worked out from equation (2.346) with $n=1$. Thus

$$1 = a^2 + b^2 + 2ab \cos(\alpha + \beta) \tag{2.348}$$

$$\alpha = \cos^{-1}(1 - b^2 + a^2)/2a$$
$$\beta = \cos^{-1}(1 - a^2 + b^2)/2b. \tag{2.349}$$

From these expressions some relationships can be determined. For $\alpha + \beta = \pi/2$, $a^2 + b^2 = 1$ and for $\alpha + \beta \leqslant \pi$, $a + b > 1$. Should there be any confusion, the vee block method [97] is not the same as the multiprobe method: the vee-block method is essentially a low-order lobe-measuring technique.

The signal has amplitude and phase characteristics of $A(n)$ and $\varphi(n)$:

$$A(n) = [(1 - a \cos n\alpha - b \cos n\beta)^2 + (b \sin n\beta - a \sin n\alpha)^2]^{1/2}$$
$$\varphi(n) = \tan^{-1}[(b \sin n\beta - a \sin n\alpha)/(1 - a \cos n\alpha - b \cos n\beta)]. \tag{2.350}$$

The real significance of the use of variable sensitivity in the method will become clear in the case of variable error suppression. It is interesting to note from equation (2.349) that this is the general case for three points.

If the major criterion is simply to get rid of the first harmonic caused by spindle movements, one probe and two points of contact at an angle of $x + \beta$ will in fact suffice to satisfy equations (2.344) and (2.345), that is a vee method, for example, on two skids (figure 2.139). This is simpler than the three-probe method and does not need balanced probes. However, it does not have the same flexibility as the three-probe method because $a$ and $b$ can be adjusted with respect to each other and still maintain equation (2.349). This means that the Fourier coefficient compensation equation (2.351) can be made to be much more well behaved over a wide range of $n$, so reducing numerical problems.



**Figure 2.139**  Three-point roundness with two skids.

So far, using the multiple probe technique, the out-of-roundness has been obtained by a synthesis of modified Fourier components. There are other ways. One such simple but novel method is to solve a set of linear simultaneous equations. In effect what needs to be done in the two-orientation method, for example, is to look for only that part of the signal which has moved by the angle $\alpha$. The signal which moves is identified as component out-of-roundness. The signal which remains stationary is attributed to instrument error.

Solving for the spindle and component values (here called $\mathsf{S}$) in terms of the matrix $\mathsf{M}$ and the input voltages $\mathsf{V}$

$$\mathsf{S} = \mathsf{M}^{-1}\mathsf{V}. \tag{2.351}$$

This method still suffers from exactly the same frequency suppressions as the synthesis technique. As before the effect can be reduced by making $\alpha$ small, but other problems then arise. Differences between measurements become small—the readings become correlated and the matrix inversion becomes susceptible to numerical noise. For any given $\alpha$, however, it is possible to remove the need for a matrix inversion and at the same time improve the signal-to-noise ratio. This is accomplished by repeating the shift of the specimen until a full 360° has been completed, that is having $m$ separate but equiangled orientations [88, 122]. The reduction of noise will be about $m^{-1/2}$ in RMS terms. Once this exercise has been carried out it is possible to isolate the component error from the instrument error simply by sorting the information. For example, to find the component signal it is necessary to pick one angle in the instrument reference plane and then to identify the changes in probe voltage at this angle for all the different orientations in sequence. To get instrument errors, a fixed angle on the workpiece has to be chosen instead. Before this sifting is carried out the data sets from each orientation have to be normalized. This means that the data has to be adjusted so that the eccentricity and radius are always the same. These are the two Fourier coefficients which cannot be guaranteed to be the same from one orientation to the next, because they correspond to setting-up errors and do not relate to instrument datum or component errors. Figure 2.140 shows a typical result in which a magnification of 1 million has been obtained using this method. The figure illustrates a plot of the systematic error in a typical spindle. Providing that these errors do not change in time they can be stored and offset against any subsequent data runs, therefore enabling very high magnifications to be obtained.

**Figure 2.140** Systematic error determination.

The above treatment has dealt primarily with the nature of roundness as seen by an instrument. There are other significant aspects of the part not specifically concerned with roundness but with other metrological features of the component such as concentricity, squareness, curvature, etc, and they can be evaluated and classified from data obtained with a roundness instrument. They can confuse the roundness data.

In what follows a number of these features will be identified and quantified. It will be shown that a multiplicity of ways of assessing the features all give slightly different forms depending on the nature of the assumptions made. This is particularly true of the measurement of curvature. Why these features are included here is because of the pressure to make integrated measurements of the whole component in order to cut down setting-up and calibration time. The fact that, in general, they are a different type of signal to that of roundness obviously makes characterization more difficult.

### 2.3.6.14 *Eccentricity, concentricity*

Eccentricity is simply a measurement of the difference in position between the centre of rotation of the instrument and the geometric centre of the workpiece. This is the term $e$ referred to in the text and covered extensively in [96].

Concentricity represents the roundness equivalent of taper Here the difference in the centres of circles taken from different parts of a component is the measure. Sometimes it is taken simply as $2\times$ eccentricity.

In figure 2.141 the distance $e$ represents the lack of concentricity of the two circles in the same plane, that is the eccentricity of one relative to the other. No mention is made of the difference in radius between the two circles. Such circles may be taken in different parts of a bore or shaft or inside and outside cylinders, etc.



**Figure 2.141** Concentricity determination.

Obtaining concentricity values instrumentally has been discussed elsewhere but it is desirable at this stage to mention one or two points relevant to the future discussion.

The addition of the reference circle greatly facilitates the measurement of errors of concentricity between two or more diameters. If the profile graphs are round and smooth the relationship can easily be determined by measuring the radial separation along the axis of maximum eccentricity (figure 2.142). If,

**Figure 2.142** Eccentricity assessment: eccentricity = $(M - N)/2 \times 1/$magnitude, where $M$ and $N$ are in inches or millimetres.



**Figure 2.143** Eccentricity assessment.

however, the graphs are of poor shape then it is a great advantage to have the least-squares circles which are automatically plotted on the graph as the basis for measurement (figure 2.143).

Remember that the centre positions of such circles are defined by the first harmonic of the Fourier series.

In the measurement of concentricity of cross-sections in separated planes it is first necessary to establish a reference axis aligned to the axis of rotation of the turntable. The relationship of all other cross-sections may then be compared with respect to this defined axis. The surfaces chosen to define the reference axis will depend largely on the configuration of the workpiece, but in most cases it can generally be established from either two cross-sections along the workpiece, or from one cross-section and a shoulder or end face.

If two cross-sections along the workpiece are chosen they should be cross-sections of functional surfaces (i.e. bearing surfaces), where good roundness and surface roughness quality may be expected. For the greatest possible accuracy in setting up the reference axis, the two surfaces should be as widely spaced as possible.

If the shaft has two separated bearing surfaces which happen to be the only suitably finished surfaces from which to define the reference axis, and which in themselves are to be measured for concentricity, the procedure would be to define the axis in the most widely spaced cross-sections in the two surfaces and then to measure the relationship of the required intermediate cross-sections.

Spragg has evaluated the errors that can arise when two probes are used to evaluate concentricity [104].

Measurements can sometimes be speeded up by the use of two probes. However, care must be used. If the perfect part shown is eccentric, the differential signal is $(e^2/2R)(1 - \cos 2\theta)$. There is an apparent ellipse present, so the eccentricity should be closely controlled if two probes are used (figure 2.144).

**Figure 2.144** Problems with probe positions, (*a*) pick-up styli in line with centre of rotation, (*b*) pick-up styli not on centre line; (*c*) pick-ups on same side of part but not aligned to axis of rotation.

A further geometrical problem arises if the pick-up styli are not accurately aligned to the centre of rotation (by g, say). Then the output is

$$\frac{e^2}{2R}(1 - \cos 2\theta) + \frac{2eg}{R}(1 - \cos\theta). \tag{2.352}$$

This shows that the ellipse produced by having an eccentric part is modified to a kidney shape (limaçon) by a factor *g*. If the pick-ups are on the same side of the component, lack of alignment again gives an error. This time it is

$$\frac{e}{R}(e + g)(1 - \cos\theta). \tag{2.353}$$

The cylinder illustrated in figure 2.145(*a*) is shown to have a misalignment between the axis of the outside surface and the axis of the bore. To determine the amount of misalignment it is necessary to define a reference axis from the outside surface and align this axis to the axis of rotation of the turntable in such a way that the profile

graphs from surfaces at A and B are concentric. Misalignment of the bore axis may then be measured by transferring the stylus to the bore and taking graphs at C and D, although movement of the pick-up position along the component does not in any way affect the alignment of the component relative to the selected axis.



**Figure 2.145** (*a*) Cylinder with misaligned bore; (*b*) determination of squareness.

### 2.3.6.15 *Squareness*

Squareness can be measured with a roundness instrument, for example by measuring the circles at A and B to establish a reference axis. In squareness measurement and alignment two measurements other than the test measurement are needed to establish a reference axis. For eccentricity one measurement other than the test measurement is needed.

In figure 2.145(*b*) the squareness can be found relative to axis AB by measuring the eccentricity of the graph on the shoulder C and, knowing the position at which it was made, *r*, the angular squareness $e/r$ can be found.

In all such measurements great care should be taken to ensure that the workpiece is not moved when the probe is being moved from one measuring position to another.

### 2.3.6.16 *Curvature measurement from roundness data*

So far only departures from roundness have been considered or immediate derivatives like eccentricity or squareness. There is, however, a growing need to measure many features simultaneously with one set-up. This saves time and calibration and in general it reduces errors. One such multiple measurement is that of measuring the curvature of the component at the same time as the roundness or sometimes in spite of it.

There is one major problem of estimating curvature on a workpiece by least-squares techniques and this is the fact that the normal equations are non-linear. Operational research techniques allow some leeway in this area. An ideal way of getting out of this problem is to linearize the equations by making suitable approx-

imations. In the case of the circle the limaçon approximation provides just this basic link, but it does rely on the fact that the $e/R$ ratio has to be small, therefore allowing second-order terms to be ignored. Measuring circular workpieces with a roundness instrument is, in effect, linearizing the signal mechanically.

When attempting to measure the curvature of workpieces having a restricted arc the same limitation arises. Measuring the workpiece using the formulae for instantaneous curvature can be very dangerous if there is noise present in the nature of form or roughness, because differentiation tends to enhance the noise at the expense of the long-wavelength arc making up the curve. Also, limaçon approximations rely on a natural period which cannot be assumed from the length of the wave, so two factors have to be considered from the metrology point of view: one is linearization (if possible) and the other is noise stability. Obviously there are a large number of ways in which both can be achieved. Two are given below to illustrate the different techniques.

From figure 2.146, if the coordinate axis starts at O, then the equation of the circle is

$$(a - x_i)^2 + (y_6 + b)^2 = r^2. \tag{2.354}$$



**Figure 2.146** Form and texture — integrated method.

It is required to find the best estimates of $a$, $b$ and $r$ taken from a Cartesian-coordinate-measuring system $x, y$. It is unlikely that the limaçon form for partial arcs would be suitable because there is no easy way of estimating the fundamental period. The $x$ values are not necessarily equally spaced and are subject to error, and the $y$ values are likely to be contaminated with roughness information.

Assume that the data has the correct units, that is the magnification values removed. Let the observed quantities be $X_1, X_2, \ldots, Y_1, Y_2, Y_3. \ldots$ and the values of the true curve be $x_1, x_2, \ldots, y_1, y_2, y_3. \ldots$ Let the weighting of each of the data points be $w_x$ for the $x$ and $w_y$ for the $y$ (here assumed to be the same for all the $x$ and $y$ values but not necessarily the same as each other, i.e. $w_y \neq w_x$). Let the residues between the observed and adjusted values be $U_i$ and $V_i$. Thus $U_i = X_i - x_i$ and $V_i = Y_i - y_i$.

An assumption is made that the observed values can be expressed in terms of the true values and the residuals by the first two terms of a Taylor series. This means that only first differentials are used. Thus

$$F(X_1 \ldots X_n, Y_1 \ldots Y_n, a_o b_o r_o) = F(x_1, x_2, \ldots, x_n, y_1, y_2, \ldots, y_n, a, b, r)$$

$$+ \sum_{i=1}^{N} U_i \frac{\partial F}{\partial x_i} + \sum_{i=1}^{N} V_i \frac{\partial F}{\partial y_i} + A \frac{\partial F}{\partial a} + B \frac{\partial F}{\partial b} + R \frac{\partial F}{\partial r} \tag{2.355}$$

where $F$ is some function in this case dependent on that of a circle.

(a) *Least squares*

The nature of minimization is such that

$$S = (w_x U_i^2 + w_y V_i^2)$$ (2.356)

is a minimum, that is $\Sigma w(\text{residues})^2$ is a minimum with respect to the adjusted values. In equation (2.356) if the values of $x$ are not in error then the residuals $U_i$ will be zero; this implies that the weight $w_x$ is $\infty$. Equation (2.356) will therefore involve the minimization of $S$ such that $S = (w_y U_i)$ is a minimum only and *vice versa* with $x$.

(b) *Curve fitting*

Here not only are observations involved but also estimates $a_0$, $b_0$, $r_0$ of the unknown parameters of the curve to be found.

In general, for $n$ points of data $D$ conditions or conditional equations can be used to help define the adjusted values of $x$, $y$ and $a$, $b$, $r$. Thus

$$\left.\begin{array}{l} F_1(x_1 x_2 \ldots x_n y_n; a, b, r) = 0 \\ F_2( \qquad\qquad\cdot\qquad\qquad ) = 0 \\ F_D( \qquad\qquad\cdot\qquad\qquad ) = 0 \end{array}\right\} D \quad \text{equations for } D \text{ conditions.}$$ (2.357)

$F_1$, $F_2$ are conditional functions and have to be chosen such that, when equated to zero, they force the conditions that have to be imposed on the adjusted coordinates. Note that all these would automatically be satisfied if the true values of $x_1$, $x_2$, $y_1$, $y_2$, $a$, $b$, $r$ were available.

The derivatives of these condition functions which are to be used are

$$\frac{\partial F}{\partial x_i} \qquad \frac{\partial F}{\partial y_i} \qquad \frac{\partial F}{\partial a} \qquad \frac{\partial F}{\partial b} \qquad \frac{\partial F}{\partial r}.$$ (2.358)

The function $F$ satisfying the above is obtained from equation (2.297), yielding

$$F = (a - x)^2 + (y + b)^2 - r^2 = 0.$$ (2.359)

Suitable estimates of the derivatives can be obtained by using the observed quantities and estimates of the parameters $a_0$, $b_0$, $r_0$. Thus

$$\frac{\partial F}{\partial x_i} = 2(x_i - a_0) \qquad \frac{\partial F}{\partial y_i} = 2(y_i + b_0) \qquad \frac{\partial F}{\partial a} = 2(a_0 - x_i)$$

$$\frac{\partial F}{\partial b} = 2(y_i + b_0) \qquad \frac{\partial F}{\partial r} = -2r_0.$$ (2.360)

The terminology is such that

$$\frac{\partial F}{\partial b} = 2(y_i + b_0).$$ (2.361)

A point on the estimation of $a_0$, $b_0$, and $r_0$ will be made later. The values of equation (2.360) are required in the calculation. The nearest to the true number has to be obtained at every point by using the observed values and the estimated parameters.

To force the adjusted data to lie on each point of the true curve it would be necessary for $F(x_i, y_i; a, b, r) = 0$ to be true for all $i$. Here again estimates can be used:

$$F(X_i,\ Y_i;\ a_0,\ b_0,\ r_0) \neq 0 \qquad \text{for all } i \text{ in general} \tag{2.362}$$

but the number should be close to zero in each case if possible.

Thus equation (2.360) provides $n$ equations. Using this Taylor assumption there are as many conditions as points.

Defining

$$L_i = \frac{1}{w_x}\left(\frac{\partial F}{\partial x_i}\right)^2 + \frac{1}{w_y}\left(\frac{\partial F}{\partial u_i}\right)^2 \tag{2.363}$$

which is called the reciprocal of the weight of the condition functions, when the observed values are used it is possible to write out the normal equations for curve fitting in matrix form. This method by Whitehouse [106] is based upon the Deming approach [105]:

$$
\begin{pmatrix}
L_i & 0 & 0 & 0 & \frac{\partial F_1}{\partial a} & \frac{\partial F_1}{\partial b} & \frac{\partial F_1}{\partial r} \\
 & & & & & \vdots & \vdots \\
0 & L_2 & 0 & 0 & \frac{\partial F_2}{\partial a} & \cdots & \cdots \\
 & & & & \vdots & \vdots & \vdots \\
0 & 0 & L_3 & 0 & \cdots & \cdots & \cdots \\
 & & & & \vdots & \vdots & \vdots \\
0 & 0 & 0 & L_n & \frac{\partial F_n}{\partial a} & \frac{\partial F_n}{\partial b} & \frac{\partial F_n}{\partial r} \\
\frac{\partial F_1}{\partial a} & \frac{\partial F_2}{\partial a} & \cdots & \frac{\partial F_n}{\partial a} & 0 & 0 & 0 \\
\frac{\partial F_1}{\partial b} & \frac{\partial F_2}{\partial b} & \cdots & \frac{\partial F_n}{\partial b} & \cdot & \cdot & \cdot \\
\frac{\partial F_1}{\partial r} & \frac{\partial F_2}{\partial r} & \cdots & \frac{\partial F_n}{\partial r} & 0 & 0 & 0
\end{pmatrix}
\times
\begin{pmatrix}
\lambda_1 \\ \\ \lambda_2 \\ \\ \\ \\ \lambda_n \\ A \\ B \\ R
\end{pmatrix}
=
\begin{pmatrix}
F_1 \\ \\ F_2 \\ \\ \\ \\ F_n \\ 0 \\ 0 \\ 0
\end{pmatrix}
\tag{2.364}
$$

where

$$S = \lambda_1 F_1 + \lambda_2 F_2 + \lambda_3 F_3 \quad \text{etc}$$
$$S = \sum \lambda F_i \tag{2.364a}$$

that is, the minimum squares sum $S$ is expressible in terms of the Lagrange multipliers. It is not necessary to compute them as far as $S$ is concerned, nor the residuals!

Equation (2.364) can be immediately reduced to a set of equations containing only $A$, $B$, $R$ by eliminating the Lagrangian multipliers $\lambda_1$, $\lambda_2$, etc. Thus

$$\lambda_1 = \frac{1}{L_i}\left(F_i - \frac{\partial F_i}{\partial a}A - \frac{\partial F_i}{\partial b}B - \frac{\partial F_i}{\partial c}C\right) \qquad i = 1 \ldots n. \tag{2.365}$$

Note that $a = a_0 - A$, $b = b_0 - B$, $r = r_0 - R$ and $A, B, C$ are the residuals in the parameters. Thus the final equation to be solved is

$$
\begin{pmatrix}
\sum\left(\dfrac{\partial F_i}{\partial a}\right)^2 \dfrac{1}{L_i} & \sum\dfrac{\partial F_i}{\partial a}\dfrac{\partial F_i}{\partial b}\dfrac{1}{L_i} & \sum\dfrac{\partial F_i}{\partial a}\dfrac{\partial F_i}{\partial r}\dfrac{1}{L_i} \\
\sum\dfrac{\partial F_i}{\partial b}\dfrac{\partial F_i}{\partial a}\dfrac{1}{L_i} & \sum\left(\dfrac{\partial F_i}{\partial b}\right)^2 \dfrac{1}{L_i} & \sum\dfrac{\partial F_i}{\partial b}\dfrac{\partial F_i}{\partial r}\dfrac{1}{L_i} \\
\sum\dfrac{\partial F_i}{\partial r}\dfrac{\partial F_i}{\partial a}\dfrac{1}{L_i} & \sum\dfrac{\partial F_i}{\partial r}\dfrac{\partial F_i}{\partial b}\dfrac{1}{L_i} & \sum\left(\dfrac{\partial F_i}{\partial r}\right)^2 \dfrac{1}{L_i}
\end{pmatrix}
\times
\begin{pmatrix} A \\ B \\ R \end{pmatrix}
=
\begin{pmatrix}
\sum F_i \dfrac{\partial F_i}{\partial a}\dfrac{1}{L_i} \\
\sum F_i \dfrac{\partial F_i}{\partial b}\dfrac{1}{L_i} \\
\sum F_i \dfrac{\partial F_i}{\partial r}\dfrac{1}{L_i}
\end{pmatrix}.
\tag{2.366}
$$

This equation can be expressed directly in terms of the observed values and estimated parameters. Assuming that the $x$ values are accurately known, $w_x = \infty$, $U_i = 0$ and $w_y = 1$, then

$$
L_i = \left(\frac{F}{y_i}\right)^2 = \left(\frac{F}{b}\right)^2 = 4(b_0 + y_i)^2.
\tag{2.367}
$$

If $n$ is the number of ordinates, all summations will be from 1 to $n$. Equation (2.307) becomes

$$
\begin{pmatrix}
\sum\left(\dfrac{a_0 - x_i}{b_0 + y_i}\right)^2 & \sum\left(\dfrac{a_0 - x_i}{b_0 + y_i}\right) & -r_0\sum\left(\dfrac{a_0 - x_i}{(b_0 + y_i)^2}\right) \\
\sum\left(\dfrac{a_0 - x_i}{b_0 + y_i}\right) & n & -r_0\sum\left(\dfrac{1}{(b_0 + y_i)}\right) \\
-\sum\left(\dfrac{a_0 - x_i}{(b_0 + y_i)^2}\right) & -\sum\left(\dfrac{1}{(b_0 + y_i)}\right) & -r_0\sum\left(\dfrac{1}{(b_0 + y_i)^2}\right)
\end{pmatrix}
\times
\begin{pmatrix} A \\ B \\ R \end{pmatrix}
$$

$$
=
\begin{pmatrix}
\frac{1}{2}\sum[(a_0 - x_i)^2 + (b_0 + y_i)^2 - r_0^2]\dfrac{(a_0 - x_i)}{(y_i + b_0)^2} \\
\frac{1}{2}\sum[(a_0 - x_i)^2 + (b_0 + y_i)^2 - r_0^2]\dfrac{1}{(b_0 + y_i)^2} \\
-\frac{1}{2}\sum[(a_0 - x_i)^2 + (b_0 + y_i)^2 - r_0^2]\dfrac{1}{(b_0 + y_i)^2}
\end{pmatrix}.
\tag{2.368}
$$

From this the values of $A, B, R$ are found by a simple matrix inversion. Notice that the rank is only three for a circle and four for a sphere so that the inversion is not complicated.

Having found $A, B, R$, the true parameters $a, b, r$ can then be found from $a = a_0 - A$, $b = b_0 - B$, $r = r_0 - R$.

(c) *Estimation of $a_0$, $b_0$, $r_0$*

These can be estimated very easily from the data subject to certain constraints such as $b_0 < r_0$ which follows from the nature of the mechanical set-up. Three methods are proposed to illustrate the flexibility which is usually available:

1. Selecting three points on the raw data $P_1$, $P_2$, $P_3$ at points $X_1$, $Y_1$, $X_2$, $Y_2$, $X_3$, $Y_3$ and forcing the condition equations to be true, therefore yielding three simultaneous equations from which $a_0$, $b_0$, $r_0$ can be found. Thus

$$Y_1 = [r_0^2 - (x_1 - a_0)^2]^{1/2} - b_0$$
$$Y_2 = [r_0^2 - (x_2 - a_0)^2]^{1/2} - b_0 \qquad (2.369)$$
$$Y_3 = [r_0^2 - (x_3 - a_0)^2]^{1/2} - b_0$$

from which a simple estimate of the parameters can be found. The values $P_1$, $P_2$, $P_3$ could be obtained from a smoothed part of the arc.

2. Taking the data, splitting it into three parts, finding the average points $\bar{P}_1, \bar{P}_2, \bar{P}_3$ of coordinates $(\bar{X}_1, \bar{Y}_1)$, $(\bar{X}_2, \bar{Y}_2)$, $(\bar{X}_3, \bar{Y}_3)$ and the same procedure fitted as in equation (2.379).

3. Using the spherometer formula to give estimates of $y_0$ and $r_0$ assuming the value $x_0$ is obtained (figure 2.147).



**Figure 2.147**   Pilot estimates.

This gives good estimates for small arc length or where an independent estimate of $y_0$ is available, which it usually is from the graph; $y$ is the average of the ordinates. Thus

$$y_0 = \bar{y} + \frac{1}{2nR} \sum (x - x_0)^2 \qquad (2.370)$$

where

$$R = \frac{n \sum (x_0 - x_i)^4 - \left[\sum (x_0 - x_i)^2\right]^2}{2\left[\sum y_i \sum (x_0 - x_i)^2 - n \sum y_i (x_0 - x_i)^2\right]} \qquad (2.371)$$

A number of simplifications are immediately possible because in most cases a good estimate of $x_0$ would be $L/2$ where $L$ is the length of the chart.

Also in equation (2.371) no direct use has been made of the fact that the $x$ values are evenly spaced. The use of orthogonal functions should enable some simplifications to be made.

(*d*) *Results*

Based on preliminary results it seems best to use method 2 or 3 to get a pilot estimate of the radius. Only a single iteration is necessary to get -a final value accurate to 1 in $10^3$. It seems that this value could well be improved to the required value of 1 in 5000 which is a specification aim for useful evaluation.

(*e*) *Procedure*

1. Work out estimates $a_0$, $b_0$, $r_0$ from the data.
2. Work out point for point the $F$ values from the same data:

$$F_1 = (x_1 - a_0)^2 + (y_1 + b_0)^2 - r_0^2$$
$$F_2 = (x_1 - a_0)^2 + (y_2 + b_0)^2 - r_0 \qquad (2.372)$$
etc.

3. Work out the differential values at every point of $\partial F/\partial x$, $\partial F/\partial a$, $\partial F/\partial b$, $\partial F/\partial r$ that is

$$\frac{\partial F_i}{\partial x} = 2(x_i - a_o) \qquad \text{etc.} \tag{2.373}$$

Here some simplification is possible because $\partial F/\partial r = 2r_o$ is constant.

4. Work out the L values at every point. Thus

$$
\begin{aligned}
L_i &= \left(\frac{\partial F_i}{\partial x_i} i\right)^2 \frac{1}{w_x} + \left(\frac{\partial F_i}{\partial x_i} i\right)^2 \frac{1}{w_y} \\
&= \frac{[2(x_i - a_0)]^2}{w_x} + \frac{[2(y_i + b_0)]^2}{w_y}.
\end{aligned}
\tag{2.374}
$$

Here again if the $x$ values are very accurately known then

$$w_x = \infty \quad and \quad L_i = \left(\frac{\partial F_i}{\partial y_i} i\right)^2 \frac{1}{w_y} \tag{2.375}$$

where $w_y$ can be made unity.

5. Evaluate the sums over all the values $i = 1$ to $n$ of the products for each point of the terms in the matrix, that is

$$\sum_{i=1}^{n} \frac{\partial F_i}{\partial b} \frac{\partial F_i}{\partial a} \frac{1}{L_i}. \tag{2.376}$$

Note that many of the differential terms can be removed if $L_i$ is a function of just $x$ or $y$ separately.

6. Fill up the matrix with these numbers.

7. Invert the matrix to get the parameter residues $A, B, R$.

8. Evaluate the true circle parameters.

#### (f) Conclusion

It seems that the method outlined, based on that used by Deming, avoids the difficulty found when treating curves which have non-linear coefficients such as circles. These non-linearities make evaluation very difficult using conventional least squares. This method also has the advantage of being able to take into account the accuracy of the data in the $y$ and $x$ directions.

A point to emerge is that, for accurate results, spatial sampling of the data rather than time sampling will be necessary.

It is possible to compensate to some extent for the problems associated with linear approximation in other ways than that given above. Two such a methods, which do not use the Taylor's theorem approximation are outlined below.

#### (i) Method 1

The origin of the measurements will be taken at point O in figure 2.146 where the datum line of the transducer intersects a radius normally; the datum line will be taken as the $x$ axis. The $y$ axis will be the radial line whose origin is at O and which is at an unknown distance $b$ from the centre of the component, T.

If the radius of the component is $r$, any point $(x, y)$ on the circle will satisfy the equation

$$x^2 + (y + b)^2 = r^2 = c^{-2} \tag{2.376a}$$

where $c$ here is taken to be the curvature.

There is a method [109] that gives a closed form for the solution of the radius of curvature equation. A profile is defined as a set of data points in a fixed plane that lie on the surface of the component to be measured. Let these be $(x_i, y_i)$, $i = 1$ $n$. It is assumed that the radius and coordinates of the centre of the true circle are not known. The problem is then to obtain estimates of these particular parameters. The technique used is a best-fit curve-fitting method.

The equation of a circle of radius $R$ and centre at $(a, b)$ is

$$(x - a)^2 + (y - b)^2 = R^2. \tag{2.377}$$

Because $(x_i, y_i)$ do not lie exactly on a circle, but will have some machined error incorporated into them, it follows that

$$(x_i - a)^2 + (y_i - b)^2 - R^2 = 2E_i \quad i = 1, \dots, n. \tag{2.378}$$

The criterion of best fit is to choose those values of $R$, $a$, $b$ that minimize $\sum_{i=1}^{n} E_i^2$. Rearranging (2.378) gives

$$x_i^2 + y_i^2 = R^2 - a^2 - b^2 + 2ax_i + 2by_i + 2E_i \tag{2.379}$$

or

$$z_i = c + ax_i + by_i + E_i \tag{2.380}$$

where $z_i = (x_i^2 + y_i^2)/2$ and $c = (R^2 - a^2 - b^2)/2$.

Note that this is now in the form of a linear equation and thus has the desired values of $c$, $a$, $b$ by least squares, which is the best-fit criterion.

The normal equations in the least squares become

$$\begin{pmatrix} n & \sum x_i & \sum y_i \\ \sum x_i & \sum x_i^2 & \sum x_i y_i \\ \sum y_i & \sum x_i y_i & \sum x_i^2 \end{pmatrix} \begin{pmatrix} \hat{c} \\ \hat{a} \\ \hat{b} \end{pmatrix} = \begin{pmatrix} \sum z_i \\ \sum z_i x_i \\ \sum z_i y_i \end{pmatrix} \tag{2.381}$$

which can easily be solved in the usual way to obtain the estimates of $a$, $b$ and $c$, namely $\hat{a}$, $\hat{b}$, $\hat{c}$. From $c$ an estimate of the radius $\hat{R}$ is obtained as follows:

$$\hat{R} = (2\hat{c} + \hat{a}^2 + \hat{b}^2)^{1/2} \tag{2.382}$$

Having now obtained estimates of the radius and coordinates of the centre of the 'true' circle they can be extracted from the profile to leave a residual profile, that is a profile of the surface texture and form error of the original workpiece. Thus

$$x_i^* = \hat{R} \tan^{-1}[(y_i - \hat{b})/(x_i - \hat{a})]$$
$$y_i^* = (x_i - \hat{a})^2 + (y_i - \hat{b})^2 - \hat{R} \tag{2.383}$$

where $(x_i^*, y_i^*)$ forms the new profile.

This trick of changing a non-linear problem into a linear one by considering the *squared terms* as simple linear parameters is not new. Recourse to this approach is often found in problems with roundness and cylindricity etc, so it is not an academic exercise. It is an elegant, but not entirely formal, way of getting around the problem of non-linearity often found in engineering metrology.

The advantage of this method over many other methods is the fact that it is non-iterative and the estimator is found in closed form.

This technique can be generalized to include any general conic section (i.e. ellipses, hyperbolics, parabolics, etc), but this is beyond the scope of the present section.

### 2.3.6.18 Estimation of radial slope

Out of roundness error or, more often, 'roundness' is not the only parameter which affects the performance of circular objects. One example is to be found in the ball bearing industry where noise is a problem. This has been attributed to impulsive forces on the balls inflicted by errors in the race (fig. 2.148).

Thus, if $m$ is the mass of the ball, the impulse resulting from a change in the race radius $Sr$

$$m\frac{\delta r}{\delta t} = m\frac{\delta r}{\delta \theta} \cdot \frac{\delta \theta}{\delta t} = mw\left(\frac{dr}{d\theta}\right)$$

(2.384)

the term $\frac{dr}{d\theta}$ is not a conventional roundness term but is important nevertheless. It may be that the impulsive character of the signal generated by the $\delta r$ term relates more to pressure variation i.e. acoustic noise than does error or roundness.



**Figure 2.148**

$\frac{dr}{d\theta}$ can be measured directly using a conventional roundness instrument or it can be found from the harmonic spectrum of the raceway (the lobing). Thus if $F(w)$ is the spectrum of $r(\theta)$ the roundness signal then $wF(w)$ is the spectrum of $\frac{dr}{d\theta}$

i.e.

$$\frac{dr(\theta)}{d\theta} \Rightarrow wF\left(w\right)$$

(2.385)

Hence given $F(w)$ multiply by $w$ and taking the inverse Fourier transform gives $\frac{dr}{d\theta}$ directly.

An instrumental method of determining $\frac{dr}{df}$ will be given in chapter 4. Usually the values of $\frac{dr}{d\theta}$ used are the average modulus within a given value of $\delta\theta$ (about 5° to correspond with the Hertzian elastic zone diameter on the ball) and the maximum value taken around the circumference of the race. Incidentally, the inner ring is considered to be more important than the outer race.



**Figure 2.149**   Measurement of radial slope.

One way of measuring the radial slope of the round workpiece is to use a conventional radius-measuring instrument having an accurate spindle, and to differentiate the signal emanating from the probe either by analogue or digital methods. The digital formulae for such a measurement are explained in the chapter on processing.

Another way is to measure the rate of change of radius. This is very important in noise generation in bearings. One method is to use a spatial differentiator comprising two probes (figure 2.149). These two probes will in fact provide the very simplest tangent estimate. However, it is important to remember that the workpiece has to rotate in order for all the periphery to be examined and this in itself could require the use of a precision spindle. This can be overcome to some extent as in some of the other examples of form measurement by using a combination of a probe and vee system as shown above. The vee suppresses the first-order movement or eccentricity of the workpiece relative to the centre of the probes, thereby considerably reducing irrelevant measurements of $dr/d\theta$ caused when even a perfectly circular part is eccentric to the axis of symmetry of the instrumentation.

It can be shown that the effect of local slopes on the vee faces will cause problems on the $dr/d\theta$ measured by the probes to the extent that the ratio of a given harmonic perturbation on the workpiece as seen on the vee relative to the $dr/d\theta$ measured by the probe is

$$\frac{\cos(n\beta/2)\ \sin\alpha}{\cos(\beta/2)\ \sin n\alpha} \quad \text{for} \quad n > 1. \tag{2.386}$$

It is useful to consider whether there is an effective transfer function for the system. The effect of a given harmonic from the ring on the probe output takes no account of the fact that they may be out of phase. This gives an effective transfer function TF given by

$$\text{TF} = \frac{\sin n\gamma\ \sin\alpha}{\sin\gamma\ \sin n\alpha}. \tag{2.387}$$

The theoretical validity of the method depends on the degree to which the periphery of the ring can be regarded as a good reference. This can be worked out and it can be shown that the $dr/d\theta$ error values are at most 80% out. This is not good enough for an absolute roundness instrument but certainly good enough for a screening instrument.

### 2.3.6.19 Assessment of ovality and other shapes

Ovality is the expected basic shape in workpieces generated from round stock rolled between two rollers. It is one of the commonest faults in rings such as ball races. It is, however, different from other faults in that it is most likely to disappear when the rings are pressed into a round hole or over a round shaft and when functioning. This is also true for piston rings. They are often made deliberately oval so as to be able to take the load (along the major axis) when inserted into the cylinder. On the other hand it is a bad fault in balls and rollers and leads to a two-point fit of a shaft in a hole and hence to sideways wobble.

There is, in principle, a clear distinction to be made between ovality and ellipticity. The term ovality can be applied to any more or less elongated figure and has been defined as the maximum difference in diameter between any cross-section regardless of whether the axes are at right angles. The ellipse on the other hand has to have the axes at right angles. In most cases the method of production is such that the two are synonymous and the terms ovality and ellipticity are taken as being the same.

For this reason the second harmonic of the Fourier series of the roundness profile is commonly called the ellipticity and ovality. If the maximum diametral difference is divided by 2 this gives the peak-to-valley radial height. Division by 6.28 gives the $Ra$ (CLA, AA) departure from the least-squares reference circle.

The shapes so far examined concern purely circular forms. Concentration upon this does not imply that other forms based on circles are not met with in engineering practice — they are, and often many different types of gear are based on centre symmetry, for example. Because of their extensive treatment elsewhere and their obvious prominence in the measurement of pitch, angle, etc, rather than the less important departures from true shape, they will not be included here.

Another sort of example is the class of shapes under the heading of trochoids. These are finding some use now in motor engineering. By the very nature of their use in rotary engines departures from the ideal form, especially in terms of waviness and roughness, can be important.

The shape used in the stator of such rotary engines has the form of an epi-trochoid. This is a general form of an epi-cycle, that is the locus of one point within a circle as it is rotated around another (figure 2.150).

The equation of this point in terms of $\theta$, the position of one circle centre O′ relative to the other centre O, is

$$x = (a+b)\,\cos\theta - \lambda b\,\cos\left(\frac{a+b}{b}\right)\theta$$

$$y = (a+b)\,\sin\theta - \lambda b\,\sin\left(\frac{a+b}{-b}\right)\theta$$

(2.388)



Epi-trochoid

**Figure 2.150** Measurement of Wankel motor stator.

from which the radial distance $r$ between O and a fixed point Q on a circle radius $b$ located $\lambda b$ from O is given by

$$r = [(a + b)^2 + \lambda^2 b^2 - 2\lambda b(a + b)\cos(a/b\theta)]^{1/2} \tag{2.389}$$

which reduce, when $a=2b$, to

$$x = 3b \ \cos\theta - \lambda b \ \cos 3\theta$$
$$x = 3b \ \sin\theta - \lambda b \ \sin 3\theta \tag{2.390}$$

which is the case for the Wankel stator.

It is in those regions where $y = 0$ and $x \sim b(3 - \lambda)$ (i.e. in the regions where $r$ is a minimum) that care has to be taken to ensure that waviness is excluded. It can be shown that in these regions the form can be expressed by a parabola

$$y^2 = 3bx - b(3 - \lambda) \tag{2.391}$$

from which deviations due to waviness can be found by simple instrumental means.

Similar curves can be generated by circles rolling within circles; in these cases they are hypotrochoids (figure 2.151). One of particular interest in engineering is the case as above where $a = 2b$. The locus of the point Q distance $\lambda b$ from O' is in fact an ellipse.

Problems of measuring these types of curve and even faithful data on the effect of misgeneration are not yet readily available.



Hypotrochoid

**Figure 2.151** Hypotrochoid measurement.

*Aspherics*

Aspheric surfaces are put into optical systems to correct for spherical aberration. Spherical aberration always occurs when polishing optical surfaces using random methods. This is because of the central limit theorem of statistics which asserts that the outcome of large numbers of operations, irrespective of what they are, will result in Gaussian statistics of the output variable, say $y$. So $p(y) = \dfrac{\exp(-y^2/2\sigma^2)}{\sqrt{2\pi}.\sigma}$. In the case of optics in three dimensions the $p(x, y, z)$ is of the form $\exp\left(\dfrac{x^2 + y^2 + z^2}{2\sigma^2}\right)$ from which the form of the geometry is obviously $x^2 + y^2 + z^2$, which is spherical.

Luckily the sphere is not a bad starting point for optical instruments. In certain regions the curves look similar. See Figure 2.152.

The general form for a curve which can be developed to an aspheric form is equation 2.392. For example consider the general conic equation for $Z$. By adding a power series (i.e. $A_1x + A_2x^2$) the aspheric can be generated. In practice the number of terms is about twelve but often up to 20 can be used. This may be to avoid patent problems.

Thus the general form

$$Z = \frac{(\text{shape factor}) \cdot x^3}{R + \sqrt{R^2 - (1 + K)x^2}} + A_1x + A_2 \ldots . \tag{2.392}$$

**Figure 2.152**

$R$ is the base radius and the shape factor is +1 for concave and -1 for convex curve.



**Figure 2.153** Conical shapes.

Obviously in figure 2.153 the shape factor is always −1.



**Figure 2.154** System correction – Schmidt plate.

Figure 2.155 shows the parameters which are often needed to be specified.



**Figure 2. 155** Parameters of residuals from aspheric shape.

When residuals have been calculated from best fit aspheric form the following parameters can be determined.

| | |
|---|---|
| Fig (figure) | The vertical distance between the least squares best-fit line and the profile height. |
| $R^a$ | Average absolute deviation of residuals from best-fit line. |
| $R^t$ | Maximum peak to valley error |
| $X_p$ | Distance of residual peak from aspheric axis. |
| $X_v$ | Distance of residual valley from aspheric axis |
| $S_{mx}$ | Maximum surface slope error |
| $S_{mn}$ | Mean surface slope error. |

### 2.3.6.20 *Three-dimensional measurement — sphericity*

In the same way that straightness was extended into flatness, roundness can be extended into sphericity.

Sphericity may be taken as the departure of a nominally spherical body from truly spherical shape, that is one which is defined by the relationship

$$R^2 = (x - a)^2 + (y - b)^2 + (z - c)^2. \tag{2.393}$$

Because of the inherently small size of these departures (in real workpieces such as ball bearings etc) relative to the dimension $R$ itself, the same procedure as for roundness can be used. As in waviness the lack of suitable instrumentation has inhibited much work in this subject. Some attention will be paid to this point here (figures 2.156 and 2.157).



**Figure 2.156** Coordinate system for sphericity.

**Figure 2.157** Coordinate system for sphericity.

Using the same nomenclature as in the case of roundness

$$\rho = R + \frac{xa}{R} + \frac{yb}{R} + \frac{zc}{R} \tag{2.394}$$

making the same assumptions as before. Thus equation (2.394) is in effect a three-dimensional limaçon.

Because of the nature of the object to be measured, spherical coordinates should be used ($\rho$, $\theta$, $\alpha$). Thus

$$x = \rho \cos\theta \cos\alpha \qquad \text{assuming } \frac{x}{\rho} \simeq \frac{x}{R} \qquad \frac{y}{\rho} \simeq \frac{y}{R} \qquad \frac{z}{\rho} \simeq \frac{z}{R}$$

$$y = \rho \sin\theta \cos\alpha \tag{2.395}$$

$$z = \rho \sin\alpha$$

$\theta$ corresponds to the longitude angle and $\alpha$ to the latitude. From this equation

$$\rho = R + a \cos\theta \cos\alpha + b \sin\theta \cos\alpha + c \sin\alpha \tag{2.396}$$

may be written. This is the nature of the signal to be measured.

(*a*) *Best-fit minimization*

Assuming that, in a real case, the raw data is $r(\theta, \alpha)$, then it is required to minimize the integral Z, which is given by

$$I = \int_{\alpha_1}^{\alpha_2} \int_{\theta_1}^{\theta_2} [r(\theta,\alpha) - \rho(\theta,\alpha)]^2 \, d\theta \, d\alpha$$

$$= \int_{\alpha_1}^{\alpha_2} \int_{\theta_1}^{\theta_2} [r(\theta,\alpha) - (R + a \cos\theta \cos\alpha + b \sin\theta \cos\alpha + c \sin\alpha)]^2 \tag{2.397}$$

dropping the argument of *r* and the limits of integration for simplicity. This becomes

$$I = \int\int [r - (R + a \cos\theta \cos\alpha + b \sin\theta \cos\alpha + c \sin\alpha)]^2 \, d\theta \, d\alpha$$

$$= \int\int (r^2 - 2rR - 2ra \cos\theta \cos\alpha - 2rb \sin\theta \cos\alpha - 2rc \sin\alpha$$

$$+ R^2 + 2Ra \cos\theta \cos\alpha + 2Rb \sin\theta \cos\alpha + 2Rc \sin\alpha + a^2 \cos^2\theta \cos^2\alpha$$

$$+ 2ab \cos\theta \cos\alpha \sin\theta \cos\alpha + 2ac \cos\theta \cos\alpha \sin\alpha + b^2 \sin^2\theta \cos^2\alpha$$

$$+ 2bc \sin\theta \cos\alpha \sin\alpha + c^2 \sin^2\alpha) \, d\theta \, d\alpha. \tag{2.398}$$

For best-fit minimization

$$\frac{\partial I}{\partial R} \qquad \frac{\partial I}{\partial a} \qquad \frac{\partial I}{\partial b} \qquad \frac{\partial I}{\partial c} = 0. \qquad (2.399)$$

Separately leaving off $d\theta$, $d\alpha$, rationalizing the trigonometry and letting $\theta = \theta_2 - \theta_1$ and $\alpha = \alpha_2 - \alpha_1$, then

$$-\iint r + R\theta\alpha + a\iint \cos\theta\,\cos\alpha + b\iint \sin\theta\,\cos\alpha + c\iint \sin\alpha = 0 \quad (2.400)$$

$$-\iint r\,\cos\theta\,\cos\alpha + R\iint \cos\theta\,\cos\alpha + a\iint \cos^2\theta\,\cos^2\alpha + b\iint \cos\theta\,\sin\theta\,\cos^2\alpha \qquad (2.401)$$
$$+c\iint \cos\theta\,\cos\alpha\,\sin\alpha = 0$$

$$-\iint r\,\sin\theta\,\cos\alpha + R\iint \sin\theta\,\cos\alpha + a\iint \cos\theta\,\sin\theta\,\cos^2\alpha + b\iint \sin^2\theta\,\cos^2\alpha \qquad (2.402)$$
$$+c\iint \sin\theta\,\cos\alpha\,\sin\alpha = 0$$

$$-\iint r\,\sin\alpha + R\iint \sin\alpha + a\iint \cos\theta\,\cos\alpha\,\sin\alpha + b\iint \sin\theta\,\cos\alpha\,\sin\alpha + c\iint \sin^2\alpha = 0. \quad (2.403)$$

From these four equations $R$, $a$, $b$, $c$ can be evaluated for best fit for partial spheres in limitations for either $\alpha$ or $\theta$.

(b) *Full sphere*
In the case where $\theta_1 - \theta_2 = 2\pi$ and $\alpha_2 - \alpha_1 = 2\pi$

$$a = \frac{4}{4\pi^2}\int_0^{2\pi}\int_0^{2\pi} r(\theta,\alpha)\,\cos\theta\,\cos\alpha\;d\theta\;d\alpha$$
$$b = \frac{4}{4\pi^2}\int_0^{2\pi}\int_0^{2\pi} r(\theta,\alpha)\,\sin\theta\,\cos\alpha\;d\theta\;d\alpha \qquad (2.404)$$
$$C = \frac{2}{4\pi^2}\int_0^{2\pi}\int_0^{2\pi} r(\theta,\alpha)\,\sin\theta\;d\theta\;d\alpha$$
$$R = \frac{1}{4\pi^2}\int_0^{2\pi}\int_0^{2\pi} r(\theta,\alpha)\;d\theta\;d\alpha.$$

(The numerical equivalents of these equations are derived in chapter 3.) These are the best-fit coefficients for a full sphere and could be evaluated from a point-to-point measuring system whose sensitive direction is always pointing to the common origin O (figure 2.156).
    The numerical equivalents of these equations are

$$a = 4\sum(x/N) \qquad b = 4\sum(y/N) \qquad c = 2\sum(y/N) \qquad R = \sum(r/N). \qquad (2.405)$$

After getting the centre and radius the amount of sphericity is obtained by calculating all the radii with respect to the centre and finding their maximum difference.

(*c*) *Partial sphere: case A*

Letting $\theta_2 - \theta_1 = 2\pi$, $\alpha_2 - \alpha_1 = \delta\alpha$ and mean angle $\alpha = \alpha$ (i.e. one latitude) the normal equations become, using the mean value theorem

$$-\delta\alpha \int_0^{2\pi} r(\theta, \alpha) \ d\theta + R\delta\alpha 2\pi + 2\pi c\delta\alpha \ \sin \alpha = 0 \tag{2.406}$$

or

$$-\int_0^{2\pi} r(\theta, \alpha) d\theta + 2\pi R + 2\pi c \ \sin \alpha = 0$$

$$-\delta\alpha \ \cos \alpha \int_0^{2\pi} r(\theta, \alpha) \ \cos \theta \ d\theta + \frac{a}{4} 2\pi\delta\alpha (1 + \cos 2\alpha) = 0$$

or

$$-\cos \alpha \int_0^{2\pi} r(\theta, \alpha) \ \cos \theta \ d\theta + a\pi \cos^2 \alpha = 0 \tag{2.407}$$

or

$$-\int_0^{2\pi} r(\theta, \alpha) \ \cos \theta \ d\theta + a\pi \cos \alpha = 0.$$

Hence

$$a = \frac{1}{\pi \ \cos \alpha} \int_0^{2\pi} r(\theta, \alpha) \ \cos \theta \ d\theta$$

$$b = \frac{1}{\pi \ \cos \alpha} \int_0^{2\pi} r(\theta, \alpha) \ \sin \theta \ d\theta \tag{2.408}$$

(the best value being when $\cos \alpha = 1$, i.e. $\alpha = 0$)

$$-\delta\alpha \ \sin \alpha \int_0^{2\pi} r(\theta, \alpha) \ d\theta + R\delta\alpha 2\pi \ \sin \alpha + c\delta\alpha 2\pi \ \sin(2/\alpha) = 0 \tag{2.409}$$

or

$$-\int_0^{2\pi} r(\theta, \alpha) \ d\theta + 2\pi R + c2\pi \ \sin \alpha = 0.$$

This equation (2.409) shows that it is only possible to get a true value of $R$ when $\sin \alpha = 0$, that is around an equatorial trace. Under these conditions.

$$R = \frac{1}{2\pi} \int_0^{2\pi} r(\theta, 0) \ d\theta. \tag{2.410}$$

For $\alpha = 2\pi$ and $\theta_2 - \theta_1 = \delta\theta$, mean angle $\theta$, equation (2.406)

$$-\iint r(\theta,\alpha) \ \mathrm{d}\theta \ \mathrm{d}\alpha + R\iint \mathrm{d}\theta \ \mathrm{d}\alpha + a\iint \sin \alpha \ \mathrm{d}\theta \ \mathrm{d}\alpha = 0 \qquad (2.411)$$

becomes

$$-\delta\theta \int r(\theta,\alpha) \ \mathrm{d}\alpha + R\delta\theta 2\pi = 0$$

or

$$R = \frac{1}{2\pi} \int_0^{2\pi} r(\theta,\alpha) \ \mathrm{d}\alpha.$$

Equation (2.407) becomes

$$-\delta\theta \ \cos \ \theta \int_0^{2\pi} r(\theta,\alpha) \ \cos \ \alpha \ \mathrm{d}\alpha + a\pi \ \cos^2 \ \theta \ \delta\alpha + \frac{\pi}{2} b\delta\theta \ \sin \ 2\theta = 0 \qquad (2.412)$$

$$-\int_0^{2\pi} r(\theta,\alpha) \ \cos \ \alpha \ \mathrm{d}\alpha + a\pi \ \cos \ \theta + \pi b \ \sin \ \theta = 0.$$

Equation (2.408) becomes

$$- \ \delta\theta \ \sin\theta \int_0^{2\pi} r(\theta,\alpha) \ \cos\alpha \ \mathrm{d}\alpha + a\pi \ \delta\theta \ \sin\theta \ \cos\theta + b\pi \ \sin^2\theta = 0. \qquad (2.413)$$

Equations (2.412) and (2.413) show that $a$ and $b$ cannot be determined using great circles of constant $\theta$ values:

$$(2.414)$$

$$-\iint r(\theta,\alpha) \ \mathrm{d}\theta \ \mathrm{d}\alpha \ \sin \alpha + \frac{c}{2} 2\pi\delta\theta = 0$$

$$-\delta\theta \int_0^{2\pi} r(\theta,\alpha) \ \sin \alpha \ \mathrm{d}\theta \ \mathrm{d}\alpha + \frac{c}{2} 2\pi\delta\theta = 0$$

$$c = \frac{1}{\pi} \int_0^{2\pi} r(\theta,\alpha) \ \sin \alpha \ \mathrm{d}\alpha.$$

Hence all measurements will help with the evaluation of $R$ and $c$ with equal weight, which is not the case for constant $\alpha$ and variable $\theta$, which have to be weighted; only the great circle value $\alpha = 0$ enables $R$ to be obtained.

### 2.3.6.21 *Interpretation of results of equations (2.406)—(2.414)*

Consider equations (2.406) — (2.410) and figure 2.158. The probe will move only a distance $r$ in response to a shift of $a$ in the $x$ direction. Hence all estimates of $a$ thus derived from the raw $r$ values will be smaller by the factor cos $\alpha$. In order to get the true estimate of $a$ and $b$ from lesser circles (i.e. those where $\alpha \neq 0$) a

weighting factor of $1/\cos \alpha$ has to be applied to compensate for this. For high latitudes (declinations) the factor will be so great as to make the estimate meaningless in terms of experimental error.



**Figure 2.158** Direction of probe movement.

To get the best estimates of $a$ and $b$ the average of a number of latitudes should be taken.

Similarly equations (2.410)–(2.414) show that the measurement of great circles by keeping $\theta$ constant (i.e. circles going through the poles, figure 2.159) will always give equal estimates of $c$ and $R$ but none will give estimates of $a$ and $b$. This is only possible if more than one value of $\theta$ is used. Similarly to get a measure of $R$ and $c$ from equation (2.409) at least two values of $\alpha$ should be used. For a full exposition of the partial sphericity problem refer to Murthy *et al* [111].



**Figure 2.159** Longitudinal tracks.

In a cylindrical coordinate measurement system the probe sensitivity direction is pointing the wrong way for good spherical measurement. All sorts of complications would be involved in compensation for this. To build up a composite sphere using cylindrical coordinates $(r, \theta, z)$ the value of $r$ would have to be measured very accurately for each latitude, far better in fact than is possible using normal techniques. A $1 \mu m$ resolution would not be sufficient. This problem is shown in figure 2.160.



**Figure 2.160** Cylindrical coordinates for sphericity.

Note:

A general rule for surface metrology instrumentation is that if the measuring system is matched to the component shape in terms of coordinate system, the number of wide range movements in the instrument which require high sensitivity can be reduced by one.

Measurement and checking of sphericity using a zonal technique rather than best-fit least squares is likely to produce errors in the estimation of the centre positions because it is difficult to ensure that peaks and valleys are always related.

(*a*) *Partial sphericity*

As in most engineering problems the real trouble begins when measuring the most complex shapes. Spheres are never or hardly ever complete. For this reason estimates of sphericity on partial spheres — or even full spheres — are made using three orthogonal planes. This alternative is only valid where the method of manufacture precludes localized spikes. Similarly estimates of surface deviations from an ideal spherical shape broken down in terms of deviations from ideal circles are only valid if the centres are spatially coincident — the relation between the three planes must be established somewhere! With components having very nearly a spherical shape it is usually safe to assume this if the radii of the individual circles are the same [112].

In the case of a hip prosthesis the difficult shape of the figure involves a further reorganization of the data, because it is impossible to measure complete circles in two planes. In this case the partial arc limaçon method proves to be the most suitable. Such a scheme of measurement is shown in figure 2.161.



**Figure 2.161**  Prosthetic head — partial sphere.

Similar problems can be tackled in this way for measuring spheres with flats, holes, etc, machined onto or into them.

The display of these results in such a way as to be meaningful to an inspector is difficult, but at least with the simplified technique using orthogonal planes the three or more traces can all be put onto one polar or rectilinear chart. Visually the polar chart method is perhaps the best, but if, for example, wear is being measured on prosthetic heads it is better to work directly from Cartesian records in order to measure wear (by the areas under the charts) without ambiguity.

Assessment of the out-of-sphericity value from the least-squares centre is simply a matter of evaluating the maximum and minimum values of the radial deviations of the measured data points from the calculated centre (*a, b, c*) and radius *R*.

(*b*) *Other methods*

The minimum zone method of measuring sphericity is best tackled using exchange algorithms [113]. Murthy and Abdin [114] have used an alternative approach, again iterative, using a Monte Carlo method which, although workable, is not definitive.

The measurement of sphericity highlights some of the problems that are often encountered in surface metrology, that is the difficulty of measuring a workpiece using an instrument which, even if it is not actually unsuitable, is not matched to the component shape.

If there is a substantial difference between the coordinate systems of the instrument and that of the component, artifacts can result which can mislead and even distort the signal. Sometimes the workpiece cannot be measured at all unless the instrument is modified. An example is that of measuring a spherical object with a cylindrical coordinate instrument. If the coordinate systems are completely matched then only one direction (that carrying the probe) needs to be very accurate and sensitive. All the other axes need to have adjustments sufficient only to get the workpiece within the working range of the probe. This is one reason why the CMM has many basic problems: it does not match many shapes because of its versatility, and hence all axes have to be reasonably accurate.

The problem of the mismatching of the instrument with the workpiece is often true of cylindricity measurement, as will be seen. Special care has to be taken with cylinder measurement because most engineering components have a hole somewhere which is often a critical part of the component.

### 2.3.7 Cylindricity

The cases of flatness and sphericity are naturally two-dimensional extensions from straightness and roundness, whereas cylindricity and conicity are not. They are mixtures of the circular and linear generators. The number of engineering problems which involve two rotations is small but the combination of one angular variable with one translation is very common hence the importance attached to cylindricity and, to a lesser extent, conicity. There is little defined in non-Euclidean terms. In what follows 'cylindricity' will be taken as departures from a true cylinder.

Many misconceptions surround cylindricity. Often a few acceptable measurements of roundness taken along a shaft are considered a guarantee of cylindricity. This is not true. Cylindricity is a combination of straightness and out-of-roundness. Worst of all, any method of determining cylindricity must be as independent of the measuring system as possible. Thus, tilt and eccentricity have to be catered for in the measurement frame.

To go back to the assessment of shaft roundness, the argument is that if the machine tool is capable of generating good roundness profiles at different positions along the shaft then it will also produce a straight generator at the same time. This method relies heavily on the manufacturing process and the machine tool and is not necessarily true (figure 2.162). To be sure of cylindricity capability these roundness graphs should be linked linearly together by some means independent of the reliance on manufacturing. Alternatively, linear measurements could be tied together by roundness graphs, as shown in figures 2.163(*a*) and (*b*).



**Figure 2.162** Cross-section of shaft with roundness graphs.

There is another possibility which involves the combination of *a* and *b* in Fig 2.163; to form a "cage" pattern. This has the best coverage but takes longer.



**Figure 2.163** Methods of measuring cylinders (*a*) radial section method; (*b*) generatrix method; (*c*) helical line method, (*d*) points method.

Yet another suggestion is the use of helical tracks along the cylinder (figure 2.163(*c*)). In any event some way of spatially correlating the individual measurements has to be made. Whichever one is used depends largely on the instrumentation and the ability to unravel the data errors due to lack of squareness, eccentricity, etc. Fortunately quite a number of instruments are available which work on what is in effect a cylindrical coordinate system, namely $r, z, \theta$ axes, so that usually work can be carried out on one instrument.

In those cases where component errors are large compared with the instrument accuracy specifications, for example in the squareness of the linear traverse relative to the rotational plane, the instrument itself will provide the necessary spatial correlation.

Unfortunately, from the surface metrology point of view there is still a serious problem in displaying the results obtained, let alone putting a number to them.

The biggest problem is to maintain the overall impression of the workpiece and at the same time retain as much of the finer detail as possible. The best that can be achieved is inevitably a compromise. The problem is often that shape distortions produced by tilt and eccentricity mask the expected shape of the workpiece. The instrumental set-up errors are much more important in cylindricity measurement than in sphericity measurement. For this and other reasons cylindricity is very difficult to characterize in the presence of these unrelated signals in the data.

Another problem interrelated with this is what measured feature of cylindricity is most significant for the particular application. In some cases such as in interference fits, it may be that the examination and specification of the generator properties are most important, whereas in others it may be the axis, for example in high-speed gyroshafts.

Depending on what is judged to be most suitable, the most informative method of display should be used. Because of the fundamental character of cylindrical and allied shapes in all machines these points will be investigated in some detail.

### 2.3.7.1 Methods of specifying cylindricity

As was the case in roundness, straightness, etc, so it is in cylindricity. There is a conflict amongst metrologists as to which method of assessment is best — zonal or best fit.

There is a good case for defining cylindricity as the smallest separation $c$ which can be achieved by fitting two coaxial sleeves to the deviations measured (figure 2.164($d$)). This corresponds to the minimum zone method in roundness. But other people argue that because only the outer sleeve is unique the minimum circumscribing sleeve should be used as the basis for measurement and departures should be measured inwardly from it.



**Figure 2.164** Methods of defining a cylinder: ($a$) least-squares axis (LSC); ($b$) minimum circumscribed cylinder (MCC); ($c$) maximum inscribed cylinder (MIC); ($d$) minimum zone cylinder (MZC).

Yet again there is strong argument for the use of a best-fit least-squares cylinder. Here the cylindricity would be defined as $P_1+V_1$ (remember that figures of roundness show highly magnified versions of the outer skin of the real workpiece which will be considered later). The advantage of the best-fit method is not only its uniqueness but also its great use in many other branches of engineering. Using the minimum zone or other zonal methods can give a distortion of the axis angle of the cylinder without giving a substantially false value for the cylindricity measurement. This is the case where the odd big spike (or valley) dominates the positioning of the outer (inner) sleeve. Least-squares methods would take note of this but would be unlikely to give a false axis angle. For this reason the conventional way of examining the least-squares cylinder will be dealt with shortly. It should be remembered here that it is a vastly more difficult problem than that of measuring either roundness or straightness. Interactions occur between the effect of tilt of the cylinder and the shape distortion introduced necessarily by the nature of cylinder-measuring machines — the limaçon approach. How

these interact will be considered shortly. Before this some other different methods will be considered to illustrate the difficulty of ever giving 'one-number' cylindricity on a drawing with ease.



**Figure 2.165** Cylindrical data set (Goto [103]).

Another difficulty arises when the concept of the 'referred cylinder' is used. This is the assessment of the out-of-cylindricity of the measured data from a referred form — the cylinder that is the perfect size to fit the data. The problem is that the referred cylinder has to be estimated first from the same raw data!

### 2.3.7.2 Assessment of cylindrical form

Ideally any method of assessment should isolate errors of the instrument and the setting-up procedure from those of the part itself.

One attempt at this has been carried out [103] (figure 2.165). Taking the coordinate system of the instrument as the reference axes, $\theta$ for rotation and $z$ for vertical translation they expressed the profile of a cylinder by

$$r(\theta, z) = \sum_{j=0}^{n} A_{0j} P_j(z) + \sum_{i=1}^{m} \sum_{j=0}^{n} A_{ij} P_j(z) \ \cos i\theta + B_{ij} P_j(z) \ \sin i\theta \qquad (2.415)$$

where $P$ and $A, B$ are orthogonal polynomials and Fourier coefficients respectively. The choice of function in any of the directions is governed by the shape being considered. In the circular direction the Fourier coefficients seem to be the most appropriate because of their functional significance, especially in bearings. $P$ is in the $z$ direction and $A$, $B$ are the Fourier coefficients in any circular plane, $j$ represents the order of the vertical polynomial and $i$ the harmonic of the Fourier series. $r$ denotes the observation of the deviation of the probe from its null at the $l$th direction and on the $k$th section. Thus

$$A_{ij} = \frac{\sum_{k=1}^{t} \sum_{l=1}^{u} r_k P_j(z_k) \ \cos i\theta_1}{\sum \sum P_j^2(z_k) \ \cos^2 i\theta_2}$$

$$B_{ij} = \frac{\sum_{k=1}^{t} \sum_{l=1}^{u} r_k P_j(z_k) \ \sin i\theta_1}{\sum \sum P_j^2(z_k) \ \cos^2 i\theta_2}.$$

$$(2.416)$$

The advantage of such a representation is that some geometric meaning can be allocated to individual combinations of the functions. The average radius of the whole cylinder, for instance, is taken as $(0, 0)$, and the dc term in the vertical and horizontal directions. $(0, j)$ terms represent variations in radius with $z$ and $(i, 0)$ represent geometrical components of the profile in the $z$ direction, that is the various samples of the generator.

The effect of each component of form error can be evaluated by the sum of squares

$$S_{ij} = \frac{u}{2}(A_{ij}^2 + B_{ij}^2)\left(\sum_{k=1}^{t} P_j^2(z_k)\right). \tag{2.417}$$

Taper due to the workpiece and lack of squareness of the axes is given by the coefficient $(l, 0)$ as will be seen from the polynomial series.

More complex forms are determined from the way in which the least squares polynomial coefficients change, for example the Legendre polynomial

$$
\begin{aligned}
P_0(z) &= 1 \\
P_1(z) &= z \\
P_2(z) &= \frac{3z^2}{2} - \frac{1}{2} \qquad \text{etc.}
\end{aligned}
\tag{2.418}
$$

In this way, taking higher-order polynomials often enables complex shapes to be specified.

Extending this concept has to be allowed because often the behaviour of the axis of the part has to be considered, not simply the profile (or generator). In these instances the function representing the cylinder may be better described using three arguments representing the three coefficients. Thus $F(P, A, B)$ describes the shape, where the centre coordinates for a given Z are $A_1, B_1$.

Plotting the curve representing $F(z, 1, 1)$ describes the behaviour of the $x$ and $y$ axes with height. Also, plotting $F(z, 0, 0)$, the way in which the radius changes with height, can be readily obtained (figure 2.166).



**Figure 2.166** Method of specifying cylindrical error.

The questions of what order of polynomial is likely to be met with in practice to describe the axis change with $z$ or how $R$ changes with $z$ have not been determined mainly because of the lack of availability of suitable instrumentation. This situation has now been remedied. However, it seems that these changes as a function of $z$ could be adequately covered by simple quadratic curves. Examining these curves gives some of the basic information about the three-dimensional object. However, showing these curves demonstrates the difficulty of displaying such results, as in figure 2.167.

One obvious way is to develop the shape, but this has the disadvantage that it does not retain a visual relationship to the shape of the workpiece.

As in the measurement of flatness, Lagrangian multipliers can be used to reduce error from the readings. Here, however, there are more constraining equations than in flatness.

**Figure 2.167** Development of cylinder surface.

Other more general shapes such as conicity can be equally well dealt with (figure 2.168). Exactly the same analysis can be used except for scaling factors in the angular direction.



**Figure 2.168** Conality development.

Note that, in all these analyses, the validity of the Fourier analysis relies on the fact that the signal for the circular part is a limaçon. The use of the Fourier coefficients to identify and separate out set-up errors from the genuine form errors depends on this approximation.

An example of the problems involved is shown in figure 2.169. Assessing the departure from a true cylinder or any similar body by using a single parameter in the separation of radii, as in figure 2.169, as the basis of measurement is prone to ambiguity. Figure 2.170 shows the classical shapes of taper, or conicity, bowing, concave and convex 'barrelling'. Each of the figures illustrated in figure 2.170 has the same nominal departure. 'One-number cylindricity' is obviously not sufficient to control and specify a cylindrical shape. The figures obviously have different forms and consequently different functional properties. Only by summing some type of specification, which includes, for example, the change of apparent radius with height and/or the way in which the least-squares axis changes with height, can any effective discrimination of form be achieved.

Figure 2.171 shows various combinations of the forms on a nominal cylinder.



**Figure 2.169** One-number cylindricity — minimum separation of two cylinders.

The fact that the axial deviations and the polar deviations are usually regarded as independent, at least to a first order suggests that the best way to measure a cylinder is by means of a cylindrical-based coordinate-measuring machine, or its equivalent, in which a linear reference and a rotational reference are provided at

Figure 2.170 Three types of error in cylindrical form, typical examples (*a*) axis distortion, (*b*) generatrix deviations, (*c*) cross-section form deviations.



Figure 2.171 Errors in cylindrical form basic types of deviations (*a*) axial form error; (*b*) overall shape; (*c*) radial form error; (*d*) combination of errors.

the same time. Furthermore, the argument follows that these two features should not be mixed in the measurement a set of individual roundness data linked with straightness should be obtained. This idea is generally used as shown in figure 2.163. However, the spiral method does mix them up. The advantage is purely instrumental, both drive motors for the reference movement are in continuous use, giving high accuracy if the bearings are of the hydrodymamic type and also some advantage in speed. Figures 2.172–2.176 show other definitions, such as run-out coaxiality, etc, and the effect of sampling.

Cylindricity or, more precisely, deviations from it, is much more complicated than roundness, straightness or sphericity. This is not only because it is three dimensional but also because it is defined relative to a mixture of coordinate systems, that is polar and Cartesian rather than either one or the other. Cylindricity is also much more important functionally than sphericity because of its central role in bearings and shafts in

machines. The figures highlight some of the problems. One of the most important is the realization that considerable thought has to be put into defining the axes of cylinders. Whereas the actual 'one-number' peak-to-valley deviation estimating cylindricity is not too sensitive to the choice of algorithm used to estimate the referred cylinder, the position of the best axis very definitely is. The extra large peak or valley can completely swing the direction of the axis, as seen in figure 2.172. Also, the amount of data required to cover completely the surface of a cylinder is likely to be high. It is sensible to use hundreds of data points per revolution in order to catch the spikes. This sort of cover is not required for the harmonic content of the basic shape, but it is required to find the flaw (see figure 2.176). Mismatching of shapes, such as in figure 2.173, is also a problem. Trying to fit a cone into a cylinder requires more constraint than if the data set of the test piece is nominally cylindrical. In the case shown a preferred direction would have to be indicated, for example.



**Figure 2.172** Effect of asperities. The MIC, MZC, MCC axes are very sensitive to asperities and where possible should not be used for the datum axis.



**Figure 2.173** Some factors which cause errors in cylindricity measurement.



**Figure 2.174** Total run-out, some definitions of cylindrical parameters. Total run-out is similar to 'total indicated reading' (TIR) or 'full indicated movement' (FIM) as applied to a roundness or two-dimensional figure, but in this case it is applied to the complete cylinder and is given as a radial departure of two concentric cylinders, centred on a datum axis, which totally enclose the cylinder under test.

There is also the problem of incompatible standards as shown in figure 2.175 for coaxiality. Whichever is used has to be agreed before queries arise. If in doubt the ISO standard should always be used. Again, if in doubt the least-squares best-fit algorithm should be used. This may not give the smallest value of cylindricity, neither may it be the most functional, but it is usually the most stable and consequently less susceptible to sample cover and numerical analysis problems, as seen in figure 2.176.

**Figure 2.175** Some definitions of cylindrical parameters. Coaxiality is the ability to measure cylindricity, and to set an axis allows the measurement of coaxiality and relates the behaviour of one axis relative to a datum axis.



With only 50 data points some of the surface detail is lost

**Figure 2.176** Some factors which cause errors in cylindricity measurement: the effect of insufficient data points per plane.

Many other factors are important, such as the definition of run-out and inclined cylinders, as will be seen. Many have not yet been formally defined but are unfortunately being asked for.

The proposed method for defining cylindricity described above relies on the usual reference limaçon, at least at first sight. It seems that this method was first adopted more to demonstrate the use of least squares by which the parameters can be found than for metrological reasons, but it is a method of describing the surface of which more use could probably be made. The most comprehensive discussion of the problem of cylindricity is due to Chetwynd [113]. Harmonic analysis in roundness measurement has been variously proposed, and the extension to a polynomial axis seems natural. It is, however, still restricted by the need to interpret part of the 'profile error' as caused by residual misalignment of the workpiece to the instrument coordinate system. To do this the first harmonic of each cross-section (e.g. $A$ and $B_1$) and the linear polynomial along the axis are the only terms caused by misalignment. Thus the limaçon approximation is being applied at each cross-section to account for eccentricity there and the least-squares straight line through the centres of these limaçons is taken as the tilt error between the workpiece and the instrument.

Other workers have implicitly assumed the use of a 'reference cylinder', which is in fact a limaçon on each cross-section perpendicular to the $z$ axis, with the centres of these limaçons lying on a straight line. This is true even of methods which do not actually measure such cross-sections, such as schemes using a helical trace around the workpiece. Virtually all reported work is concerned with least-squares methods. One partial exception is an attempt to discover the minimum zone cylinders from the least-squares solution. Many search methods have been proposed but are considered to be inefficient and an alternative is proposed

which uses a weighted least-squares approach, in which the weights relate to the residuals of an unweighted least-squares solution so that the major peaks and valleys are emphasized. This method is an estimation of the minimum zone cylinders rather than a solution. It still relies upon the validity of the limaçon approximation at every cross-section.

The measurement of cylindricity will be examined in more detail from the viewpoint that it will be required to produce extensions of the roundness standards and the methods are required for the solution of least-squares, minimum zone, minimum circumscribing and maximum inscribing cylinders in instrument coordinates.

### 2.3.7.3  Reference figures for cylinder measurement

None of the literature describing work on the measurement of cylinders makes use of cylindrical reference figures. Nearly always the same implicit assumption is made, namely that the cross-sectional shape in a given plane is unaltered as the alignment of the workpiece is altered. The reason for this constancy of approach probably arises from the nature of the instruments used in the measurement. In effect, they produce profiles representing sections of a cylinder on planes perpendicular to the $z$ axis of the instrument coordinate frame. The cylinder is represented by a series of circles of the same radius placed perpendicular to the $z$ axis and having their centres lying on a straight line. In practice these circles are almost inevitably approximated by limaçons (figure 2.177). The distinction of these different forms is important, so the distinct terminology used by Chetwynd will be adopted here. 'Cylinder' will be reserved strictly for describing a figure in which all cross-sections perpendicular to its axis are identical with respect to that axis. Unless specifically stated otherwise, a right circular cylinder is implied. Other cylinder-like figures which do, however, have a different geometry will be called 'cylindroids', again following the terminology of Chetwynd. Distinction is also made between 'tilt', in which all points of a figure are rotated by the same amount relative to the coordinate system, and 'skew', in which the axis of the figure is so rotated but the cross-sections remain parallel to their



Skewed circular cylinder

Tilted circular cylinder

**Figure 2.177** Tilted cylinders. Should a cylinder be tilted when measured, the result will appear as an ellipse. Therefore it is essential that the levelling of the cylinder is performed before measurement. However, it may be that if a second cylinder is being measured relative to the first (e.g. for coaxiality), relevelling is not practical (since the priority datum will be lost). In this case it is possible for the computer to correct for tilt by calculating the tilt and orientation of the axis and noting the radius of the second cylinder, and to compensate by removing the cylinder tilt ovality for each radial plane prior to performing the cylinder tilt.

Removal of the second-harmonic term or applying a 2 UPR filter is not adequate, as any true ovality in the component will also be removed.

original positions. In the case of a cylindroid, shape is described in terms of the cross-section parallel to the unskewed axis. The reference figure commonly used in cylinder measurement is then a skew limaçon cylindroid. It should be noted that, since the axis is tilted (skewed), the eccentricity at different heights will vary and so the skew limaçon cylindroid does not have a constant cross-sectional shape (figure 2.177).

An investigation of reference figures suitable for measuring cylindricity must start from a statement of the form of a true cylinder oriented arbitrarily in the space described by a set of instrument coordinates. The circular cylindrical surface is defined by the property that all its points have the same perpendicular distance (radius) from a straight line (the axis). The following analysis follows the work of Chetwynd conveniently described using direction cosines and a vector notation. The axis is fully defined by a set of direction cosines $l_1$ and a point $X_0$ through which it passes. The perpendicular distance of a general point $X$ from this line is given by

$$p = \left| X - X_0 \right| \sin \alpha \qquad (2.419)$$

where a is the angle between the axis and the line joining $X$ to $X_0$.

The direction cosines $l$ of this joining line will be

$$l = \frac{X - X_0}{[(X - X_0)^T (X - X_0)]^{1/2}} \qquad (2.420)$$

and the angle $\alpha$ is then found:

$$\cos \alpha = l_1^T l. \qquad (2.421)$$

Substituting equations (2.421) and (2.420) into (2.419) and working, for convenience, with $p^2$, since $p$ is a scalar length, gives.

$$p^2 = (X - X_0)^T (X - X_0) - ((X - X_0)^T l_1)^2. \qquad (2.422)$$

To define the cylinder, all points having $p = R_0$ are required, so a complete description is

$$R_0^2 = (X - X_0)^T (l_3 - l_1 l_1^T)(X - X_0) \qquad (2.423)$$

where $l_3$ is the three-square identity matrix.

Within the context of normal cylindricity measurement, a less generally applicable description of the cylinder can be used to give a better 'feel' to the parameters describing it. Also, experience of two-dimensional roundness measurement shows the type of operations likely to be needed on the reference (e.g. linearizations) and the forms of parameterization which are convenient to handle. It may be assumed that the axis of a cylinder being measured will not be far misaligned from the instrument $Z$ axis (i.e. the axis of the instrument spindle) and so its description in terms of deviation from that axis has advantages. In practical instruments this distance can be taken to be small. In a direct parallel with the description of eccentricity, Cartesian components of these deviations are used. The intersection of the axis with the $A = 0$ plane will be at $(A_0, B0)$ and the slopes from the $Z$ axis of the projections of the cylinder axis into the $XZ$ and $YZ$ planes will be $A_1$ and $B_1$. Any point on the axis is then defined by the coordinates $(A_0 + A_{1Z}, B_0 + B_{1Z}, Z)$. The slopes $A_1$ and $B_1$ relate simply to the direction cosines so that

$$(l_3 - l_1 l_1^T) = \frac{1}{(1 + A_1^2 + B_1^2)} \begin{pmatrix} 1 + B_1^2 & -A_1 B_1 & -A_1 \\ -A_1 B_1 & 1 + A_1^2 & -B_1 \\ -A_1 & -B_1 & A_1^2 + B_1^2 \end{pmatrix} \qquad (2.424)$$

and, on multiplying out equation (2.424), gives

$$R_0 = \frac{1}{(1 + A_1^2 + B_1^2)^{1/2}}[(X - A_0)^2(1 + B_1^2) + (Y - B_0)^2(1 + A_1^2) + Z^2(A_1^2 + B_1^2)$$
$$-2(X - A_0)(Y - B_0)A_1B_1 - 2(X - A_0)A_1Z - 2(Y - B_0)B_1Z]^{1/2}.$$

(2.425)

The conversion of equation from Cartesian to cylindrical polar coordinates gives the equation of a tilted cylinder as

$$R(\theta, Z) = \left(\frac{(A_0 + A_1Z + A_0B_1^2 - A_1B_0B_1)\cos\theta}{1} + \dots \frac{(B_0 + B_1^2 + B_0A_1^2 - A_0A_1B_1)\sin\theta}{B_1\cos\theta - A_1\sin\theta)^2}\right)$$
$$+ \frac{R_0(1 + A_1^2 + B_1^2)^{1/2}}{[1 + (B_1\cos\theta - A_1\sin\theta)^2]^{1/2}}\left(1 - \frac{[(A_0 + A_1Z)\sin\theta - (B_0 + B_1Z)\cos\theta]^2}{R_0^2[1 + (B_1\cos\theta - A_1\sin\theta)^2]}\right)^{1/2}.$$

(2.426)

In this form both the similarity to and the differences from the simple eccentric circle can be seen. The cross-section in a plane of constant $Z$ is an ellipse with minor semi-diameter $R_0$ and major semi-diameter $R_0(1 + A_1^2 + B_1^2)^{1/2}$, its major axis having the same direction as the cylinder axis projected onto the $XY$ plane. The direction of the ellipse does not correspond to the direction of eccentricity in the plane since this latter value includes the contribution of $A_0$ and $B_0$.

To allow analytical solutions to reference fitting and to allow the practical need to work with radius-suppressed data, a reference figure linear in its parameters is desired. This may be found either by the direct application of Taylor expansions (it is easier to work with equation (2.425) and then convert the result to polar coordinates) or by the removal of relatively small terms from equation (2.426) in a manner akin to the truncation of the binomial series in deriving the limaçon from the circle. The linearization of the cylinder about the point of perfect alignment ($A_0 = B_0 = A_1 = B_1 = 0$) is shown to be the skew limaçon cylindroid

$$R(\theta, Z) = (A_0 + A_1Z)\cos\theta + (B_0 + B_1Z)\sin\theta + R_0.$$

(2.427)

A comparison of equations (2.426) and (2.427) shows how much information is totally disregarded by this linearization. In particular, there is no remaining term concerned with the ellipticity of the cross-section. For small parameter values, the differences between equations (2.426) and (2.427) will be dominated by the second-order term of the power series expansion, namely

$$\frac{R_0}{2}(A_1\cos\theta + B_1\sin\theta)^2 - \frac{1}{2R_0}[(A_0 + A_1Z)\sin\theta - (B_0 + B_1Z)\cos\theta]^2.$$

(2.428)

The nature of these error terms is emphasized if they are re-expressed as

$$\frac{\tan^2\alpha R_0}{4}[1 + \cos 2(\theta - \varphi_a)] - \frac{E^2(Z)}{4R_0}[1 - \cos 2(\theta - \varphi_E(Z)]$$

(2.429)

where $\alpha$ is the angle of the axis to the $Z$ axis and $\varphi_a$ and $\varphi_E$ are the directions of tilt and total eccentricity in the $XY$ plane. The eccentricity terms $E$ and $\varphi_E$ depend upon $Z$ whereas the terms due to pure tilt do not. The acceptability of the model depends upon the maximum value of eccentricity ratio which occurs at any plane (which will be at one end of the axis length over which measurements are taken) and also upon the magnitude of the tilt compared with the absolute radius. As written above, the first term in the error can be identified with the representation of the tilted cylinder in terms of a skew circular cylindroid, while the second term relates to the approximation of the circular cross-sections of that cylindroid by limaçons.

The above discussion is naturally also of concern to the measurement of roundness profiles on cylindrical objects. It is quite common for tilt to be the major cause of eccentricity in a reading, particularly when using fixtures that cannot support the workpiece in the plane of measurement; under such conditions the phases $\varphi_a$ and $\varphi_E$ will be broadly similar so that the possible sources of second-harmonic errors reinforce each other. On the other hand, the error in the radial term could be rather smaller than would be expected simply from the limaçon approximation.

### 2.3.7.4 *Practical considerations of cylindroid references*

The development of the skew limaçon cylindroid from the cylinder is a parameter linearization. Thus the immediate consequence to measurement practice is that exactly the same assessment techniques may be used as have been used here for roundness assessment. The cylindroid's behaviour under radius suppression is exactly the same as that of the limaçon since the suppression operates in directions perpendicular to the $Z$ axis. The magnification usually associated with the translation to chart coordinates has one extra effect on the cylindroid since, generally, it would be expected that different values of magnification would be applied in the radial and axial directions. The slope of the cylindroid axis from the measurement axis will be multiplied by the ratio of the magnifications in these directions.

The shape difference between a limaçon cylindroid and a cylinder is subject to more sources of variation than that between a limaçon and a circle, but again similar methods can be used to control them. The amplitude of the second harmonic of the horizontal section through the cylinder will be, under practical measurement conditions, the effective error in that particular cross-section of the cylindroid. A worst condition for its size is that the harmonics generated by tilt and eccentricity are in phase when the combined amplitude will be $R_{0/4}(\tan^2 \alpha + \gamma^2 Z)$, $\gamma(Z)$ being the eccentricity ratio at the cross-section. Thus a quite conservative check method is to use $(\tan^2 \alpha + \gamma_{max}^2)^{1/2}$ as a control parameter in exactly the manner that $e = 0.001$ is used for roundness measurement. It should be stressed that the values of $\alpha$ likely to be encountered within current practices are very small. The total tilt adjustment on some commercially available instruments is only a few minutes of arc, so values of $\tan \alpha = 0.001$ would not be regarded as particularly small. In the majority of situations the limit on tilt will come from its effect on the allowable eccentricity: if the axial length of cylinder over which the measurement is performed is $L_0$, there must be at least one plane where the eccentricity is at least $L_{\alpha/2} \tan \alpha$, so $\gamma_{max}$ will exceed $\tan \alpha$ whenever the length of cylinder exceeds its diameter (as it may, also, if this condition is not satisfied).

The ellipticity introduced by tilting a cylinder is difficult to account for in reference figure modelling since, apart from the problems of working with a non-linear parameterization, there are other causes of elliptical cross-sections with which interactions can take place. Using, for example, best-fit 'ellipses', probably modelled by just the second harmonic of the Fourier series, on cross-sections will not usually yield information directly about tilt. This relates to the observation that, while every tilted cylinder can be described alternatively, and equivalently, as a skew elliptical cylindroid, the vast majority of elliptical cylindroids do not describe tilted circular cylinders. Given a good estimate of the cylinder axis and knowledge of the true radius, the amplitude and phase of the elliptical component can be calculated and could be used in a second stage of determining the reference.

### 2.3.7.5 *Limaçon cylindrical references*

#### (*a*) *Least-squares cylindroids conventional measurement*
The skew limaçon cylindroid is linear in its parameters and so the least-squares solution for residuals $\delta_i$,

$$\delta_i = R_i - [(A_0 + A_1 Z_i) \cos \theta + (B_0 + B_1 Z_i) \sin \theta_i + R_L] \qquad (2.430)$$

can be stated directly.

In matrix form, the parameter estimates are given by the solution of

$$
\begin{pmatrix}
\sum \cos^2\theta & \sum \sin\theta\cos\theta & \sum Z\cos^2\theta & \sum Z\sin\theta\cos\theta & \sum \cos\theta \\
\sum \sin\theta\cos\theta & \sum \sin^2\theta & \sum Z\sin\theta\cos\theta & \sum Z\sin^2\theta & \sum \sin\theta \\
\sum Z\cos^2\theta & \sum Z\sin\theta\cos\theta & \sum Z^2\cos^2\theta & \sum Z^2\sin\theta\cos\theta & \sum Z\cos\theta \\
\sum Z\sin\theta\cos\theta & \sum Z\sin^2\theta & \sum Z^2\sin\theta\cos\theta & \sum Z^2\sin^2\theta & \sum Z\sin\theta \\
\sum \cos\theta & \sum \sin\theta & \sum Z\cos\theta & \sum Z\sin\theta & M
\end{pmatrix}
\times
\begin{pmatrix}
A_0 \\ B_0 \\ A_1 \\ B_1 \\ R_1
\end{pmatrix}
$$

$$
=
\begin{pmatrix}
\sum R\cos\theta \\
\sum R\sin\theta \\
\sum RZ\cos\theta \\
\sum RZ\sin\theta \\
\sum R
\end{pmatrix}
\tag{2.431}
$$

where, to save space, indices have been omitted: $R$, $\theta$ and $Z$ all have subscript $i$ and all summations are over $i = 1$ to $N$.

The added complexity of the three-dimensional problem means that there is even higher motivation than with the simple limaçon for choosing measurement schemes that allow simplication of the coefficient matrix. This is unlikely to be possible on incomplete surfaces and so only full cylindrical surfaces will be considered. For these it is probable that a sampling scheme having a high degree of uniformity would be used for instrumental as well as arithmetic convenience. Since, also, on a roundness-measuring instrument it is normally advisable for best accuracy to keep the spindle rotating constantly throughout the measurement, two patterns of measurement are suggested: a series of cross-sections at predetermined heights $Z_i$ or a helical traverse.

If a series of cross-sections is used and each sampled identically, the summations over all the data in equation (2.431) can be replaced by a double summation over the points in each plane and the number of planes, for example

$$
\sum_{i=1}^{N} Z_i \cos\theta_i = \sum_{k=1}^{m} Z_k \sum_{j=1}^{n} \cos\theta_{jk}
\tag{2.432}
$$

where there are $m$ sections each of $n$ points, $mn = N$. Now, if the sum over $j$ satisfies the fourfold symmetry identified earlier for the simplification of the least-squares limaçon solution, at each plane the summations over $\cos\theta$, $\sin\theta$ and $\sin\theta\cos\theta$ will be zero and so also will the sums of these terms over all the planes. The matrix of coefficients then becomes quite sparse:

$$
\begin{pmatrix}
\sum\sum \cos^2\theta & 0 & \sum Z\sum \cos^2\theta & 0 & 0 \\
0 & \sum\sum \sin^2\theta & 0 & \sum Z\sum \sin^2\theta & 0 \\
\sum Z\sum \cos^2\theta & 0 & \sum Z\sum \cos^2\theta & 0 & 0 \\
0 & \sum Z\sum \sin^2\theta & 0 & \sum Z^2\sum \sin^2\theta & 0 \\
0 & 0 & 0 & 0 & mn
\end{pmatrix}
\tag{2.433}
$$

Noting that those terms involving $\cos^2\theta$ correspond to $A_0$ and $A_1$ and, similarly, $\sin^2\theta$ to $B_0$ and $B_1$, further interpretation of this matrix is possible. The radius of the least-squares limaçon cylindroid is the mean value of all the radial data points and its axis is the least-squares straight line through the centres of the least-squares limaçons on the cross-sectional planes.

The measurement scheme has, apart from computational simplicity, two advantageous features: the information on both axial straightness and cylindricity measurements is produced simultaneously and, depending upon exactly what is to be measured, there is considerable scope for data reduction during the course of the measurement.

There are other ways of selecting measuring schemes which lead to simplifications similar to, but not as complete as, the above when using measurement in cross-section. No details of them will be given here.

### (b) Least-squares cylindroids helical measurement

The helical traverse method is attractive from an instrumentation point of view. However, computationally, it loses the advantage of having $Z$ and $\theta$ independent and so evaluation must be over the whole data set in one operation. It would be expected that samples would be taken at equal increments of $\theta$ and, since $Z$ depends linearly on $\theta$, this allows various schemes for simplifying equation (2.433) quite considerably. Only one scheme will be discussed here. If it can be arranged that the total traverse encompasses an exact even number of revolutions of the workpiece and that there is a multiple of four samples in every revolution, then defining the origin such that $Z = 0$ at the mid-point of the traverse will cause all summations of odd functions of $Z$ and $\theta$ to be zero, as will all those in simply $\sin\theta$, $\cos\theta$ or $\sin\theta \cos\theta$. The coefficient matrix in equation (2.431) then becomes

$$
\begin{pmatrix}
\sum \cos^2\theta & 0 & 0 & \sum Z \sin\theta \cos\theta & 0 \\
0 & \sum \sin^2\theta & \sum Z \sin\theta \cos\theta & 0 & 0 \\
0 & \sum Z \sin\theta \cos\theta & \sum Z^2 \cos^2\theta & 0 & 0 \\
\sum Z \sin\theta \cos\theta & 0 & 0 & \sum Z^2 \sin^2\theta & \sum Z \sin\theta \\
0 & 0 & 0 & \sum Z \sin\theta & N
\end{pmatrix}.
\tag{2.434}
$$

The original set of five simultaneous equations is therefore reduced to a set of two and a set of three, with considerable computational saving.

One failing of the helical traverse relative to the measurement of cross-sections is that no information relating directly to axial straightness is produced. Overall, it would seem that there need to be fairly strong instrumental reasons for a helical traverse to be used, particularly as there would appear to be more types of surface discontinuity that can be excluded from the measurement by the judicious choice of cross-sectional heights than from the choice of helix pitch.

A point concerning coordinates should be made here. Given only the provision that the $Z$-axis scaling is unchanged, the cylindroid parameters can be used in chart or instrument coordinates by applying magnification and suppressed radius in the normal way.

One property of the limaçon fit which does not apply to the cylindroid is the observation that the estimate for the centre is exact. Reference to equation (2.431) reveals that there are additional terms that contribute slightly to the odd harmonics in the case of the cylindroid. Taking the second-order term of the binomial expansion of the first part of the equation suggests that the fundamental is changed only by about $1 + \tan^2(\alpha/4)$: 1 so that the estimate of the axis from the cylindroid should still be good in practice. This is, however, a further warning that there is a greater degree of approximation between cylinder and cylindroid than between circle and limaçon. Although still a good approximation, the cylindroid can stand rather less abuse than the simpler situations.

This section has been concerned with the definition and expression of the cylinder form as seen realistically from a typical instrument measurement point of view. Least-squares methods have been examined and the various linearizations necessary have been described. The use of zonal methods, that is the minimum zone cylinder, maximum inscribed cylinder and minimum circumscribed cylinder, to put a 'number' to the deviations from a perfect cylinder will be considered in the section on exchange techniques in chapter 3.

### 2.3.7.6 Conicity

Conicity is slightly more complicated than cylindricity in that a taper term is included. There are many ways of specifying a cone, including a point on the axis, the direction of the axis and the apex angle. Together they make an independent set of six constraints, as opposed to five in cylindricity. Thus if the limaçon (or linearization) rules apply then the equations can be written taking $n$ to be the taper. Thus a 6 x 6 matrix results,

$$
\begin{pmatrix}
- & - & - & - & EZ_i \cos \theta_i & - \\
- & - & - & - & EZ_i \sin \theta_i & - \\
- & - & - & - & EZ_i^2 \cos \theta_i & - \\
- & - & - & - & EZ_i^2 \sin \theta_i & - \\
EZ_i \cos \theta_i & EZ_i \sin \theta_c & EZ_i^2 \cos \theta_i & EZ_i^2 \sin \theta_i & EZ_i^2 & EZ_i \\
- & - & - & - & EZ_i & -
\end{pmatrix}
\times
\begin{pmatrix}
A_0 \\ B_0 \\ A_1 \\ B_1 \\ n \\ R_1
\end{pmatrix}
=
\begin{pmatrix}
\\ \\ \\ \\ ER_iZ \\
\end{pmatrix}
\tag{2.435}
$$

where the dashes are in the respective places as in equation (2.431).

Equation (2.435) represents the equation of the best-fit cone.

Despite being able to devise plausible formulae for measurement philosophy, there are still plenty of specimens which, for practical reasons, simply cannot be measured properly even today. Examples of such workpieces are long, slender, thin shafts, rolls of paper presses or the driving shafts for video recorders. They cannot be measured either owing to their physical dimensions or the accuracy required.

## 2.4    Comparison of definitions for surface metrology and coordinate-measuring machines

The definitions for geometrical features of engineering objects are not necessarily the same for surface metrology instruments and three-dimensional coordinate-measuring machines.

Traditionally the methods and equipment for the two types of measurement have been almost completely independent of one another. The result has been that definitions have developed separately. However, with today's technology the distinction between the two is becoming less clear. Coordinate-measuring machines with position sensors now have resolutions that are approaching those of surface metrology instruments and, although still less sensitive, are capable of detecting the major shape errors of surfaces. So it is imperative that in those areas where either type of instrument could be used to measure the feature, the definitions are consistent with each other.

A classic case is that of determining the position of the centre of a circular part or hole from measurements taken around the periphery, as shown in figure 2.178.

In the figure it is required to find the circle centred at $a, b$. So

$$
\sum [(x_i - a)^2 + (y_i - b)^2] = S
\tag{2.436}
$$

is a minimum, from which, minimizing $S$ with respect to $a$ and $b$ separately,

$$
a = \frac{1}{N} \sum x_i \qquad b = \frac{1}{N} \sum y_i.
\tag{2.437}
$$

The first squared term in equation (2.436) represents the residuals in the $x$ direction and the second those in the $y$ direction.

The equivalent equation in figure 2.178($b$) is

$$
\sum (r_i - \rho)^2 = S = \sum \{r_i - [R^2 - e^2 \sin^2(\theta - \varphi)]^{1/2}\}^2 = S.
\tag{2.438}
$$

**Figure 2.178** Divergence in two metrology procedures (*a*) coordinate-measuring scheme; (*b*) surface metrology scheme.

The summation term on the right-hand side has to be modified because those terms which are under the square root are not what is actually measured by a roundness machine. The data is distorted: it is not the best-fit circle but the best-fit limaçon which is being minimized; the roundness instrument does not regard a perfectly round workpiece as circular but as if it is a limaçon.

So although equation (2.438) is exactly the same as (2.436), equation (2.439) is not:

$$\sum (r_i - [R + e \, \cos(\theta - \varphi)])^2 = S. \tag{2.439}$$

This equation is more truly representative of the data picked up from the surface metrology instrument than (2.436), which is

$$\sum [r_i - (a \, \cos \theta + b \, \sin \theta + R)]^2 = S. \tag{2.440}$$

Differentiating with respect to *a* and *b* gives

$$a = \frac{2}{N} \sum r_i \cos \theta = \frac{2}{N} \sum x_i$$
$$b = \frac{2}{N} \sum r_i \sin \theta = \frac{2}{N} \sum y_i. \tag{2.441}$$

It is because of the radius suppression that the definition becomes changed. Also, in the usual form of the CMM, the signal is obtained over a relatively null angle from the centre of the coordinate system, whereas it is almost an even function for the roundness instrument. Actually both instruments are measuring the same thing and both formulae are correct. This simply means that the part looks different to the surface metrology instrument the more closely the skin of the workpiece is being examined.

This is still true if the origin of the three-dimensional coordinate-measuring machine is contained within the part geometry. In this case equation (2.438) would be the relevant equation again and the results of (2.436) would be valid. From this situation there is a transition from (2.436) to (2.441) as the degree of radius suppression increases. Considerations of symmetry in the angular positions of the data points have been considered elsewhere [90].

This example has been brought in to show that there is no inconsistency in the definitions used for the different systems. It highlights the fact that the user should be aware of the different nature of the signals being investigated. It is a direct result of the fact that, in some sorts of metrology, there is an interaction between the dimensional part (the absolute $R$ value) and the value of the surface skin itself (in this case $R - L$).

Other considerations have to be taken into account when comparing the results of 3D measuring machines and surface metrology instruments. Perhaps the most important is the fact that, often, the number and spacing of data points is different for the two approaches. Consider again the measurement of the centre position of a hole or circle: three, four or five points are usual in 3D machines, whereas a typical number in roundness measurement is 256 or 512 so that the variability of the centre is very different for the two.

### 2.4.1 Other differences

There are other causes for the differences found between surface metrology instruments and coordinate-measuring machines that are not basically concerned with the nature of the signal but with general usage.

One example of this is in the assessment of the centre position of holes. At one time this could be achieved quite satisfactorily using a three-dimensional measuring machine with a tapered plug as a probe (figure 2.179). This automatically gave an estimate of the position of the centre if pushed firmly down: of course, after this the hole could be distorted. As tolerances become tighter this type of probe is being replaced by a small touch probe which is programmed to sample the side of the hole at a limited number of places from which the centre and size of hole can be estimated. Unfortunately, this approach is only possible when the form errors are small compared with the tolerances in size and position. Figure 2.179 shows how the estimate of position can change depending on the orientation of samples taken on, for example, a three-lobed figure.

Similarly, the estimation of form error itself and straightness, roundness, etc, used to be isolated from the effects of surface roughness by the use of a suitable probe, but again this approach is now being ques-



Centre registered  Centre registered  Centre registered

**Figure 2.179** Past and present procedures: (*a*) plug gauge; (*b*) and (*c*) different sample possibilities.



**Figure 2.180** Use of mechanical integration.

tioned because of the smallness of the objects under examination. For example, in miniature precision ball bearings the track radius can be 20 times smaller than the stylus radius used to measure a conventional track.

Take as an example the effect of roughness on the assessment of out-of-roundness; this is not unlike estimation of the hole centre because it involves the positioning of a centre to start with. The effect of the roughness has, to a large extent, up to the present time, been minimized by the use of a special stylus (figure 2.180). This stylus is, in effect, a hatchet shape whose axis when in use is parallel to that of the hole or circular object [7]. The hatchet shape spans across the surface texture marks and acts as a mechanical filter. Its job is to remove, or at least reduce, the effect of surface texture. The extent to which it does depends on the ratio of the hatchet radius relative to the tool feed, as shown by Reason [96] for periodic surfaces, and the axial correlation length of the texture for random finishing processes as determined by Whitehouse [115].

In both cases this ratio has to be large in order to get effective integration. This in turn means that the physical size of the hatchet relative to the texture spacing is big. It is this factor which is causing problems because, even though workpieces are getting smaller, the average roughness spacing is not: the same class of grinding wheel is used for large and small workpieces. Because of the physical limitations on the size of the hatchet and/or ball styli there is a growing trend to use a very sharp stylus either because there is no option, such as in the case of the miniature ball bearings mentioned above, or with the intention of measuring both roughness and form at the same time. The difficulty of achieving this latter aim is only just beginning to be appreciated [107]. Failure to take account of the errors which can result can be very serious in the assessment of critical miniature workpieces which are now being made. These include gyroscopes, fuel injectors and precision bearings. The problem being addressed here is that of measuring the centre of circular workpieces and roundness, and its associated parameters, which depends critically on the accurate assessment of the position of the centre and radius of the circle from which to measure. It is not an insignificant issue because, as in most metrological investigations, errors build on each other. The position of the centre and the radius of the best-fit circle affects concentricity, sphericity, etc, as shown earlier. However, owing to the increasing tendency to use sharper styli on 3D machines differences between techniques are creeping in because of the small data sample usually taken by the 3D machine.

It can easily be shown that if the surface roughness in the circumferential direction has a correlation $\rho$ between sampled ordinates, the estimate of variance of the eccentricity is changed.

Thus if $\sigma_b^2$ is the variance using uncorrelated data and $\sigma_{bc}^2$ is that with correlated data, then

$$\sigma_{bc}^2 = \sigma_b^2/(1-\rho). \tag{2.442}$$

The degree of correlation is determined by how the data is collected. In the case of the rubber bung method the data is completely integrated and the variability $\sigma_b^2$ is close to zero. The centre is unambiguous. However, using a point probe, the sharper the stylus the less the correlation. So the question resolves itself to get as many independent data points from the feature as possible. As shown in figure 2.182 three or four points are used, which is totally inadequate if there is any considerable roundness or roughness present.

## 2.5    Characterization of defect shapes on the surface

### 2.5.1    General

There is as yet no absolute agreement as to what constitutes a flaw or defect, but the problem of defects is so important, especially as the surface roughness gets finer, that some mention of it has to be included. Very little in the way of quantification has been laid down. In fact there is still no agreement on the visual attributes.

In what follows a list of possible defects is shown, together with pictures of them, related to the German DIN 8785. It should be clear that the defects are not related to the surface roughness or waviness.

A defect is generally defined as an element, irregularity or groups of elements and irregularities of the real surface unintentionally or accidently caused during manufacture, storage or use of the surface. Such ele-

ments differ considerably from those constituting the 'roughness'. Whether this defect is desirable or undesirable depends on the function of the part.

Often, apart from generally classifying the type of defect there is another parameter which is specified. This is the number of defects per unit area on the surface. Usually this count can only take up to a certain number before the part is rejected. The number of defects is called SDN (surface defect number). Parameters like this are often given in the sheet steel or aluminium industries.

### 2.5.2 Dimensional characteristics of defects

A reference surface is usually specified onto which defect characteristics are projected. This reference surface is defined as passing through the highest peak of the real surface, excluding the defects, and is equidistant from the mean surface determined by the least-squares method.

The size of this reference surface is chosen to be large enough to assess the defect yet at the same time be small enough so as not to be influenced by form deviations — it usually coincides with the area adjacent to the defect.

The dimensional characteristics are as follows:

(1) defect length — greatest dimension of the defect;
(2) defect width — the greatest dimension across the defect measured parallel to the reference surface;
(3) defect depth — the greatest depth measured perpendicular to the reference surface;
(4) defect height — the greatest height of the defect;
(5) defect area — the area of a single defect projected onto the reference surface;
(6) defective area — the area equal to the sum of the area of individual defects.

### 2.5.3 Types of defect

Figure 2.181 shows typical recognition patterns taken as the standard characterization [124]. Also included are shrinkage holes — similar to blowholes and buckle, a depression on the surface of sheet material caused by local bending.



Groove          Scratch          Crack

Pore          Blowhole          Fissure

Cleavage          Wane          Dent

**Figure 2.181** Some types of defect: 1, groove; 2, scratch; 3, crack; 4, pore; 5, blowhole; 6, fissure, chink, crevice; 7, cleavage, flaking; 8, wane; 9, dent *(continued overleaf)*

**Figure 2.181** *(continued)* 10, scale; 11, inclusion; 12, burr; 13, raising, wart; 14, blister; 15, scoring; 16, weld edge; 17, lap; 18, deposit; 19, ship rest; 20, skidding; 21, scaring; 22, erosion; 23, corrosion; 24, pitting; 25, crazing; 26, streak; 27, discoloration.

Defects such as those shown above are relatively easy to see by eye but very difficult to assess by pattern recognition methods. Sometimes more than one type of defect is present, in which case not only do the defects have to be separately recognized but their functional importance has to be assessed.

In practice the problem of defect quantification is one which is very application conscious. About the only general quantitative measure so far in general use is the count of defects per unit area. This is often the only parameter worth using and it is usually used in a cosmetic context, where the presence or absence is the main issue rather than the type of defect.

More problems arise when defects are not easily visible. This situation occurs when an object is passed as satisfactory after initial processing, only to turn out faulty at a later stage owing to an original defect. This sort of thing happens on mirror mouldings of the bodywork of cars. It is only when the final coating is applied that, say, ripple, orange peel or bumps are observable. Under these circumstances the initial blemish (i.e. of a long wavelength) is present but masked by the surface texture. It is then a question of highlighting

the defect marks at an early stage before the expensive finishing processes have been applied. One such technique is speckle holography; another is 'D' imaging. A complementary processing technique may be the use of wavelet theory (See chapter 3[57]). These will be considered in chapter 4.

## 2.6    Summary

This chapter has been developed mainly in a chronological way to highlight the nature of surfaces and of signals obtained from instruments. From this the past, present and possible future investigations have been picked out. Because it is a dynamic scene this format should enable a framework of understanding to be built up from which future possibilities can be judged.

This chapter in many ways reflects the growth of the subject of surface metrology. Its development has never been smooth. On the contrary it has progressed very unevenly. This is largely due to the fact that the understanding of what the surface is has been almost completely determined by the ability to measure it.

The measuring techniques have grown up in a way that best optimizes the gathering of the signal. Unfortunately, in some cases, the result of this is to weaken the understanding of the surface itself, sometimes because the nature of the signal obtained from the instrument measuring it can confuse rather than clarify the situation. Roundness measurement is a good example. Understanding the nature of the signal therefore is every bit as important as understanding the nature of the surface, so comprehension of the surface is definitely a two-stage operation: first understand the signal and then understand the surface. The nature of any distortions or misleading features must be understood before there is any hope of sensibly characterizing the surface or putting a number on it. Unfortunately many researchers and users have not been aware that there is a difference. In the next chapter some of the methods of characterization will be expanded on, particularly the application of digital techniques to the processing of the surface data. This is precursive to examining ways in which the surface is examined, which is covered in chapter 4 on instrumentation. Because of the sometimes involved nature of the text, chapter 3 could be omitted in a first reading of the book — it is intended mainly for the practising instrument engineer or the student of data processing. Nanometrology does not have the same characterization problems because, often, the nature of the signal, although of the same scale of size as surface metrology, is difficult and brings its own problems. See chapters 7 and 8.

## References

[1]   Schlesinger G 1942 *Surface Finish* (London: Institute of Production Engineers)
[2]   Reason R E, Hopkins M R and Garrod R I 1944 *Report on the measurement of surface finish by stylus methods* Rank Organisation
[3]   Perthen J 1949 *Prufen und Messen der Oberflachengestalt* (Munich: Carl Hanser)
[4]   Page F S 1948 *Fine Surface Finish* (London: Chapman and Hall)
[5]   Schorsch H 1958 *Gutebestimmung an technischen Oberflachen* (Stuttgart: Wissenschaftliche)
[6]   Reason R E 1970 The measurement of surface texture *Modern Workshop Technology; Part 2* ed. Wright Baker (London: Macmillan)
[7]   Reason R E 1966 *Report on the measurement of roundness* Rank Organisation
[8]   Abbott E J and Firestone F A 1933 Specifying surface quality *Mech. Eng.* **55** 569-72
[9]   Pesante M 1963 Determination of surface roughness typology by means of amplitude density curves *Ann. CIRP* **12** 61
[10]  Ehrenreich M 1959 The slope of the bearing area as a measure of surface texture *Microtecnic* XII
[11]  Norton A E 1942 *Lubrication* (New York: McGraw-Hill)
[12]  Lukyanov V S 1962 Local irregularities of the surface and their influence on surface roughness *Wear* **5** (182)
[13]  Spragg R C and Whitehouse D J 1970/1 A new unified approach to surface metrology *P. I. Mech. Eng.* **185** 47–71
[14]  Myres N O 1962 Characterisation of surface roughness *Wear* **5** (182)
[15]  Whitehouse D J and Von Herk P 1972 A survey of reference lines for use in surface finish *Ann. CIRP* **21** 267
[16]  Whitehouse D J and Reason R E 1965 *The equation of the mean line of surface texture found by an electric wavefilter* Rank Organisation
[17]  Whitehouse D J 1967/8 An improved wavefilter for use in surface finish measurement *P. I. Mech. Eng.* **182** pt 3K 306–18
[18]  Von Weingraber H 1956 Zur definition der Oberflachenrauheit Werk Strattstechnik *Masch. Bau* 46

[19] Von Weingraber H 1957 Uber die Eignung des Hullprofils als Bezugslinie für Messung der Rauheit *Microtechnic* **11** 6–17; 1956 *Ann. CIRP* **5** 116–28

[20] Radhakrishnan V 1971 On an appropriate radius for the enveloping circle for roughness measurement in the E. System *Ann. CIRP* 20

[21] Radhakrishnan V and Shumugan M S 1974 Computation of 3D envelope for roundness *Int. J. Mach. Tool Design Research* **14** 211–6

[22] Fahl C 1982 Handschriftlicher *Ber. Verfasser* 28 (6)

[23] Reason R E 1964 The bearing parameters of surface topography Proc. *5th MTDR Conf.* (Oxford: Pergamon)

[24] Al-Salihi T 1967 *PhD Thesis* University of Birmingham Rice

[25] R O 1944 Mathematical analysis of random noise *Bell System Tech. J.* **23** 282

[26] Longuet-Higgins M S 1957 Statistical analysis of a random moving surface. *Proc. R. Soc.* **A249** 966

[27] Peklenik J 1967 Investigation of the surface typology *Ann. CIRP* **15** 381

[28] Whitehouse D J and Archard J F 1969 Properties of random processes of significance in their contact *ASME Conf. on Surface Mechanics* (*Los Angeles, 1969*) p 36–57

[29] Wormersley J R and Hopkins M R 1945 *J. Stat. Surf.* **135**

[30] Nayak P R 1971 Random process model of rough surfaces *Trans. ASME* **93** 398

[31] Linnik Y and Khusu A P 1954 Mathematico-statistical description of surface profile irregularity in grinding *Inzh. Sborn* **20** 154

[32] Nakamura T 1959 *J. S. P. M. J.* **25** 56; 1960 *J. S. P. M. J.* **26** 226

[33] Whitehouse D J 1973 Review of topography of machined surfaces *Proc. Int. Conf. on Surf. Tech., SME, Pittsburgh* pp 53–73

[34] Sankar T S and Osman M O M 1974 Profile characteristics of manufactured surfaces using random function excursion technique *ASME J. Eng. Ind. Paper* **74**–D 1–7

[35] Peklenik J and Kubo M 1968 A basic study of a three-dimensional assessment of the surface generated in a manufacturing surface *Ann. CIRP* **16** 257

[36] Kubo M and Peklenik J 1968 An analysis of microgeometric isotropy for random surface structures *Ann. CIRP* **16** 235

[37] Whitehouse D J 1983 An easy to measure average peak-valley parameter for surface finish *P. I. Mech. Eng.* **197**C 205

[38] Williamson J B P 1967/8 The microtopography of solid surfaces *P. I. Mech. Eng.* **182** pt 3K

[39] Sayles R C and Thomas T R 1976 Mapping a small area of a surface *J. Phys. E: Sci. Instrum.* **9** 855

[40] Hunter A J 1972 On the distribution of asperity heights of surfaces *B. J. Appl. Phys.* **23** 19

[41] Greenwood J A and Williamson J B P 1986 Contact of nominally flat surfaces *Proc. R. Soc.* **A295** 300

[42] Izmailov V V and Kourova M S 1980 *Wear* **59** 409–21

[43] McAdams H T, Picciano L A and Reese P A 1968 A computer method for hyposometric analysis of a braded surface *Proc. 9th MTDR Conf.* (Oxford: Pergamon) p 73

[44] Whitehouse D J and Phillips M J 1978 Discrete properties of random surfaces *Philos. Trans. R. Soc.* **290** 267–98

[45] Phillips M J 1984 *Math. Methods in Appl. Sci.* **6** 248–61

[46] Whitehouse D J and Phillips M J Two-dimensional discrete properties of random processes *Philos. Trans. R. Soc.* **A305** 441–68

[47] Lukyanov V S 1979 Evaluation of the autocorrelation functions used when investigating surface roughness *J. Mech. Eng. Sci.* **21** 105

[48] Bhushant B 1984 Prediction of surface parameters in magnetic media *Wear* **85** 19–27

[49] Greenwood J A 1984 A unified theory of surface roughness *Proc. R. Soc.* **A393** 133–57

[50] Staufert G 1979 Characterisation of random roughness profiles *Ann. CIRP* **28** 431

[51] Wirtz A 1971 *CIRP* (lay) **19** 381; *CIRP* **17** 307–15

[52] Mandelbrot B B 1977 *The Fractal Geometry of Nature* (New York: Freeman)

[53] Stout K J and Sullivan P J 1993 The developing methods for the characterisation of roughness in 3D *Phase II report vol 2* EC contract 3374/1/0/170/90/2

[54] Ledocq J H M 1977 The anisotropic texture of machine components *Proc. 18th MTDR Conf.* p 641

[55] Whitehouse D J 1978 Beta functions for surface typology *Ann. CIRP* **27** 491–7

[56] Davies N, Spedding T A and Watson W 1980 Auto regressive moving average processes with no normal residues *Time Ser. Anal.* **2** 103–9

[57] Murthy T S R, Reddy G C and Radhakrishnan V 1982 Different functions and computations for surface topography *Wear* **83** 203–14

[58] DIN 4776 (1985) Rauheitsmessung, Kenngrosen, $R_k$, $R_{pk}$, $R_{vk}$, $M_{r1}$, $M_{r2}$ zur Beschreibung des Materialanteils in Raugheitprofil

[59] Box G P and Jenkins G M 1970 *Time Series Forecasting and Control* (London: Holden-Day)

[60] Pandit S M, Nassirpour F and Wu S M 1977 *Trans. ASME* **90B** 18–24

[61] Watson W, King T G, Spedding T A and Stout K J 1979 The machined surface — time series modelling *Wear* **57** 195–205

[62] Watson W and Spedding T H 1982 Time series modelling of non-Gaussian engineering processes *Wear* **83** 215–31

[63] DeVries W R 1979 Autoregressive time series modelling for surface profiles characterization *Ann. CIRP* **28** (1)

[64] Wold H 1938 *A Study of the Analysis of Stationary Time Series* (Uppsala: Almquist and Wilksell)

[65] Yolles M I, Smith E H and Walmsley W M 1982 Walsh theory and spectral analysis of engineering surfaces *Wear* **83** 151

[66] Harmuth H F 1972 *Transmission of Information by Orthogonal Functions* (Vienna: Springer)

[67] Whitehouse D J and Zheng K G 1991 The application of Wigner functions to machine tool monitoring *Proc. Inst. Mech. Cngrs.* **206** 249–64

[68] Boulanger J 1992 The motifs methods *Int. J. Mach. Tool Manuf.* **32** 203

[69] Qian S and Chen D 1992 Orthogonal like discrete Gabor expansion *Proc. 26th Conf. on Inf. Sci. and Systems, Princeton University, March*

[70] Sayles R and Thomas T R 1978 Surface topography as a non-stationary random process *Nature* **271** 431–4

[71] Berry M V 1973 The statistical properties of echoes diffracted from rough surfaces *Philos. Trans. R. Soc.* **A273** 611

[72] Meakin P 1987 Fractal scaling in thin film condensation and material surfaces *Crit. Rev. Solid State Mater. Sci.* **13** 147

[73] Viscec T 1986 Formation of solidification material in aggregation models *Fractals in Physics* ed. Pietronero and Tossat pp 247–50

[74] Termonia Y and Meakin P 1986 Formation of fractal cracks in a kinetic fracture model *Nature* **320** 429

[75] Majundar A and Tien C L 1990 Fractal characterisation and simulation of rough surfaces *Wear* **136** 313–27

[76] Majundar A and Bhusan B 1990 Role of fractal geometry in roughness characterisation and contact mechanics of surfaces *ASME J. Tribology* **112** 205

[77] Mulvaney D J and Newland D E 1986 A characterisation of surface texture profiles *P. I. Mech. Eng.* **200** 167

[78] McCool 1984 Assessing the effect of stylus tip radius and flight on surface topography measurements *Trans. ASME* **106** 202

[79] Tlusty J and Ismail F 1981 Basic non-linearity in machining chatter *Ann. CIRP* **30 (1)** 299–304

[80] Scott P J 1989 Non linear dynamic systems in surface metrology surface topography **2** 345–66

[81] Reason R E 1966 The measurement of waviness *Proc. Conf. MTDR*

[82] Reason R E 1962 The report on reference lines for roughness and roundness *Ann. CIRP* **2** 96

[83] Shunmugam M S and Radhakrishnan 1974 2 and 3D analysis of surfaces according to E system *Proc. Inst. Mech. Cngrs.* **188(59)** 691–9

[84] Whitehouse D J 1982 Assessment errors of finishing processes caused by skid distortion *J. Phys. E: Sci. Instrum.* **15** 1337

[85] Leroy A 1972 *Ann. CIRP* Analyse statistique et restitution spectral des profile de surface techniques **22**

[86] Peklenik J 1973 envelope characteristics of machined surfaces and their functional importance *Proc. Int. Symp. Surfaces, Pittsburgh, SME*

[87] Rank Organisation 1970 The measurement of optical alignment

[88] Whitehouse D J 1976 Error separation techniques in surface metrology *J. Phys. E: Sci. Instrum.* **9** 531–6

[89] Hume K J 1951 *Engineering Metrology* (London: McDonalds)

[90] Chetwynd D J 1980 *PhD Thesis* University of Leicester

[91] Dietrich C F 1973 *Uncertainty, Calibration and Probability* (Bristol: Hilger)

[92] Miyazaki M 1965 *Bull. Jpn. Soc. Mech. Eng.* **8** 274–80

[93] Forbes A B 1989 NPL publication. Least squares best fit geometric elements *Report* DITC 140/89, April

[94] Cardon, Bouillon and Tremblay 1973 *Microtechnic* XXVI (7)

[95] Birch K G, Cox M G and Harrison M R 1973 and 1974 Calculation of the flatness of surfaces *Reports Mom 5 and Mom 9*, National Physics Laboratory, London

[96] Reason R E 1966 *Report on the measurement of roundness* Rank Organisation

[97] 1967/68 Bendix gauge *Inst. Process Mech. Eng* **182** pt 3K

[98] Steger A 1974 *PhD Thesis* University of Dresden

[99] Gyorgy H 1966 Gepalkatreszek Koralakpointtossaganak, Merese *Finomechanika,* EVF, p257

[100] Whitehouse D J 1973 A best fit reference line for use in partial arcs *J. Phys. E: Sci. Instrum.* **6** 921–4

[101] Spragg R C 1964 Methods for the assessment of departures from roundness BS 3730

[102] Chien A Y 1982 Approximate harmonic models for roundness profiles with equivalent mean square energy value *Wear* **77** 247–52

[103] Iizuka K and Goto M 1974 *Proc. Int. Conf. Prod. Eng. Res., Tokyo* **part 1** p451

[104] Spragg R C 1968 Eccentricity — a techniques for measurement of roundness and concentricity *The Engineer* Sept. p 440

[105] Deming W E 1938 *Statistical Adjustment of Data* (New York: Dover)

[106] Whitehouse D J 1973 The measurement of curved parts in the presence of roughness *Rank Organisation Tech. Rep.* T .52

[107] Phillips M J and Whitehouse D J 1977 Some theoretical aspects of the measurement of curved surfaces *J. Phys. E: Sci. Instrum.* **10** 164–9

[108] Graybill F A 1969 *Introduction to Matrices with Applications in Statistics* (Belmont, CA: Wadsworth)

[109] Scott P 1980 Private communication, Rank Organisation

[110] Whitehouse D J 1987 Radial deviation gauge *Precis. Eng.* **9** 201

[111] Murthy T S R, Rao B and Abdin S Z 1979 Evaluation of spherical surfaces *Wear* **57** 167–85

[112] Chetwynd D G and Siddall G J 1976 *J. Phys. E: Sci. Instrum.* 9 537–44, 1975 *P. I. Mech. Eng.* **4** 3–8

[113] Chetwynd D G and Phillipson P H 1980 An investigation of reference criteria used in roundness measurement *J. Phys. E: Sci. Instrum.* **13** 538

[114] Murthy T S R and Abdin S Z 1980 Minimum zone evaluation of surfaces *J. M. T. D. R.* **20** 123–36

[115] Whitehouse D J 1983 Aspects of errors introduced into the measurement of form due to miniaturisation *J. Phys. E: Sci. Instrum.* **16** 1076–80

[116] Tsukada T, Anno Y, Yanagi K and Suzuki M 1977 An evaluation of form errors of cylindrical machined parts by a spiral tracing method *Proc. 18th MTDR Conf.* pp 529–35

[117] Boudrean B D and Raja J 1992 Analysis of lay characteristics of 3D surface maps *Int. J. Mech. Tool Manuf.* 32 171

[118] Provisional French NORM as atlas of lay characteristics

[119] Whitebouse D J 1985 Assessment of surface branch profiles produced by multiple process manufacture *Proc. Inst. Mech. Cngrs.* 199B (4) 263–70

[120] Scott P ISO TC 57 WG 3 93/701187

[121] ISO 8785 Surface defects

[122] Aoki Y and Ozono S 1966 *J. Jpn Soc. Process. Eng.* 32 27–32

[123] Bendat J S 1958 Principles and applications of random noise theory. (New York, Wiley)

[124] ISO Standard Handbook Limits and Fits and Surface properties Geneva ISBN 92-67-10288-5, (1999)

[125] Whitehouse D J 1999 Identification of two functionally different peak forms in a random process. *Proc. Inst. Mech. Eng. Sc.* Vol. 213 p 303

[126] Whitehouse D J 1971 Ph. D Thesis Leicester University

[127] Whitehouse D J 1999 Some theoretical aspects of structure functions fractal parameters and related subjects. *Proc. Inst. Mech. Eng.* Part 5 **vol 213** 303

[128] Scott P J Areal topological characterization of a functional description of surfaces. European Surface Workshop. France et Silette June (1998)

[129] Maxwell J C 1870 "on hills and vales" *Phil. Mag. J. Science 4* **40** p 421

[130] Dong W P and Stout K J 1995 Two dimensional FFT and power spectrum for surface roughness in 3 Dimensions *Proc. Inst. Mech. Engrs.* **209** 81

[131] Gabor D 1946 Theory of communication *Proc. Inst. Elec. Engrs.* **93** 429

[132] Raja J 2000 Filtering of Surface profile - past present and future 1st National *Conf. Proc. Eng.* Chenai India p 99

[133] Haze Winkel M 2000 *Wavelets understand Fractals Wavelets vol 1* Ed. Koornwinder T H p 209

[134] Daubechies I 1988 Orthonormal bases of compactly supported wavelets *Com. Pure and Applied Maths* **41** 909

[135] Whitehouse D J 2001 Fractal or friction *Wear* **249** 345–353

[136] Hasegawa M, Lui J, Okuda K, Nunobiki M 1996 Calculation of the fractal dimesions of machined surface profiles *Wear* **192** 40–45

[137] Davies S and Hall P 1988 Fractal analysis of surface roughness by using spatial duty *Royal Stat. Soc.* B **61** part 1

[138] He L Zhu J 1997 The fractal character of processed metal surfaces *Wear* **208** 17–24

[139] Brown C A, Johnson W A and Rutland R M 1996 Scale sensitive fractal analysis of tuned surfaces duals *CIRP.* **45** 515

[140] Svitkin M 2002 Technomach Co. Ltd. *Roundness Technologies*

# Chapter 3
# Processing

Although the subject of engineering surfaces covers a large number of interdisciplinary subjects ranging from, say, economics on the one hand to chemistry on the other, there are a number of areas in which a degree of unity can be seen. One of these, perhaps the most important, is in the processing of measured information to enable it to be used most effectively for control of production or design purposes. Processing information is taken to mean effecting a series of changes on the original data to isolate, separate or identify features which may be of particular significance. Certain operations which keep on recurring will be highlighted. In this section a few of these will first be described and then illustrated with particular reference to some of the issues pertinent to the subject matter being discussed.

Using all or any combination of these 'processing' methods should allow a researcher, designer, or other user not only to understand more fully the background to the subject but also to progress the technique of surface analysis and measurement further.

The signal being processed is that obtained from the surface by any means. This could mean digital or analogue information obtained from a measuring instrument. It can, however, equally be a signal obtained by utilizing a mathematical model of the surface.

Because some of the contents of this chapter are concerned more with the theory of the subject rather than its routine application, it is envisaged that it will be of most use to the instrument researcher or tribologist or perhaps a student of metrology.

The first of these fundamental operations to be considered is that of numerical techniques, and the second is filtering in one form or another. These are singled out because of their special importance. Obviously the techniques have overlapping regions so that the discussion of filtering will necessarily involve some discrete concepts. However, the repetition of such important concepts is not a redundant exercise.

Filtering is examined in the same sense as that used in system vibration analysis; that is to say that the frequency characteristics of the input signal are operated on by the same method whether electrical, mechanical or computational. One reason for doing this is to introduce the important concepts of system response and convolution. These terms have wider mathematical significance than equivalent expressions used in either of the individual disciplines, for example they happen to have equivalents in random process theory. Other subjects will be included wherever deemed necessary.

The other point about the numerical analysis of random surfaces is that, to some degree, this is exploring the properties of discrete random surfaces. To this extent there is some overlap with the methods of characterizing surfaces given in chapter 2.

## 3.1  Digital methods

In what follows it should be assumed that digital methods are being used unless specified. Various optical methods will be considered in chapter 4. The present chapter ends with some graphical methods, which are useful if existing recorded data has to be checked.

Digital methods arise in two major areas: the first is in the analysis of the various waveforms, and the second is in the use of computers to control instruments, machines, etc. Here various aspects of the analysis of waveforms using digital methods will be considered; control features will be discussed in instrumentation. This topic has been singled out because of the increasing use of such methods in surface technology. The increase in use is because digital techniques tend to be more flexible and more accurate than analogue ones.

The starting point is taken to be a typical waveform which could have been obtained by any number of methods. Three basic considerations are taken into account in order to get useful digital results. These are the sampling, quantization and the numerical technique. Each is briefly considered together with typical examples. Emphasis is given to those problems unique to surface metrology. Topics such as fast Fourier transformation will be considered later.

### 3.1.1 Sampling

The operation of sampling both in time (or space) and frequency is shown in figure 3.1 in which (*a*) represents a typical waveform $z(t)$ to be sampled and (*b*) its frequency characteristic which is shown limited at a frequency *B*. In what follows the word 'time' will be synonymous with 'space'.

The process of taking a single sample of a time signal at time *t* is equivalent to multiplying the time function by a unit impulse $\delta$ at $t_1$ and integrating. Thus

$$z(t_1) = \int_{-\infty}^{\infty} z(\tau)\delta(t_1 - \tau) \, \mathrm{d}\tau \qquad (3.1)$$

where $\tau$ is a dummy variable. This operation is called the sampling property of impulses. Sampling at regular intervals in time *h* is equivalent to multiplying the waveform by an impulse train, where each impulse is separated by *h,* and then integrating. The equivalent operation in frequency is shown in figure 3.1(*b*)*,* (*d*) and (*f*). The Fourier transform of an impulse train is itself an impulse train whose spacing is 1/*h* as shown in figure 3.1(*d*). Because time sampling involves a multiplication, the equivalent operation in frequency is a convolution. Thus figure 3.1(*f*) shows the effect of convoluting the frequency impulse train with that of the frequency characteristic of the waveform (shown symmetrical about the zero-frequency axis).

The criterion for good sampling is that all information should be recoverable. From figure 3.2 it is obvious that passing this sampled signal through a low-pass filter whose low-frequency cut is higher than *B* will remove the other bands of frequency introduced by the sampling, namely *A, B, C,* etc. But this is only possible if the other bands do not encroach into the band around zero frequency and this is only possible providing 1/*h*>2*B*, otherwise the situation shown in figure 3.1(*h*) arises in which an overlap occurs. Cutting out the high-frequency bands, even with an infinitely sharp cut filter still does not isolate the original frequency because some degree of scrambling has taken place. The extent of this scrambling can be seen by the cross-hatched area shown in figure 3.1(*h*). It is in effect a folding back of the frequency characteristic, on itself, about the mid-line. Problems arising from this folding will be described shortly. However, one important fact emerges: in order to be sure of preserving all the information in an analogue signal of frequency *B* it is necessary to sample in time at a spacing which is a maximum of 1/2*B* long. This is called the Nyquist criterion.

Unfortunately signals in the real world do not have a distinct cut-off in frequency *B* as shown in figure 3.2. Various insignificant frequencies are invariably present therefore precluding satisfactory sampling and opening the way to the folding (or aliasing) effect. Consequently, it is usual to decide on the highest frequency of real interest, to filter the signal by analogue means to remove higher frequencies and then to sample.

Notice that the Nyquist criterion has to be relaxed slightly depending upon how sharp the analogue filter can cut. Simply sampling at 1/2*B* will cause a foldover at *B* due to the inability of filters to attenuate infinitely sharply. A guard band *G* of frequency is usually allowed to cater for the finite drop-off of the filter.

**Figure 3.1** Graphical representation of Nyquist sampling theorem showing *(a)* adequate sampling and *(b)* sampling too rarely causing aliasing.



**Figure 3.2** Use of filter to reduce bandwidth.

$G$ is taken to be about $0.5B$ so that the sample rate becomes $\sim 1/3B$ in time. It is possible to use a digital filter for this preprocessing only if the data is previously sampled at a much higher rate than the Nyquist rate for the frequency of interest. But this is sometimes wasteful in effort. It has the advantage that artificial filters can be used (described in the section on filters).

To illustrate the sort of misleading results which can occur when there is an interaction between sampling rate and signal frequency, consider figure 3.3. This shows that by sampling at a distance slightly different from that of the true wavelength a false wavelength appears. This is similar to 'beats' between waves.



**Figure 3.3**  Aliasing.

'Aliasing' is a similar effect. Apparent low frequencies are introduced by the folding of frequencies around the sampling frequency. It becomes impossible to detect whether a signal of frequency $f$ is genuine or whether it is the result of an aliased signal $2f_s - f_2$, where $f_2 - f_s = f_s = f_s - f_1$ as shown in figure 3.4.

Other forms of sampling can reduce these problems. Second-order sampling is still periodic in nature but within each period two measurements are taken, usually close together. This sort of sampling has been used on signals having a bandpass frequency characteristic. Random sampling has also been used where fewer samples need to be taken but the problem then arises of unscrambling the data afterwards.

Summarizing, samples should be taken at about $3 \times$ the rate of the highest frequency required and known to be present. Sampling much more often than this can be wasteful and only results in highly correlated data which can give biased results and lead to numerical problems.



**Figure 3.4**  Folded frequency response.

In surface texture measurement the signal representing the roughness waveform has already been smoothed relative to the true surface profile because of the finite size of the stylus tip, which acts as a mechanical filter, so that the problem of trying to decide the highest frequencies has to some extent been solved prior to measurement. In practical instruments the tip is about 2.5 $\mu$m, which implies that sampling should take place every micrometre or so. If flaws are being looked for this tip size is very necessary but for average-type parameters such as $R_a$ 5 $\mu m$ is usually adequate.

### 3.1.2   Quantization

This is not concerned with the way in which the analogue signal is turned into a time series of values. It is concerned with the actual conversion of the analogue values of the waveform into a digital form. This always means a choice between two levels, the separation being determined by the discrimination of the A/D convertor.

A point on the analogue signal at P in figure 3.5 will have to be given the value of level A or level $B$, whichever is the nearer. Having to take discrete values is the process of digitization.



**Figure 3.5** Quantization.

This breaking down of the continuous signal into discrete levels can introduce errors known as quantization errors. They do not refer to instrumental accuracy and such errors are usually small. For instance, some idea can be obtained by using Sheppard's grouped data result. With this it is easy to show that, if $q$ is the separation of levels, then an RMS noise $\varepsilon$ of $q/\sqrt{12}$ will be introduced into any assessment of the RMS value of the digital signal above that of the signal itself. Normally in metrology the quantization interval, expressed as a percentage of the signal value, is about 0.1%, and hence $\varepsilon = q/\sqrt{12} = 0.03\%$ which is negligible. It only becomes significant if the separation of levels, that is the quantization interval, becomes comparable with the signal size. In almost all practical cases in metrology the quantization intervals are equal over the whole range, but use has been made of unequal intervals in the measurement of autocorrelation, for example.

### 3.1.3   Numerical analysis—the digital model

One of the main reasons why the results obtained by different people often do not agree—even given the same data—is that not enough attention is given to numerical techniques. In this section some of the basic operations needed in analysis are examined, followed by particular problems more specific to engineering surfaces. No attempt is made to cover the subject as a whole. Suitable references enable the operations usually employed—including differentiation, integration, interpolation, extrapolation, curve fitting—to be separately studied.

### 3.1.3.1 Differentiation

There is a tendency amongst people versed in analogue ways to take a very simple formula for the first differential (figure 3.6(a)). Thus, the differential between points $z_1$ and $z_{-1}$ at $z_0$ is

$$\frac{z_1 - z_{-1}}{2h}. \tag{3.2}$$

In fact, this is only the tangent. More than just the two ordinates are needed to get a good estimate of the first differential. One usually adequate formula involves the use of seven ordinates, equally spaced by $h$.



**Figure 3.6** Numerical differentiation.

Thus the differential

$$\frac{dz}{dx}\bigg|_{z=z_0} = \frac{1}{60h}\bigg|z_3 - 9z_2 + 45z_1 - 45z_{-1} + 9z_{-2} - z_{-3}\bigg|. \tag{3.3}$$

The errors in this are of the order of $(1/140)\,\mu\delta^7 z_0$, where $\mu$ is the averaging operator between ordinates and $\delta$ is the central differencing operator. These are very small providing that ordinates outside $z_3$ and $z_{-3}$ are well behaved. A similar error in the three-point formula given above is of the order of $(1/6)\mu\delta z_0$, which turns out to be

$$\frac{1}{3h}(z_2 - 2z_1 + 2z_{-1} - z_{-2})$$

when expressed in terms of ordinates. These error formulae show that if the ordinates outside those used are significantly different then errors can creep in [1].

The numerical formulae (3.2) and (3.3) are examples of Lagrangian formulae.

By choosing a formula encompassing a wide range of ordinates the chances of rogue ordinates effecting the true derivatives are reduced. Hence the need for seven rather than three-point analysis. Similarly, to use

$$\frac{d^2 z}{dx^2}\bigg|_{z=z_0} = \frac{1}{h^2}(z_1 + z_{-1} - 2z_0) \tag{3.4}$$

as a formula for the second differential is sometimes dangerous. The errors here are of the order of

$$\frac{1}{12} \delta^4 z_0 \tag{3.5}$$

that is

$$\frac{1}{12}(z_2 - 4z_1 + 6z_0 - 4z_{-1} + z_{-2}). \tag{3.6}$$

An equivalent seven-point formula that reduces noise is

$$\left.\frac{\mathrm{d}^2 z}{\mathrm{d}x^2}\right|_{z=z_0} = \frac{1}{180h^2}(2z_3 - 27z_2 + 270z_1 - 490z_0 + 270z_{-1} - 27z_{-2} + 2z_{-3}) \tag{3.7}$$

with error $(1/560)\ \delta^8 z_0$.

Note the fundamental point about these formulae: it is still the central three ordinates that are dominant; the adjacent ordinates merely apply some degree of control over the value obtained should $z_{-1}$, $z_0$ or $z_{+1}$ be freaks or in error. Similarly, the $z_3$ and $z_{-3}$ values act as constraints on $z_2$ and $z_{-2}$ and so on.

Alternative formulae to these exist. It is possible to extend the number of ordinates on either side indefinitely, their effect getting smaller and smaller. It is also possible to evaluate the differentials in terms of backward and forward differences rather than central differences. For example,

$$\left.h\,\frac{\mathrm{d}z}{\mathrm{d}x}\right|_{z=z_0} = (\Delta + \tfrac{1}{2}\Delta - \tfrac{1}{6}\Delta^3 \ldots)z_0$$

where $\Delta$ is the forward difference, whence

$$\left.\frac{\mathrm{d}z}{\mathrm{d}x}\right|_{z=z_0} = \frac{1}{12h}(z_3 - 6z_2 + 8z_1 - 10z_0 - 3z_{-1}). \tag{3.8}$$

The derivative at $z_0$ has been evaluated by ordinates obtained later in the series, similar formulae can be obtained for second derivatives etc. The only usual reason for using these is at the beginning and ending of a set of data points. This enables all the waveform to be differentiated leaving a gap at the front or the end of the data rather than leaving gaps at both ends as is the case for central difference.

### 3.1.3.2   Integration

The tendency when attempting to integrate is to regard the sum of ordinates multiplied by their interval as equal to the value of the integral, that is the area under a curve. Although this is true enough when the number of terms is large it is not necessarily so when the number of ordinates is small. Figure 3.7(a) shows a typical metrological example. Often it is required to evaluate convolution integrals, especially in the context of filtering. In these cases a weighting function representation of the filter is used. This weighting function is convoluted with the profile signal to get a filtered output. In the case of a steady-state (DC) input signal, all that is required is the weighting function to be integrated, that is the area under the function between limits

**Figure 3.7** Numerical integration: (*a*) simple addition of ordinates; (*b*) modified addition.

along the axis evaluated (i.e. a definite integral). Summing ordinates or weighting factors derived from the function directly is equivalent to measuring the area shown as squares.

Clearly, the first ordinate is contributing about twice the area needed. The intended area is on the left-hand side of the *z* axis, so using just the height factors of the function gives too large a value for the area,

A better solution is shown in figure 3.7(*b*). Here the ordinates of the weighting function can be grouped. Thus, if the ordinates are $b_0, b_1, b_2, \ldots, b_n$ starting at $t = 0$, then a grouping of

$$\frac{b_0 + b_1}{2}, \ \frac{b_1 + b_2}{2} \tag{3.9}$$

will produce an area representation of figure 3.7(*b*)

In a convolution with profile ordinates $a_0, a_1, a_2$, etc, the integral would be

$$k\left( a_0 \frac{(b_0 + b_1)}{2} + a_1 \frac{(b_1 + b_2)}{2} \ \ldots \right) \tag{3.9a}$$

where *k* is a constant.

Now, however, each profile ordinate becomes associated with the succesive mid-ordinate positions of the weighting function (e.g. in between $b_0$ and $b_1$) with the result that there is a phase shift of half an ordinate spacing between the profile and the operative part of the filter which cannot always be ignored.

The numerical formulae for any function being integrated therefore are not simply the sum of the evaluated points; the values can only be added after modification. Thus, in the case above, the factors would be $b_0/2$, $b_1, b_2, \ldots, b_n/2$. This is no more than a statement of the trapezoidal rule for numerical integration and it corresponds, in fact, to joining the function values by straight lines instead of each ordinate being represented by a block. The area is made up of trapezoids rather than rectangles. It is a special case of Gregory's formula.

Further improvement can be achieved by fitting curves between the points of evaluation. The quadratic curve gives rise to Simpson's rule for numerical integration. Interpreted in a purely geometric way this gives the sum of the areas under second-degree parabolas that have passed through the points $b_0, b_1, b_2 \ldots$

In this case if there are $n + 1$ pairs of given values, where *n* is even, then

$$\int_{x_0}^{x_n} b \ \mathrm{d}x = \frac{h}{3}[b_0 + b_n + 4(b_1 + b_3 + \ldots + b_{n+1}) + 2(b_2 + b_4 + \ldots + b_{n-2})] \tag{3.10}$$

for $n + 1$ pairs of values, and if $n$ is a multiple of three then

$$\int_{x_0}^{x_3} b \; dx = \frac{3h}{8}[b_0 + 3(b_1 + b_2) + b_3] \tag{3.11}$$

which can be extended to take into account the general case by grouping the ordinates four at a time, giving

$$\int_{x_0}^{x_n} b \; dx = \frac{3h}{8}[b_0 + b_n + 3(b_1 + b_2 + b_4 + b_5 + \ldots + b_{n-2} + b_{n-1}) + 2(b_3 + b_6 + \ldots + b_{n-3})]. \tag{3.12}$$

Other more elaborate formulae exist involving unequal spacing of the function values. These require fewer values for a given accuracy for the reasons given below. However, it is well known that for most engineering applications the equal-interval formulae show a surprising accuracy, especially Simpson's rule. If, however, more accuracy is required then other techniques may have to be adopted. Lagrange interpolation formulae involve fitting a polynomial through any set of points not necessarily equally spaced. This polynomial therefore represents the function as a whole. Even when the intervals are equal, Lagrangian techniques have the advantages of permitting interpolation without the need to construct a difference table. They have the disadvantage common to all polynomial curve-fitting methods of requiring some knowledge of the degree of the polynomial needed to achieve any real accuracy.

In integration, use can also be made of the unequal interval in the Gauss' formulae for numerical integration. Although these and other formulae exist they are most useful in those situations where a limited amount of well-understood data is present. Usually there is enough data to get the required accuracy for the basic operations using equal intervals and first or second-degree interpolation.

### 3.1.3.3 Interpolation, extrapolation

In many respects these two areas of digital analysis are the most fundamental because in the process of sampling continuous signals have been transformed into discrete values in time (space). Reconstitution of the original signal in between sampled values is of prime importance, as is the ability to project the signal outside its immediate vicinity. Fortunately in surface metrology a lot of data is usually available, and consequently the use of such techniques is rarely required. In essence a polynomial is fitted to the measured or listed points and made to conform with the fixed points. Examination of the polynomial behaviour between and also outside the fixed values is then possible. In particular, the most used polynomial is called the Lagrange interpolating polynomial.

Well-known formulae exist for interpolation and extrapolation, in particular those due to Everett and Bessel [1] for equally spaced data and Lagrange for unequally spaced data. An example of the use of interpolation formulae in surface topography comes in the field of contact, where mechanical models of surface peaks are used. Using interpolation it is possible to find the position of the maximum value of a peak even if it is not touched by the sampled points. For instance, if the ordinates are $z_{-1}$ and $z_{+1}$ where $z_1 \neq z_{-1}$ then the apex is not at $z_0$ but at a distance $V$ from $z_0$ given by

$$V = \frac{(h/2)(z_{-1} - z_1)}{2z_0 - z_1 - z_{-1}}. \tag{3.13}$$

It is further possible to work out curvatures at the apex of this peak. Interpolation can also help to cut down the number of mean line points that need to be calculated for the mean line assessment.

## 3.2 Digital properties of random surfaces

### 3.2.1 Some numerical problems encountered in surface metrology

Some examples will be given by way of illustration of some of the typical problems encountered in surface metrology. Those chosen will not be exhaustive in content but should give a fairly representative cover of the types of problem likely to be encountered.

There are a limited number of features of most interest. From the section on surface characterization much emphasis was placed on the peak and slope measurement of a profile and, in particular, in two dimensions (or areal). These features can be extended to include peak height distributions, curvature distributions and how they vary with height, slope and associated parameters.

It is only in the last twenty years that digital methods have become available to make possible the measurement of these important parameters. However, simply theorizing about parameters is not enough—they have to be measured. This apparently straightforward task is fraught with problems as section 3.2.2 will show. In the past parameters have been restricted to those that could be measured with simple analogue circuitry. This in itself imposed natural constraints upon the wavelengths which could be measured. Recorders, amplifiers and meters have bandwidth limitations which cannot be overcome but the restriction on digital methods is less distinct; the experimenter is met head-on with the sampling, quantization and numerical model problem. Sometimes these considerations are outside the experience or training of the investigator. Unfortunately, the parameters of use to the surface metrologist are just those parameters that are difficult to measure. In what follows it will become apparent that the correlation function of the surface is of fundamental importance in assessing the change in value of such parameters with sampling. Instrumental limitations such as the stylus tip or optical resolution will be incorporated more fully in chapter 4.

### 3.2.2 Definitions of a peak and density of peaks [2]

The problem of defining peaks and assessing peak properties often arises in metrology, especially in contact theory. For instance, it is often of interest to know the density of peaks of a surface profile. The question that has to be posed is how this count of peaks depends on the digitizing interval, the quantization interval and the definition of the peak. All three can affect the count. This is one of the reasons why it is so very difficult to get agreement between researchers, even given the same data.

Take, for example, the problem of the definition. One of the most used definitions is a three-point model as shown in figure 3.8(a). If the central ordinate of three consecutive ordinates is the highest then the three together constitute a peak. An alternative one is also shown in the figure in which four ordinates are used, the central two being higher than the others for a definition. Many similar possibilities exist. In any case the number of peaks counted will be different for the same data. For example, a peak counted by the three-point method could get ignored using the four or more ordinate models.

Also, differences in definition have been used within the three-point method. Some investigators have imposed a height difference constraint on the definition. For instance, the central ordinate has to be a height $z'$ above the higher of the other two before a peak is registered, that is as in figure 3.8(a), $z_0 - z_1 > z'$. This constraint reduces the count as before.

### 3.2.3 Effect of quantization

The quantization interval can influence the count as is shown in figure 3.8(b). It can be seen that using exactly the same waveform, simply increasing the quantization interval by a factor of 2 means that, in the case of figure 3.8(b), the three-point peak criterion fails, whereas in the other case it does not.

So, even the A/D resolution can influence the count. In order to get some ideas of the acceptable quantization interval it should be a given ratio of the full-scale signal size, subject to the proviso that the interval

**Figure 3.8** (*a*) Possible definition of peaks, (*b*) effect of quantization on peaks.

chosen gives sufficient accuracy. As an example of the quantitative effect of quantization, consider a signal that has a uniform probability density. If the range of this density is split up into $m+1$ levels (i.e. $m$ blocks, figure 3.9) then it can be shown that the ratio of peaks to ordinates is given by

$$ratio = \frac{1}{3} - \frac{1}{2m} + \frac{1}{m^2} \, . \tag{3.14}$$

This makes the assumption that the samples are independent and that the three-ordinate model is used as a definition of a peak.



**Figure 3.9** Uniform distribution of intervals.

Examination of the formula shows that, when $m$ is large, the ratio is 1/3. This makes sense because, for independent samples with no quantization restriction, one would expect one-third of all samples to be peaks. Similarly, when m = 1, the ratio is zero. Again this makes sense because no information is being conveyed.

Various other values of $m$ are listed in table 3.1.

**Table 3.1**

| $m$ | Ratio peaks/ordinates | Percentage drop from continuous signal |
|-----|----------------------|----------------------------------------|
| 1   | 0.125                | 70                                     |
| 3   | 0.19                 | 45                                     |
| 4   | 0.21                 | 37                                     |
| 5   | 0.23                 | 30                                     |
| 10  | 0.25                 | 15                                     |
| 100 | 0.33                 | <1                                     |

Therefore it is obvious that even representing the signal by one decimal digit goes a long way towards giving acceptable results. An extra 15% can be obtained by going to the second decimal digit. On this basis, taking measurements to the third place seems unnecessary. Taking them to the fourth place certainly is.

A similar formula to the one for rectangular (or uniform) distributions can be obtained for the very important Gaussian distribution. Again using the three-point analysis the probability of an ordinate being a peak is given at a level between quantization levels of $n\Delta z$ and $(n-1)\Delta z$ by

$$\left| \int_{(n-1)\Delta z}^{n\Delta z} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right) dz \right| \left| \int_{-\infty}^{(n-1)\Delta z} \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-z^2}{2}\right) dz \right|^2 \tag{3.15}$$

assuming that the z values have a zero mean and unit variance.

Taking this over all possible values of interval allowed, in this case taken to be symmetrical about $z = 0$ and extending $m/2x$ blocks to either extreme of $3\sigma$ where $\sigma$ is the RMS value of the distribution, the probability is

$$\frac{1}{8} \sum_{-3\sigma \cong \Delta z}^{n=3\sigma \cong \Delta z} \left[ \mathrm{erf}\left(\frac{n\Delta z}{\sqrt{2}}\right) \mathrm{erf}\left(\frac{(n-1)\Delta z}{\sqrt{2}}\right) \right] \left| 1 + \mathrm{erf}\left(\frac{(n-1)\Delta z}{\sqrt{2}}\right) \right|^2 \tag{3.16}$$

which gives the following results:

$$m = 2 \quad \text{probability} = 0.125$$
$$4 \qquad\qquad\qquad 0.167$$
$$6 \qquad\qquad\qquad 0.211$$
$$12 \qquad\qquad\qquad 0.266.$$

These generally follow the rectangular case except for being slightly lower. In both cases the value is asymptotic to 1/3.

### 3.2.4 Effect of numerical model

While investigating the effect of quantization it is informative to investigate the effect of different models. This is especially valid in peak behaviour because it is in the region of a peak that changes in level tend to be

small and the effect of quantization can dominate the result. To combat this predicament it is sometimes useful to use a modified three-ordinate model in which, if the next-to-central ordinate is at the same digital level as the central ordinate, the judgement of whether or not a peak exists is deferred for one ordinate. If this is still not conclusive it is deferred further until an adjacent ordinate is lower.

Under these circumstances one would expect rather more peaks to be revealed than when restricting the model strictly to three ordinates. The effect of this modification can be taken account of in probability merely by noting that the probability of the situation shown in figure 3.10 is the same as for the three-point one except that the central ordinate has been repeated three times, that is the probability is $P_1^3 P_2^2$, where $P_1$ is the probability of an ordinate lying between $(n-1)\Delta z$ and $n\Delta z$, and $P_2$ is the probability of an ordinate being below $(n-1)\Delta z$.



**Figure 3.10** *n*-ordinate definition of peak.

Taking all such possibilities into account gives an additive sequence of probabilities. Thus

$$P_1 P_2^2 + P_1^2 + P_2^2 + P_1^3 P_2^2 + \ldots + P_1^n P_2^2 = \frac{P_1 P_2^2}{1 - P_1}. \tag{3.17}$$

The formulae derived previously for the three-point model definition of a peak can therefore be modified to give the following:

(1) for a rectangular distribution

$$\text{probability} = \frac{2m^2 - 2m - 1}{6m(m-1)} \tag{3.18}$$

(2) for the Gaussian distribution

$$\frac{1}{8} \sum_{n=-m/2}^{n=+m/2} \frac{[\text{erf}(n\Delta z/\sqrt{2}) - \text{erf}((n-1)/\sqrt{2})\Delta z][1 + \text{erf}((n-1)/\sqrt{2})\Delta z]^2}{1 - \frac{1}{2}[\text{erf}(n\Delta z/\sqrt{2}) - \text{erf}((n-1)/\sqrt{2})]} \tag{3.19}$$

Both of these formulae have the effect of considerably increasing the probability of a peak, especially at small values of *m* (see figure 3.11). For both in fact, if the quantization is purely linear, that is 0+1 or ±1, a complete change in the count can result depending simply upon which numerical model is chosen to represent a peak—on the same original profile.

As will now be obvious, each one of the constraints of sampling, quantization and numerical definition can aggravate the effects of the others.

**Figure 3.11** Peak density—effect of quantization, model and sampling: numerical problems in system parameter evaluation—an example of peak density count.

Again, notice that as $m \to \infty$ for both formulae the value of the probability tends to 1/3 because as $m \to \infty$ the chances of more than one ordinate lying within the infinitesimal limit become increasingly remote.

The main reason for taking a transducer range of more than one decimal digit is that it is not always possible to ensure that the signal is a significant percentage of the range of the A/D convertor. For very smooth surfaces this is often the case when the part is fitted. Another reason that will be amplified later is concerned with the best use of computers in metrology. The latest trend is to let the computer take out a large part of the manual setting-up by accepting all the signal from the workpiece—including errors of form, misalignment, waviness, etc—removing all the errors that are unwanted, digitally, and then magnifying the remainder for evaluation. This technique presupposes that the remaining signal is known with sufficient accuracy (i.e. to a sufficient resolution) to be useful. In order to ensure this, the whole of the original signal has to be digitized to a reasonably small quantization interval. This is because the process of removing the large error signal leaves only a small percentage of the total, as seen in figure 3.12.

### 3.2.5 Effect of distance between samples on the peak density value

The real issue here is not the physical distance between the discretely-measured sampled ordinates, but how this distance relates to the autocorrelation function of the surface. In order to investigate the way in which ordinate spacing affects surface parameters it is necessary to know what constraints exist on the choice of a typical correlation function.



**Figure 3.12** (*a*) Total signal, (*b*) magnified remnant.

Perhaps the obvious choice of autocorrelation function is the exponential form shown in figure 3.13. The argument for adopting this is that most finished surfaces exhibit such autocorrelation functions in regions other than at the origin. This is a direct result of the Poissonian nature of surface generation incurred because of the random cutting action.



**Figure 3.13** Exponential autocorrelation function.

There are, however, theoretical objections to the use of the exponential autocorrelation function, for example there are certain features of the exponential function at the origin which are undesirable. In particular, the density of peaks $D_p$ is given, for an ordinary, random process, by

$$D_p = \frac{1}{2\pi} \left( \frac{A^{iv}(0)}{-A''(0)} \right)^{1/2}$$

(3.20)

where $A''(0)$ is the second differential of the autocorrelation function at the origin and $A^{iv}(0)$ is the fourth. For the exponential function these are undefined because $A'(0)$ and $A'''(0)$ are non-zero. This becomes necessary because the autocorrelation function is an even function. In practice this is not a serious problem because there is always some degree of smoothing of the profile caused by the finite resolution of the measuring instrument. In what follows the exponential correlation will be assumed for simplicity unless specified otherwise.

Consider for example, as before, the measurement of peak density or the probability that an ordinate is a peak. Using the three-point model this becomes $N$, where $N$ is given by

$$N = \frac{1}{\pi} \tan^{-1} \left( \frac{3 - \rho}{1 + \rho} \right)^{1/2}$$

(3.21)

and $\rho$ is the correlation between ordinates separated by $h$. For an exponential correlation function $h = \beta \ln(1/\rho)$, as shown in figure 3.13.

The peak density is given simply by

$$N/h.$$

(3.22)

In the case of a general autocorrelation function shown, for example in figure 3.14, two correlation coefficients are needed, $\rho_1$ corresponding with the autocorrelation value at $h$ and $\rho_2$ the value at $2h$. The formula relating the peak density to the correlation function and hence to the spacing is modified somewhat to give

$$\frac{N}{h} = \frac{1}{\pi h} \tan^{-1}\left[\left(\frac{3 - 4\rho_1 + \rho_2}{1 - \rho_2}\right)^{1/2}\right]. \tag{3.23}$$

How the value of peak density changes with $h$ can easily be calculated for an exponential autocorrelation function and is shown in figure 3.15.



**Figure 3.14** General autocorrelation function.



**Figure 3.15** Peak density—effect of correlation.

Varying $\rho$ from 0 to 1 gives values of $N$ from 1/3 to 1/4 respectively, in theory at least. The value of 1/3 is to be expected because when $\rho = 0$ the sampled ordinates are independent. The separation when this occurs is usually a matter of definition. One such definition, shown in figure 3.13, is the value of $h$ such that $\rho \sim 0.1$ to measure is $\tau_{max}$. Another definition, which is more reliable but more difficult to measure, is

$$\tau_{max} = \int_0^\infty |A(\tau)|\, d\tau \tag{3.24}$$

which corresponds to a value of $\rho = 1/e$.

When the ordinates are highly correlated the density is $1 \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$ so that, in theory at least, as $\rho \to 1$, $N \sim \frac{1}{4}$.

When $h$ is small the practical value of peak density always falls to zero. This can be seen by postulating any correlation function having well-behaved derivatives at the origin and inserting the $\rho_1$ and $\rho_2$ values measured from it into the general formula for $N$ given above in equation (3.23).

For the case of the exponential function all digital features can be incorporated. Thus

$$
N = \frac{1}{8h}\left[ \mathrm{erf}\left( \frac{n\Delta q}{\sqrt{2}} \sqrt{\frac{1-\rho}{1+\rho}} \right) - \mathrm{erf}\left( \frac{(n-1)\Delta y}{\sqrt{2}} \sqrt{\frac{1-\rho}{1+\rho}} \right) \right]\left[ 1 + \mathrm{erf}\left( \frac{(n-1)\Delta y}{\sqrt{2}} \sqrt{\frac{1-\rho}{1+\rho}} \right) \right]^2 \Bigg/
$$
$$
\left\{ 1 - \frac{1}{2}\left[ \mathrm{erf}\left( \frac{n\Delta y}{\sqrt{2}} \sqrt{\frac{1-\rho}{1+\rho}} \right) - \mathrm{erf}\left( \frac{(n-1)\Delta y}{\sqrt{2}} \sqrt{\frac{1-\rho}{1+\rho}} \right) \right] \right\}. \tag{3.25}
$$

This incorporates quantization, sampling and the numerical model and is shown in figure 3.11 for $\rho$ values of 0, 0.7 and 0.9, the last two being typical of those used in practice.

From what has been said it is clear that problems of surface analysis, such as peak definition, arise from three different factors: the quantization, the sampling and the model. These three cannot be divorced from each other. Not to realize this interrelationship can be the source of much trouble when different investigators compare results. Any one of these three variables can be used as a fixed datum: the model, because it is used universally; the quantization, perhaps because of an accuracy requirement on individual measurements or subsequent analysis; and/or the sampling, because of instrumentation problems (stylus resolution) or the functional significance of the harmonic picked out.

To illustrate this interdependence, suppose that it is important to make sure that a particular numerical model is universally valid, say the simple three-point method for peaks. A criterion could be used in which the allowable loss of peak data is, for instance, 5%. If the surface is Gaussian some quantitative picture can be obtained.

From figure 3.11 it can be seen that in order to ensure that only 5% of peak information is lost using the three-ordinate model, about 50 or 60 quantization intervals are required throughout the range of the signal. Because the ordinates in this model are Gaussian it is possible to work out a relationship between the average difference of height from ordinate to ordinate and to equate this difference of adjacent ordinate heights as $E(z_1 - z_2)^2 = 2\,\sigma^2(1-(2))$, where $\sigma$ is the RMS value of the surface and $\rho$ is the value of the correlation between sets of data separated by such an interval in space that the autocorrelation is $\rho$.

The average separation $\Delta z$ is $\sqrt{2/\pi}$, the square root of the variance. Thus

$$
\Delta z = \frac{2\sigma}{\sqrt{\pi}}\sqrt{1-\rho^2} = 1.13\sigma\sqrt{1-\rho^2} \tag{3.26}
$$

For $\rho = 0$, that is independence between ordinates, $\Delta z = 1.13\,\sigma$.

If it is assumed that the range of the signal is $\pm 3\sigma$ then the acceptable relation between $z$ and $q$, the quantization interval, is

$$
\frac{1.13 \times (50 \sim 60)}{6} \sim 12
$$

that is $\Delta z \simeq 12q$. (It is not just in the measurement of peaks that problems lie. There are similar problems in slope measurement and curvature, as will be seen later in sections 3.2.6.2 and 3.2.6.3.)

Remembering that this ratio must be preserved for all other digital sample intervals in order to maintain the 95% integrity, the ratio of $q$ to $\sigma$ for any other interval can be determined. Take for example the case where $\rho = 0.7$. The value of $\Delta z$ becomes $0.58\sigma$. Hence $q = \Delta z/12 \sim 0.05\ \sigma$. Therefore, to cover the range of $6\sigma$ for the signal requires about 120 quantization intervals $q$.

Thus, in order to maintain the integrity of the model from a sampling giving $\rho = 0$ to one giving $\rho = 0.7$, the number of $q$ levels has to be increased from a nominal 60 to 120, a factor of 2. If this integrity is not preserved comparative deviations in peak count cannot be attributed to the surface profile itself. A similar quantitative result is obtained if the differences between ordinates around the peaks themselves are investigated rather than those of the profile as a whole.

Other constraints may have to be added, for instance the inability of the stylus to resolve small-wavelength undulations or the need to exclude the short undulations for functional reasons. Two things are plain: first, the three variables of numerical model, sampling and quantization are related and, second, at least one of these three needs to be fixed, preferably by using a functional requirement—which is usually a sampling consideration. So the preferred rule is to choose the sampling to match the function and work out a suitable quantization and numerical model in terms of the sampling. In a lot of work on contact wear and friction, it is the features of long wavelength picked out by independent samples that dominate performance and therefore need to be unambiguously measured.

### 3.2.6 Digital evaluation of other important peak parameters

#### 3.2.6.1 Peak height measurement

The method adopted previously in which the exponential autocorrelation function is used as a basis will be continued here. Expressions are derived which can easily be transformed into the general autocorrelation case.

The mean peak height $\bar{z}_p$ as a function of correlation between ordinates of $\rho$ is given by

$$\bar{z}_p = \frac{\sqrt{(1-\rho)/\pi}}{2\cos^{-1}[(1+\rho)/2]^{1/2}} \tag{3.27}$$

and for the general autocorrelation by

$$\bar{z}_p = \frac{[(1-\rho_1)/\pi]^{1/2}}{2N} \tag{3.28}$$

where $N$ is the probability that an ordinate is a peak; $\bar{z}_p$ is normalized with respect to $\sigma$.

Therefore, in the simple exponential case the mean peak height expressed in terms of the RMS value of the surface migrates from $1.5/\sqrt{\pi}$ to zero as the correlation changes from zero to unity. This height value would have to be multiplied by $\sqrt{\pi/2}$ if the deviations were to be expressed in terms of the $R_a$ value of the surface. This behaviour is modified slightly in the case of the general autocorrelation because although, as one would expect, the independent result is the same, the limiting behaviour as $\rho_1$ and $\rho_2$ tend to unity is different. There is a limiting value of mean peak height for $\rho_1 = \rho_2 = 1$ which depends on the particular type of autocorrelation function modelled. Thus, for a highly oscillatory autocorrelation function, the limiting value of $\bar{z}_p$ could be near to or even higher than the RMS value $R_q$ of the surface. In any event it must be positive.

The standard deviation of peak height $\sigma_p(\rho)$ is also affected in the same way. Thus for an exponential autocorrelation function

$$\sigma_p(\rho) = \left(1 + \frac{(1-\rho)\sqrt{1+\rho}}{2\sqrt{3-\rho}\,\tan^{-1}\sqrt{(3-\rho)/(1+\rho)}} - \frac{\pi(1-\rho)}{4[\tan^{-1}\sqrt{(3-\rho)/(1+\rho)}]^2}\right)^{1/2} \tag{3.29}$$

which, for $\rho = 0$, gives $\sigma_p = 0.74$ and for $\rho = 1$ gives $\sigma_p = 1$, the same value as the profile itself. The standard deviation of mean peak height for a general correlation is

$$\sigma_p(\rho_1, \rho_2) = \left[ 1 + \frac{(1 - \rho_1)}{2\pi N}\left( \frac{1 - \rho_2}{3 - 4\rho_1 + \rho_2} \right)^{1/2} - \frac{(1 - \rho_1)}{4\pi N^2} \right]^{1/2} \tag{3.30}$$

where $N$ is defined for the general case in equation 3.23.

Some idea of the effect of quantization can be obtained by using Sheppard's correction for grouped data. Thus if $\sigma_{pq}(\rho)$ is the quantized value

$$\sigma_{pq}^2(\rho) = \sigma_p^2(\rho) + q^2/12. \tag{3.31}$$

With a typical worst-case value of $q = 0.5$, $R_q$ reveals that an error of a few per cent in $\sigma_p$ will result. Obviously peak height measurements can also be in error due to numerical modelling. For example, the three-point model is only an approximation, and is in effect a vertical parabola fitted through three ordinates. What is normally assumed is that the central ordinate is at the apex of the parabola. In fact this is only true if those ordinates adjacent to the central one are of equal height; if not, the error in peak height $\delta z_p$ is given by

$$\delta z_p = \frac{1}{8} \frac{(z_{+1} - z_{-1})^2}{(2z_0 - z_{+1} - z_{-1})} \tag{3.32}$$

Taking typical values of $z_0$, $z_{-1}$ and $z_{+1}$ to be 3, $-1$ and $+1$ respectively, an error of about 3% results. Errors are generally much smaller for highly correlated data.

### 3.2.6.2 Peak curvature

The curvature of peaks is especially sensitive to extraneous effects. This is mainly due to the fact that curvature is defined in terms of derivatives of the signal which are notoriously prone to noise. To get an estimate of the size of digital effects, consider the basic formula for curvature:

$$C = \frac{d^2z/dx^2}{(1 + (dz/dx)^2)^{3/2}}. \tag{3.33}$$

Usually, because in the vicinity of a peak $dz/dx \sim 0$ the curvature can be closely approximated by the second differential, so the curvature $C$ can be most simply expressed by the three-ordinate Lagrangian formula

$$C = \frac{2z_0 - z_{+1} - z_{-1}}{h^2} = \frac{C'}{h^2} \tag{3.34}$$

where $C'$ is curvature measured in terms of ordinate differences only.

As before, the mean peak curvature can be expressed digitally as $\overline{C}'(\rho)$ where

$$\overline{C}'(\rho) = \frac{(3 - \rho)\sqrt{(1 - \rho)/\pi}}{2\cos^{-1}[(1 + \rho)/2]^{1/2}} \tag{3.35}$$

To get $\bar{C}(\rho)$ this is simply divided by $h^2$.

For the general case $\bar{C}(\rho)$ is given by

$$\bar{C}(\rho) = \frac{3 - 4\rho_1 + \rho_2}{2Nh^2[\pi(1 - \rho_1)]^{1/2}} R_q^2. \tag{3.36}$$

Comparing, for example, $\bar{C}'(\rho)$ at $\rho = 0.1$ and $\rho = 0.9$ shows a ten to one difference, but in practice the difference in $\bar{C}(\rho)$ will be larger because, not only are the differences between ordinates changing as shown in $\bar{C}'(\rho)$, but the actual value of $h$ changes and it is this which has the dominant effect. Changes in curvature of fifty to one can be obtained on ordinary surfaces by changing the sampling rates!

A simple maximum probable estimate for the errors in $C$ due to quantization is $\delta C$:

$$\delta C \simeq 4q/h^2 \tag{3.37}$$

where $h$ is, as before, the ordinate spacing and $q$ is the quantization interval. For typical values of $q$ and $h$ it would appear at first sight that the error is not likely to be serious. This is not necessarily so because near to a peak the differences between ordinates is small and can even be of the same order as $q$, with the result that very large errors can occur. Increasing the ordinate spacing $h$ in this case not only increases the ordinate differences but reduces $\delta C$ in the above expression.

It is possible to get a limit for $q$ in terms of ordinate differences above which useless estimates of curvature will be obtained. The mean value of ordinate differences is given by

$$\frac{(3 - \rho)(1 - \rho)^{1/2}}{4\pi}. \tag{3.38}$$

Hence the upper limit for $q$ could be taken as the inequality

$$q << \frac{(3 - \rho)\sqrt{1 - \rho}}{4\pi} R_q \sim \frac{\sqrt{1 - \rho}}{2\sqrt{\pi}} R_q. \tag{3.39}$$

Knowing $q$, $\rho$ can be found (assuming the exponential model of a given exponent), from which the minimum ordinate spacing can be worked out. For example, for $\rho = 0.7$ it works out that $q < 0.17 R_q$, which implies that there should be about 35 quantization intervals across the range of the signal.

Numerical model errors are also important in measuring curvature. The three-ordinate model used earlier is merely the simplest second-differential formula. It is only strictly valid to use it when the fourth-order central differences are small when compared with it. A factor of one-tenth is typical, that is

$$2z_0 - (z_{+1} + z_{-1}) \simeq 2z_0 - (z_{+1} + z_{-1}) + \frac{\delta^4}{12} z_0. \tag{3.40}$$

This in turn is most likely to be true when the correlation between ordinates is high and the profile smoothly fits between discrete sampled ordinates.

For higher accuracy better Lagrangian formulae exist; for example, the seven-point model gives

$$C = \frac{1}{180h^2}(2y_3 - 27y_2 + 270y_1 - 490y_0 + 270y_{-1} + 27y_{-2} + 2y_{-1}) \tag{3.41}$$

in which case the eighth central differences have to be small.

One way to assess the suitability of a curvature formula is to find out where it fails to be a good second differential. The way to do this is to consider the formula to be a digital filter by finding the Fourier transform of the numerical sequence and thence to find the break point where the gain does not increase as $\omega^2$. This is because a theoretically correct second differential should have a transfer function of $\omega^2$ for all $\omega$.

Hence if $C(x)$ *is* the operator on $f(x)$ to get $f''(x)$, that is

$$f''(x) = f(x) * C(x), \tag{3.42}$$

where * denotes convolution

$$C(x) = 2\delta z(x - 3h) - 27\delta z(x - 2h) + 270... \text{ etc.}$$

Thus, if $A(\omega)$ is the transform of $C(x)$

$$A(\omega) = \frac{1}{180h^2}[2\ \exp(3j\omega h) - 27\ \exp(2j\omega h) + 270\ \exp(j\omega h) - 490$$

$$+270\ \exp(-j\omega h) - 27\ \exp(-2j\omega h) + 2\ \exp(-3j\omega h)] \tag{3.43}$$

which results in

$$A(\omega) = \frac{1}{180h^2}[4\cos3\omega h - 54\cos2\omega h + 450\cos\omega h - 490]. \tag{3.44}$$

If the frequency characteristic is plotted out it will be observed that it is only when $\omega \sim 2/h$ that it behaves as a true second differential. For higher values of $\omega$ the differentiating property breaks down. For the three-point differential model, $\omega$ has to be even shorter, so even though the three-point analysis is suitable for detecting peaks without the complication of quantization, it is severely restrictive when measuring curvature and should only be used when the shortest wavelengths present are four or five times the spacing interval $h$.

### 3.2.6.3 Profile slopes

Exactly the same arguments can be used when attempting to measure slopes: the first-differential formulae referred to in an earlier section still hold. In this case, however, the breakdown of the formulae is that point where the frequency characteristic fails to increase as $\omega$ and not $\omega^2$. A similar conclusion to that for $C$ is reached [3].

In terms of the autocorrelation function and discrete ordinates, the mean absolute slope is given by

$$\overline{m} = \frac{R_q}{h}\sqrt{\frac{1-\rho}{\pi}} \qquad \text{or} \qquad \Delta q = \frac{R_q}{h\sqrt{2}}\sqrt{1-\rho} \qquad \text{from } E\frac{[z_{-1} - z_{+1}]^2}{(2h)^2}$$

or

$$\overline{m} = \frac{R_q}{h}\sqrt{\frac{1-\rho_2}{\pi}} \tag{3.45}$$

for a general correlation function. As before for curvature of peaks

$$f'(x) = m(x) * f(x). \tag{3.46}$$

Hence

$$m(x) = \frac{1}{12h}[\delta z(x - 2h) - 8\delta z(x - h) + 8\delta z(x + h) - \delta z(x + 2h)] \tag{3.47}$$

for the five-point model from which

$$m_5(\omega) = \frac{j}{6h}[8\sin\omega h - \sin 2\omega h].$$

(3.48)

Expanding this as a power series in $\sin \omega h$:

$$m_5(\omega) = \frac{j}{h}\left(\omega h - \frac{(\omega h)^5}{40} - \frac{(\omega h)^7}{252} + \cdots\right)$$

(3.49)

from which it is obvious that only the first term represents a true differentiation being proportional to $\omega$; the others are errors due to the limitation of the five-point method.

The actual error at a specific frequency can be found by using equation (3.48). Thus, as an example, consider a sine wave sampled at four points per cycle. It will have a wavelength of $4h$ and $\omega = 2\pi/4h$, the true value is $m(\omega)$

$$m(\omega) = j\omega = j\frac{\pi}{2h}$$

(3.50a)

whereas

$$m_5(\omega) = \frac{j}{6h}\left(8\sin\frac{\pi}{2} - \sin\pi\right) = j\frac{4}{3h}$$

(3.50b)

which is 15% down on the true value.

Note that the formula for the average values of slope curvature and other features in terms of the spacing and the correlation coefficients of the single, double, and in fact multiple spacings between ordinates can be extended.

There is a general pattern that can be used to get some idea of the applicability of surface parameters in terms of digital characteristics. This is shown in simple form in figure 3.16.



**Figure 3.16** Fidelity of surface parameters.

In the figure, region *A* gives the wrong results because the ordinates are too highly correlated and the small differences between them (perhaps due to the finite tip of the stylus or the limited resolution of the optical device) do not show at all, because the quantization interval $\Delta z$ does not see the differences, so region *A*

reflects instrument resolution + quantization. Region *B* tends to produce realistic values where the correlation is about 0.7. In this region the instrument limitations are minimal and the reasonably fine structure of the surface as well as the gross structure are picked up. In region C the fine structure is lost because $\rho \sim 0$, aliasing can be a problem and only the gross structure is picked up. Loss of information and misleading information are very possible. For region *D* the answer is completely dominated by the numerical model—the data is completely uncorrelated. The probability of an ordinate being a peak is 1/3 if the three-ordinate model is used and 1/5 if the five-ordinate model is used, so region *D* is numerical model limited.

The graph shown in figure 3.16 demonstrates the limited range in which some sort of result can be obtained which is not corrupted by either the instrumentation or the digital processing. The question arises, therefore, whether the data means anything anyway, because the answer appears to be dependent on the sampling even if no quantization effects are present! There are two issues: one is whether a digital technique would ever agree with an analogue one; and the other is whether surfaces are well behaved enough in the sense of differentiability to allow a unique answer for a feature which did not continually change as the sample interval and hence 'scale of size' changed. The first point will be addressed first, not because it is simpler but because it can be attempted. Even then there are two aspects to it: the profile problem and the 'areal' problem. Consider the profile first.

Take for example the case of the mean number of peaks $m_0$ in unit distance and the mean number of crossings $n_0$ in unit distance in terms of $D_2$ and $D_4$ where

$$D_r = \left. \frac{d^r A(h)}{dh^r} \right|_{h=0}.$$

(3.51)

These are as revealed in chapter 2 and section 3.2

$$m_0 = \frac{1}{2\pi} \left( \frac{D_4}{-D_2} \right)^{1/2} \qquad\qquad n_0 = \frac{1}{\pi} \left( \frac{D_2}{-D_0} \right)^{1/2}.$$

(3.52)

The results here can be compared with those given in digital form in equation (3.23) for $m_0$ and $n_0$

$$\frac{1}{\pi} \tan^{-1} \left( \frac{3 - 4\rho_1 + \rho_2}{1 - \rho_2} \right)^{1/2} \qquad \text{and} \qquad \frac{1}{\pi} \cos^{-1} \rho_1$$

(3.53)

by first expressing the autocorrelation function as a Taylor expansion and then investigating the behaviour as the sampling interval $h$ is made to approach zero when $\rho_1 = \rho(h)$ and $\rho_2 = \rho(2h)$. Thus

$$\rho(h) = 1 - (-D_2) \frac{h^2}{2!} + D_4 \frac{h^4}{4!} + O(h^4).$$

(3.54)

Using this expansion and inserting the values for $\rho_1$ and $\rho_2$ in the equations for peak distributions and other features described above, it can be shown [4] that they become the results obtained by Rice in 1944 [12] and by Bendat [6] thereby satisfactorily allowing the conclusion that the digital method does indeed converge onto that which would have been obtained from a continuous surface with perfect instrumentation. Unfortunately this happy state of affairs is subject to two massive assumptions: that equation (3.54) is allowable (i.e. $D_2$ and $D_4$ exist) and that the results obtained digitally from a profile are acceptable simplifications of the areal digital contour of a surface.

Assuming that these two conditions are acceptable—a point which will be raised later—it is possible to use the limiting case of the discrete analysis to express any parameters hitherto not obtained and yet which are of considerable importance in the continuous theory of tribology (which is in fact what counts in the macrophysical world).

Thus the joint probability density function of the peak height and curvature is given as the probability density of an ordinate being a peak of height $z_0$ and curvature $C$ shown as $p(C, z_0|2\text{peak})$ where

$$\rho(C, z_0 \mid \text{peak}) = \frac{C}{[2\pi\, D_4(D_4 - D_2^2)]^{1/2}} \exp\left|-\frac{(D_4 z_0^2 + 2 D_2 z_0 C + C^2}{2(D_4 - D_2^2)}\right| \qquad \text{for} \quad C > 0 \qquad (3.55)$$

Also, the correlation coefficient between peak height and curvature is

$$\text{corr}\,(C, z_0 \mid \text{peak}) = \frac{(-D_2)}{(D_4)^{1/2}}\left(\frac{2 - \pi/2}{1 + (1 - \pi/2)D_2^2/D_4}\right)^{1/2} \qquad (3.56)$$

amongst others, which demonstrates that useful results can be obtained by moving between the discrete and continuous cases for surface metrology.

### 3.2.6.4 Summary

What has been indicated is how many of the most useful surface parameters of surface roughness vary in the general sense as a function of the digital (discrete) parameters. Hopefully, this gives an insight into the possible relationships that might be obtained on real surfaces. However, obviously the results in the foregoing section pose the question of whether the correlation function is measured and from this whether the measured values of the various tribological parameters of interest are predicted or should the actual parameters be measured directly from the surface using a data logging system and a computer? The former is simple and less variable because the expressions given in terms of the correlation coefficients are concerned with the average behaviour of surfaces, assuming them to be nominally Gaussian in statistics. From the formulae given, which illustrate the way in which the parameters can be estimated from the correlation coefficients, it has to be said that the specific behaviour of different surfaces will depend on the actual type of correlation function. For this reason, to estimate the relationships between parameters it is necessary to assume a correlation function which approximates to the surface in question. This is a mathematical modelling problem of surfaces.

However, as a first guess the exponential correlation function is usually sufficiently adequate for values away from the origin, that is $\rho = 0.7$. It also fits in with fractal-type surfaces. None the less, this apparently severe restriction need not be. It is considered that in many applications of surfaces the function depends on two surfaces and not one: it is the properties of the gap between them that are important. Thus one advantage of specifying surfaces in terms of the correlation coefficients is that, at least to a first approximation, the gap properties (neglecting actual separation) can be estimated from the additive properties of the mating surfaces.

Thus if $\sigma_A^2$ is the variance ($R_q^2$) value of surface A and $\rho_A(h) + \rho_A(2h)$ are its coefficients, and similarly for surface $B$, then

$$\begin{aligned}
\sigma_g^2 &= \sigma_A^2 + \sigma_B^2 \\
\rho_g(h) &= [\sigma_A^2 \rho_A(h) + \sigma_B^2 \rho_B(h)\big/\sigma_g^2 \\
\rho_g(2h) &= [\sigma_A^2 \rho_A(2h) + \sigma_B^2 \rho_B(2h)\big/\sigma_g^2.
\end{aligned} \qquad (3.57)$$

Use of the density of peaks and zero crossings instead of correlation coefficients is not attractive because they are not additive. See chapter 7 for a full examination of these issues.

Using the nomenclature $n_0$ and $m_0$ as before, the comparison with (3.57) becomes

$$\sigma_g^2 = \sigma_A^2 + \sigma_B^2$$
$$n_{0g} = [(\sigma_A^2 n_{0A}^2 + \sigma_B^2 n_{0B}^2)/\sigma_g^2]^{1/2}$$
$$m_{0g} = \frac{\sigma_A^2 n_{0A}^2 m_{0A}^2 + \sigma_A^2 n_{0B}^2 m_{0B}^2}{\sigma_A^2 n_{0A}^2 + \sigma_B^2 n_{0B}^2}.$$

(3.58)

Any errors tend to be cumulative. (More will be said about this aspect of characterization in chapter 7.)

The latter method of obtaining tribological data (i.e. the formal data logging method) gives better estimates of extreme properties in the sense that the odd, very high peak or valley may be picked up, but it does tend to be more time consuming to gather the necessary data. Obviously in the former method the correlation coefficients have to be determined. Another way to be discussed is concerned with estimation by other than stylus means.

### Areal ( 3D ) Filtering

The same basic techniques can be used for areal (3D) filtering as were used in profile filtering. However, areal manipulation offers opportunities for a better understanding of the manufacturing process and the fuction of the workpiece than was possible with profile information.

### General

The starting point is the surface data $f(xy)$. This has a frequency (wavelength) content given by its Fourier transform

$$F(w,v) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x,y)\exp(-j(xw + yv))dx \ dy$$

(3.59)

$F(w, v)$ can be modified by a filter function, say, $H(w, v)$ so that the filtered areal data is now $F'(w, v)$ where

$$F'(w,v) = \int\int F(w,v)H(w,v)dx \ dy$$

(3.60)

this can be inversely transformed to give the filtered areal data $f'(x,y)$

$$f'(x,y) = \frac{1}{4\pi^2} \int\int F'(w,x)\exp(+j(xw + yv))dw \ dv$$

(3.61)

The data need not be analysed using the Fourier transform. Dong [7] has been using wavelet transforms to examine areal (3D) data.

The way to calculate the areal spectrum is by carrying out a spectrum over all the profile $y$ values of which there may be $N$ and finding the average and then carrying out an equivalent operation in the orthogonal direction and averaging [8]. The fast Fourier transform is used for these calculations. The resultant output is a two dimensional FFT.

The same arguments apply to areal filtering as to profile filtering, namely that a conventional filter does not need to have knowledge of the surface topography in order to be applied effectively. The penalty for this freedom is losing data at the boundaries because of the finite size of the impulse response of the filter. (Fig. 3.17). It is possible to play some tricks with the data in order to salvage some of the 'lost' area. One is simply to fold A to A″ and repeat the data (Fig. 3.18)



**Figure 3.17** Area available for assessment.



**Figure 3.18**

The shaded area and the repeat areas A′B′C′D′, A″B″C″D″ can now be utilized.

Unfortunately, although there is no discontinuity between the levels at A A′ and B B″ and A A‴ B B‴ etc, the differentials at the boundaries will be discontinuous so that slope and curvature data will be corrupt. The first differentials could also be made to fit on a first order surface but higher differentials cannot.

Polynomials can be effectively fitted to the areal data providing that the longer wavelengths can be matched to the order of the polynomial which presupposes prior knowledge. Failure to do this will produce spurious results.

There has been some interest in using spline functions e.g. [9]. It is instructive to carry out the spline fit for a profile and the extend it to the area fit. Following Rhadakrishnan [9] and using his nomenclature the following method can be adopted.

The equation of the spline can be written as

$$P(x_i) = \sum_{i=1}^{n} N_{ik}(x)B_i$$

(3.62)

$N_{ik}$ are the B spline basis functions given by

$$Nk = \begin{cases} 1 . x_i \le x \le x_{i+1} \\ 0 \qquad \text{otherwise} \end{cases}$$

$$N_{ik}(x) = \frac{x - x_i}{x_{i+k} - x_i} \cdot N_{ik-1}(x) + \frac{x_{i+k} - x}{x_{i+k} - x_i} \cdot N_{i+1, k-1}(x) \qquad (3.63)$$

where        $x_i$ is the 'i'th knot value

                $B_i$ is the 'i'th control value = values at knot

                $N$ is the number of control values

                $k$ is order of the curve

If $x_i \, y_i$ are data points

$$\frac{\partial}{\partial B_i} \sum [(p(x_i) - y_i)]^2 = 0$$

The canonical equations obtained match the number of unknowns so that the $B$s can be found. When simplified this becomes

$$\sum_i N_{ik}(x_i) \sum_{i_2=1}^{n} B_{i_2} N_{i_2 k_2}(n_i) = \sum y_1 N_{ik}(x_i) \qquad (3.64)$$

As $i$ is stepped to $n$, $n$ equations result from which the $n$ control values of the mean curve are obtained. For two dimensions, equation (3.64) can be extended to give [9,59]

$$P(x_i y_i) \sum_i [P(x_i y_i - Z_i)]^2 = 0 \qquad (3.65)$$

$$= \sum \sum z_i N_{i_1 l_n}(x_i) M_{j_1 \cdot k_n} y(i) \qquad (3.66)$$

Varying $i_1$ from 1 to $n$ and $j_1$ from 1 to $M$ gives a set of $m \times n$ equations which, when solved, give the control values of the mean surface. Parameters from this mean surface are extensions of the 1 dimension profile case. So,

$$S_a = \frac{1}{A} \int_A \int \left| z(\partial c, y) \right| dx \, dy \qquad (3.67)$$

$$= \frac{1}{MN} \sum_{j=1}^{N} \sum_{i=1}^{M} | z(i, j) | \qquad (3.68)$$

Other parameters can also be obtained similarly. It should be remembered that the true spline is a cubic equation and is meant to be a line of minimum potential energy. It was originally used in ship building. An elastic or bendy piece of wood was pressed against a set of $n$ stanchions. The resulting curve made contact at $n$ points. The curve in between restraints could be measured and drawn. The shape of the strip of wood should be such that it has a minimum potential energy. Mathematical splines are derived from this.

### 3.2.7  Areal ( two-dimensional )  digital measurement of surface roughness parameters

#### 3.2.7.1  General

Note here that areal is taken to be the three-dimensional picture having two independent variables.

The term 3D is sometimes used instead of areal.

In the earlier sections on the characterization of surface roughness it became clear that a major objective is to measure the areal dimensional characteristics of the surface rather than to restrict the discussion to single profiles. This expansion of subject immediately brings into focus such things as specification and measurement of lay, the relationship between the moments of the spectrum of the area versus those measured from one or more profiles. Exactly the same situation exists in the digital domain. It has already been highlighted that there are differences between the theoretical values of features for profile and areal measurements according to Longuet-Higgins [10] and Nayak [11]. However, the questions emerge as to the differences which might arise between profile discrete measurement on the one hand and the convergence of areal discrete estimates of surfaces to the theoretical values on the other.

Exactly the same sort of calculation can be carried out as before, yielding some very interesting and important results. Two attempts have been made to analyse the areal discrete problem for a random surface, the first by Whitehouse and Phillips [4] and the second by Greenwood [5]. The former results will be given here, yet the more digestible results of the latter will also be included for comparison. Both approaches can be considered to be a typology of surfaces in their own right (see chapter 2).

So far profile information has been examined. It is understood that such information is only a convenient approximation to the real surface geometry. Exploring the areal properties on a continuous basis in the way that Longuett-Higgins and Nayak have is most informative but does not actually relate directly to what is measured. All measurements are now carried out digitally. It will be observed that areal discrete measurements cannot converge to the values for a continuous waveform because of the discontinuous spatial bandwidth of any sampling pattern. In certain circumstances, for example isotropic surfaces or when the two perpendicular components are independent, measuring the spectra either in one direction in isotropy or in two directions for independent lay. These enable Nyquist criteria to be tested if the pass conversion is possible. It is the purpose of this next section to investigate the nature of the differences that can and do occur in parameters which are functionally important, such as summit curvatures and height distributions.

Some of the theory used in investigating discrete properties overlaps the theory used for surface characterization using the discrete parameters reported in chapter 2. The analytical results are necessarily complicated but are absolutely essential if the results for surfaces obtained by experimenters are ever to be made consistent with each other.

Also, whereas the idea of tribology parameter prediction using the correlation function for a profile might not seem very attractive, it most certainly is for areal estimation because of the very great difficulty in measuring over an area and maintaining geometric fidelity between traces. Simply taking disjointed profile tracks is not good enough—each trace has to be integrated into the overall mapping scheme in height and starting point. Another point which adds to the complexity is that, in areal measurement, there is not necessarily a straightforward grid system for data ordinates to be used. Before discussing the alternative methods, the analysis for the areal equivalent of the three-point method will be given.

In one dimension a continuous definition of a local maximum or peak on a profile only requires one first-order and one second-order differential, and only three ordinates are needed for the discrete definition of a peak. However, in two dimensions, the continuous definition of a local maximum (or summit, in the terminology of Nayak) requires two first-order and three second-order differentials, and a minimum of five ordinates are usually considered necessary for the discrete definition of a summit. Sayles and Thomas [13] gave two discrete definitions of a summit, one using the five nearest-neighbour ordinates, and the other nine. Here a number of possibilities will be considered; to start with, the conventional five-ordinate definition will be used. The definition states that an ordinate at height $z_0$ above an arbitrary datum is a five-point summit if it is

higher than the four neighbouring ordinates, which are at a distance $h$ from it, on the two orthogonal Cartesian axes. If the differences between $z_0$ and each of the four neighbouring ordinates are denoted by $s_1$, $s_2$, $s_3$ and $s_4$, then the condition for the ordinate $z_0$ to be a summit is that $s_1$, $s_2$, $s_3$ and $s_4$ are all positive.

The summit density is the proportion of ordinates that are summits. The summit height is the height $z_0$ of an ordinate that is a summit.

It seems natural to extend to two independent dimensions the discrete definition of Whitehouse and Archard [14] for peak curvatures, which was given by

$$C_h = (s_1 + s_2)h^2 \tag{3.69}$$

where $s_1$ and $s_2$ are defined later. This extension uses the average orthogonal peak curvatures, and gives the discrete definition of five-point summit curvature as

$$C_h^{(2)} = \tfrac{1}{2}(s_1 + s_2 + s_3 + s_4)/h^2 \tag{3.70}$$

which again is a linear combination of the ordinates.

### 3.2.7.2 The expected summit density and the distributions of summit height and curvature

The distributions of peak height and curvature for a profile in one dimension were obtained by Whitehouse and Phillips [4] for a surface with ordinates from a Gaussian (normal) distribution. (A multivariate normal distribution for a vector $Z$ will be denoted by $Z \sim N[\mu; V]$ where $\mu$ is the vector of means and $Z$ is the variance-covariance matrix, and the probability density function is given by $\varphi^{(m)}(z' - \mu'; V)$.) They obtained these results from the theory of truncated random variables. This was because the peak height distribution is the conditional distribution of $z_0$, the profile height, given that $s_1$ and $s_2$ are positive, or, in other words, the distribution of $z_0$ conditional on $s_1$ and $s_2$ being positive, where

$$s_1 = z_0 - z_{-1} \qquad s_2 = z_0 - z_1 \tag{3.71}$$

and $z_{-1}$ and $z_1$ are the preceding and succeeding ordinate values, respectively, on the profile. Similarly the distribution of peak curvature is the conditional distribution of $C_h$ given that $s_1$ and $s_2$ are positive. Hence the results of Whitehouse and Phillips [4] can be obtained by using the results for truncated random variables, with $m = 2$, $Y_0 \equiv Z_0$, $X = (2 - 2\rho_1)^{-1/2}(s_1, s_2)'$, $d = (\tfrac{1}{2} - \tfrac{1}{2}\rho_1)^{1/2}$ and

$$V = \begin{pmatrix} 1 & \tfrac{1}{2}(1 - 2\rho_1 + \rho_2)/(1 - \rho 1) \\ \tfrac{1}{2}(1 - 2\rho_1 + \rho_2)/(1 - \rho 1) & 1 \end{pmatrix} \tag{3.72}$$

where $\rho_1 = \rho(h)$, $\rho_2 = \rho(2h)$.

The derivations can also be used to obtain the more important two-dimensional (areal) distributions of five-point summit height and curvature. For this analysis the surface height measurements will be assumed to have a multivariate normal distribution (MND) and, because the surface is assumed to be isotropic, the distribution properties of a profile are invariant with respect to the direction of the profile. Hence

$$(z_0, (2 - 2\rho_1)^{-1/2}(s_1, s_2, s_3, s_4) \sim N[0; V_5] \tag{3.73}$$

where

$$\mathbf{V}_5 = \begin{pmatrix} 1 & d & d & d & d \\ d & 1 & a & b & b \\ d & a & 1 & b & b \\ d & b & b & 1 & a \\ d & b & b & a & 1 \end{pmatrix} \tag{3.74}$$

$$a = \tfrac{1}{2}(1 - 2\rho_1 + \rho_2)\big/(1 - \rho_1) \qquad \text{and} \qquad b = \tfrac{1}{2}(1 - 2\rho_1 + \rho_3)\big/(1 - \rho_1) \tag{3.75}$$

where $\rho_1 = \rho(h)$, $\rho_2 = \rho(2h)$, $\rho_3 = \rho(h\sqrt{2})$, and $\rho(t)$ is the correlation coefficient between ordinates a distance $t$ apart.

If $\gamma_5$ is the event ($s_1 > 0$, $s_2 > 0$, $s_3 > 0$, $s_4 > 0$) then the expected five-point summit density is the probability of $Y_5$ occurring, and the distribution of five-point summit height is the conditional distribution of $Z_0$ given that $Y_5$ has occurred. These can be obtained from $m = 4$,

$$\mathbf{X} = (2 - 2\rho_1)^{-1/2}(s_1, s_2, s_3, s_4)' \tag{3.76}$$

$$d = (\tfrac{1}{2} - \tfrac{1}{2}\rho_1)^{1/2} \tag{3.77}$$

and $\mathbf{V}$ obtained from $\mathbf{V}_5$ by removing the first row and column (and denoted by $\mathbf{V}_4$) so that

$$\lambda = 1 + a + 2b. \tag{3.78}$$

Then the probability density function of the five-point summit height distribution is given by

$$p(z_0 \mid Y_5) = \frac{\varphi^{(4)}(z_0)[(1 - \rho_1)\big/(1 + \rho_1)]^{1/2}; \mathbf{V}_c)\varphi(z_0)}{\varphi^4(0, \mathbf{V}_4)} \tag{3.79}$$

where

$$\mathbf{V}_c = \begin{pmatrix} 1 & a_c & b_c & b_c \\ a_c & 1 & b_c & b_c \\ b_c & b_c & 1 & a_c \\ b_c & b_c & a_c & 1 \end{pmatrix} \tag{3.80}$$

$$a_c = (\rho_2 - \rho_1^2)\big/(1 - \rho_1^2) \qquad \text{and} \qquad b_c = (\rho_3 - \rho^2)\big/(1 - \rho^2). \tag{3.81}$$

The expected (average or mean) five-point summit height is given by $E(z_0|Y_5)$, where $E$ denotes the statistical expectation:

$$E[z_0 \mid \gamma_5] = \frac{2[(1 - \rho_1\big/\pi)]^{1/2}\Phi^{(3)}(0; \mathbf{B}_4)}{\Phi^{(4)}(0; \mathbf{V}_4)} \tag{3.82}$$

where

$$
\mathbf{B}_4 = \begin{pmatrix} 1 - b^2 & a - b^2 & b(1-a) \\ a - b^2 & 1 - b^2 & b(1-a) \\ b(1-a) & b(1-a) & 1 - a^2 \end{pmatrix}.
\tag{3.83}
$$

From equation (3.82) it can be seen that the expected five-point summit height depends on two orthant probabilities $\Phi^{(3)}(0; \mathbf{B}_4)$ and $\Phi^{(4)}(0; \mathbf{V}_4)$, which have to be evaluated. From Whitehouse and Arehard [14], $\Phi^{(3)}(0; \mathbf{B}_4)$ is given by

$$
\Phi^{(3)}(0; \mathbf{B}_4) = \tfrac{1}{2} - (4\pi)^{-1} \cos^{-1}\left[(a - b^2)\big/(1 - b^2)\right]
$$

$$
- (2\pi)^{-1} \cos^{-1}\left\{ b(1-a)[(1 - a^2)(1 - b^2)]^{-1/2} \right\}.
\tag{3.84}
$$

Cheng [15] has evaluated $\varphi^{(4)}(0; \mathbf{V}_4)$, so using this result the expected five-point summit density has been found, using Plackett's [16] method. As this orthant probability only depends on the parameters $a$ and $b$, it will be denoted by $\varphi^{(4)}[a, b]$. Then

$$
pr(Y_5) = \Phi^{(4)}[a, b] = \Phi^{(4)}[a, o] + \int_0^b \frac{\delta\, \Phi^{(4)}[a, t]}{\delta t}\, \mathrm{d}t
$$

$$
= [\tfrac{1}{4} + (2\pi)^{-1} \sin^{-1} a]^2 + (2\pi)^{-1} \sin^{-1} b
$$

$$
+ \pi^{-2} \int_0^b (1 - t^2)^{-1/2}\, \sin^{-1}[t(1-a)(1 + a - 2t)^{-1}]\, \mathrm{d}t.
\tag{3.85}
$$

This result was given by Cheng [15] in terms of the dilogarithm function (see [11]).

The distribution of a five-point summit having a height $z_0$ conditional on a curvature $C_h^{(2)}$ is normal and is given by

$$
(Z_0 \mid C_h^{(2)} = C_h^{(2)}) \sim N[h^2(1 + a + 2b)^{-1} C_h^{(2)},\ 1 - 2(1 - \rho_1)\big/(1 + a + 2b)].
\tag{3.86}
$$

This is henceforth called the conditional distribution of summit height given curvature. Thus the expected five-point summit curvature is given by

$$
E[C_h^{(2)} \mid Y_5] = h^{-2}(1 + a + 2b) E[Z_0 \mid Y_5].
\tag{3.87}
$$

Hence, by the application of the theory of Gaussian (normal) truncated random variables, it has been possible to obtain the expectations of the five-point summit height, curvature and density in terms of the correlation coefficients. These are the basic tribological parameters required. The results have been arrived at by using a five-point model for summits (called the tetragonal model) and this implies a rectangular grid of sampled data. These results therefore can be regarded as the tribological parameters of a discrete random surface with values at the intersections of a rectangular grid. As long as the correlation is $\rho_1$ and $\rho_2$ at $h$ and $2h$ in both directions or scaled accordingly the results are exact. A wide variety of surfaces can be modelled by allowing $\rho_1$ and $\rho_2$ to vary.

### 3.2.7.3  The effect of the sampling interval and limiting results for the discrete surface

The distributions of five-point summit height and curvature have been derived in terms of the correlation coefficients between ordinates. These correlation coefficients are $\rho_1$ for ordinates a distance $h$ apart, $\rho_2$ for ordinates $2h$ apart and $\rho_3$ for ordinates $\sqrt{2}h$ apart. If the surface is isotropic and the autocorrelation function is $\rho(t)$ then

$$\rho_1 = \rho(h) \qquad \rho_2 = \rho(2h) \qquad \text{and} \qquad \rho_3 = \rho(h\sqrt{2}). \tag{3.88}$$

So $\rho_1$, $\rho_2$ and $\rho_3$ will vary as $h$ varies, depending on the shape of the autocorrelation function of the surface. As $h$ approaches zero

$$\lim_{h\to 0} \rho_1(h) = \lim_{h\to 0} \rho_2(h) = \lim_{h\to 0} \rho_3(h) = 1 \tag{3.89}$$

and as $h$ approaches infinity

$$\lim_{h\to \infty} \rho_1(h) = \lim_{h\to \infty} \rho_2(h) = \lim_{h\to \infty} \rho_3(h) = 0 \tag{3.90}$$

If $\rho_1$, $\rho_2$ and $\rho_3$ are plotted in three dimensions then, as $h$ varies, the curve will start at $(1,1,1)$ for $h = 0$ and end at $(0,0,0)$ for $h = +\infty$. In order that the matrix $V_4 - d^2 J$ is positive definite it is necessary for this curve to lie in the region bounded by $\rho_2 < 1$ and $-\frac{1}{2}(1 + \rho_2) + 2\rho_1^2 < \rho_3 < \frac{1}{2}(1 + \rho_2)$.

Results for the summit height have been obtained by Nayak [11] for the continuous surface, so it is possible to compare his results with those obtained for the discrete results of this chapter as the sampling interval $h$ approaches zero. The expected summit height depends on $\rho_1 \rho_2$ and $\rho_3$ through $a$ and $b$, and

$$\lim_{h\to \infty} a = -1 \qquad \text{and} \qquad \lim_{h\to \infty} b = 0.$$

The autocorrelation function of the surface can be approximated by

$$\rho(h) = 1 + D_2 h^2 / 2! + D_4 h^4 / 4! + O(h^4)$$

where $D_2$ and $D_4$ are the second and fourth derivatives of the autocorrelation function at the origin and

$$\eta = -D_2(D_4)^{-1/2} < \sqrt{\tfrac{2}{3}}. \tag{3.91}$$

The limiting value for the expected five-point summit height is given by

$$\lim E[z_0 \mid Y_5] = \frac{16}{\pi + 2 \ \sin^{-1}(1/3) + 4\sqrt{2}} (\tfrac{1}{2}\pi)^{1/2} \eta = 1.688(\tfrac{1}{2}\pi)^{1/2}\eta. \tag{3.92}$$

Nayak [11] showed that the expected continuous summit height for the areal case was

$$E[z \mid Y_N] = (4\sqrt{2\pi})(\tfrac{1}{2}\pi)^{1/2}\eta = 1.801(\tfrac{1}{2}\pi)^{1/2}\eta. \tag{3.93}$$

which is comparable with the expected continuous peak height for the profile

$$E[z \mid Y_{\mathrm{p}}] = (\tfrac{1}{2}\pi)^{1/2}\eta \tag{3.94}$$

a result given by Whitehouse and Phillips [4]. Then it can be seen that the limit of the expected five-point summit height (3.92) is 69% larger than the expected peak limit (3.94) as opposed to 80°/o larger than the expectation of the distribution of summit height for the continuous definition of Nayak [11] (3.92). However, this is only an overall reduction of about 6% and suggests that the discrete five-point definition may be adequate. Compare this with the seven-point result, equation (3.128) below.

It is possible to obtain $\Phi^{(4)}(z_0[(1 - \rho_1)/(1 + \rho_1)]^{1/2}; \mathsf{V}_c)$ by the methods of (3.10) and (3.11), and hence to obtain the probability density function of the limiting distribution of the five-point summit height as $h$ converges to zero. This is given by

$$\lim_{h \to 0} p(z_0 \mid Y_s) = 12\pi\varphi(z_0)(\pi + 2\sin^{-1}\tfrac{1}{3} + 4\surd 2)^{-1}$$
$$\times\{(1 - \eta^2)[\varphi(w) + w\Phi(w)]^2 + \tfrac{4}{3}(1 - \tfrac{3}{2}\eta^2)\Phi^{(2)}(w;\tfrac{1}{3} - \eta^2) \cong (1 - \eta^2)$$
$$-(1 - \eta^2)\Phi^{(2)}(w;0) + 4(1 - \eta^2)T^{(2)}(w,(2 - 3\eta^2)^{-1\cong 2})\} \tag{3.95}$$

where

$$T^{(2)}(w,\upsilon) = (2\pi)^{-1}\int_0^{\upsilon}(1 + x^2)^{-2}\,\exp[-w^2(1 + x^2)]\,\mathrm{d}x \tag{3.96}$$

and

$$w = \eta(1 - \eta^2)^{-1/2}z_0. \tag{3.97}$$

This probability density function is compared with the probability density function of continuous summit height, given by

$$p(z \mid Y_{\mathrm{N}}) = 3\eta(2 - 3\eta^2)^{1/2}\varphi(0)z\varphi((1 - \tfrac{3}{2}\eta^2)^{-1/2}z)$$
$$+3\surd 3\eta^2(z^2 - 1)\Phi(\eta[\tfrac{3}{2}\big/(1 - \tfrac{3}{2}\eta^2)]^{1/2}z)\varphi(z)$$
$$+2(1 - \eta^2)^{-1/2}\Phi(\{\tfrac{1}{2}\big/[(1 - \eta^2)(1 - \tfrac{3}{2}\eta^2)]\}^{1/2}\eta z)\varphi((1 - \eta^2)^{-1/2}z) \tag{3.98}$$

which was obtained by Nayak [11] for three values $\eta = 0$, $\sqrt{\tfrac{1}{3}}$ and $\sqrt{\tfrac{2}{3}}$, as shown in figure 3.19. For $\eta = 0$ both distributions are the standardized normal distribution.

When $\eta = \sqrt{\tfrac{2}{3}}$, the probability density function of the continuous summit height is

$$p(z \mid Y_{\mathrm{N}}) = 2\surd 3\{z^2 - 1 + 2\pi[\varphi(z)]^2\}\varphi(z) \qquad \text{for} \quad z > 0$$
$$= 0 \qquad \text{for} \quad z \le 0 \tag{3.99}$$

while the limiting distribution of the five-point summit height is given by

$$\lim p(z \mid Y_{\mathrm{N}}) = \frac{8\pi\varphi(z_0)[(2z_0^2 - 1)(\Phi(\surd 2 z_0) - \tfrac{1}{2}) + \surd 2 z_0\varphi(\surd 2 z_0)]}{\pi + 2\,\sin^{-1}\tfrac{1}{3} + 4\surd 2} \qquad \text{for} \quad z_0 > 0$$
$$= 0 \qquad \text{for} \quad z_0 \le 0. \tag{3.100}$$

**Figure 3.19** The probability density function (PDF) of the distribution of summit height (full line) and the limiting distribution of the five-point summit height (broken line) for $\eta = 0$, $\sqrt{1/3}$ and $\sqrt{2/3}$.

Nayak [11] used the definition of mean summit curvature $K_m$ given by Sokolnikoff as minus the average of the second partial derivatives in orthogonal directions. With this definition the conditional distribution of continuous summit height given the curvature is a normal distribution with

$$E_p[z,\ K_m,\ Y_N] = \tfrac{3}{2}\eta\, K_m (D_4)^{-1/2} \tag{3.101}$$

and

$$\mathrm{var}(z,\ K_m,\ Y_N) = 1 - \tfrac{3}{2}\eta^2 \tag{3.102}$$

This is also the limit of the conditional distribution of $z_0$ given $C_h^{(2)}$. Hence, the limit as $h$ tends to zero of the expected five-point summit curvature will be 6% smaller than the expectation of the continuous summit curvature. Thus the operational procedure of sampling in a plane and taking the sample interval to zero gives different results from the continuous values for the surface!

Greenwood [12] approached the problem of profile measurement in the same basic way using the multi-normal distribution from the joint distribution of height $z$, slope $m$ and curvature $k$:

$$p(z,\ m,\ k) = \frac{1}{2\pi^{3/2}\sqrt{1-r^2}}\exp\!\left(\frac{-1}{2(1-r_2)}(\xi^2 + t^2 - 2\xi t)\right)\exp\!\left(\frac{-s^2}{2}\right) \tag{3.103}$$

where $\xi = z/\sigma$, $s = m/\sigma_m$ and $t = -k/\sigma_k$, and $\xi$, $s$ and $t$ are the standardized height, slope and curvature $z$, $m$ and $k$ normalized to their standard deviations $\sigma$, $\sigma_m$ and $\sigma_k$ respectively. In equation (3.103) $r$ is $\sigma_m^2/\sigma\sigma_k$ and corresponds to the variable $\alpha = m_0 m_4/m_2^2$ as $r^{-2}$.

Greenwood introduces another variable, $\theta$, to replace in part the sampling interval defined by

$$\sin\theta = h\sigma_k/2\sigma_m \tag{3.104}$$

where $\sigma_m$ is the standard deviation of the slope and $\sigma_k$ that of curvature. If $h$ is the sampling interval, this transcription from $h$ and the correlation coefficients $\rho(0)$, $\rho(h)$ and $\rho(2h)$ preferred by Whitehouse and Phillips for a profile makes some of the subsequent equations simpler in form. However, the concept of sampling interval is masked, which is a disadvantage to the investigator.

Formulae for the joint distributions of height and curvature are found, as are those for summit curvature distributions (figure 3.20).

The reader is requested to consult the paper by Greenwood [12] to decide the simplest approach to the areal problem. However, it is gratifying to note that the results and predictions of both methods are compatible.



**Figure 3.20** Effect of surface character on surface features.

### 3.2.8 Patterns of sampling and their effect on discrete properties

#### 3.2.8.1 Comparison of three-, four-, five-and seven-point analysis of surfaces

An investigation of the various sampling pattern options open to a researcher will be given here to see if there is any advantage to be had by changing from the conventional rectilinear sampling pattern. The real issue is to find what sampling pattern best covers the area and picks out summits. Obviously, there are instrumental factors that have to be taken into account as well as purely theoretical ones. One such factor is the ease with which a complex sampling pattern can be achieved by hardware comprising a specimen table driven by two orthogonal, motorized slides.

One thing that emerges from these sampling procedures is that the properties so revealed are substantially different from each other.

#### 3.2.8.2 Four-point sampling scheme in a plane

A comparison of possible sampling schemes for sampling in an areal plane will now be presented. They are illustrated in figure 3.21 and in figure 4.128. In all cases the distance between measurements is $h$.

First, there is sampling along a straight line (from a profile of the surface). This sampling scheme only takes measurements in one dimension of the plane. This has been presented for completeness and because it

**Figure 3.21** Sampling patterns: (*a*) three points (digonal), $k = 2$; (*b*) four points (trigonal), $k = 3$; (*c*) five points (tetragonal), $k = 4$.

was the first case considered by Whitehouse and Archard [9] and Whitehouse and Phillips [4]. In figure 3.21(*a*) it is illustrated with $k = 2$.

Second, a sampling scheme could be used that would take measurements at the vertices of a hexagonal grid. The summit properties could be defined using four ordinates, that is the measurement at a vertex and the three adjacent ordinates at a distance $h$ from the chosen vertex. This is the case when $k = 3$ and is referred to as the trigonal symmetry case.

In order to produce such a hexagonal grid pattern on a surface it would be necessary to sample along parallel lines separated alternately by a distance of $\frac{1}{2}h$ and $h$. The spacing between ordinates along a line would be $\sqrt{3}h$ but the position at which the first ordinate is measured would be 0 for the $(4j - 3)$rd and $4j$th lines and $\frac{1}{2}\sqrt{3}h$ for the $(4j - 2)$nd and $(4j - 1)$st lines for $j \geq 1$. This is illustrated in figure 3.21(*b*). Alternatively, it would be possible to sample along parallel lines a distance of $\frac{1}{2}\sqrt{3}h$ apart, but this would involve a different position for the first ordinates and the spacing between ordinates would alternatively be $h$ and $2h$.

Third, there is sampling on a square grid. This was considered by Whitehouse and Phillips [4] and Greenwood [5] and will be referred to as the tetragonal symmetry case. It is illustrated in figure 3.21(*c*) with $k = 4$. The sampling scheme requires sampling along parallel lines separated by a distance $h$ and with a spacing between ordinates along a line of $h$.

Results will also be given for the hexagonal grid or trigonal symmetry and the hexagonal case with $k = 6$. For purposes of comparison, the cases when $k = 2, 3$ and $4$ will also be considered from earlier. The notation will be used below. If the $m$ random variables $x = (x_1, x_2, ..., x_m)$ have a joint multivariable Gaussian (normal) distribution with mean $\mu$ and variance-covariance matrix $V$ then this is denoted by $x - N[\mu, V]$. Also the convention of using an upper-case letter for a random variable and a lower-case letter for a realization of the random variable will be followed, as in the previous section.

*The hexagonal grid in the trigonal symmetry case*

Results have been obtained for the probability density function and expectation (mean) of peak (or summit) height, the density of summits and the expected peak (or summit) curvature in the cases when $k = 2$ by Whitehouse and Phillips [17] and when $k = 4$ by Whitehouse and Phillips [4] and Greenwood [5]. The results for the hexagonal grid ($k = 3$) in the trigonal symmetry and the hexagonal case $k = 6$ will now be given [14]. These can be obtained from the general results of truncated random variables in the appendix of Whitehouse and Phillips [4].

For measurements with four ordinates let $z_0$ be the height of the central ordinate and $s_1$, $s_2$ and $s_3$ be the differences between this ordinate and the three adjacent ordinates at a distance $h$. The ordinate $z_0$ will be defined to be a four-point summit if $s_1$, $s_2$ and $s_3$ are all positive. By analogy with the three-point and five-point definitions of curvature the discrete definition of four-point curvature is

$$C = \frac{2(s_1 + s_2 + s_3)}{3h^2}.$$

(3.105)

Assuming that the surface height measurements have a multivariate Gaussian distribution and that the surface is isotropic then

$$(z_0, 2 - 2\rho_1)^{-1/2}(s_1, s_2, s_3)) \sim N[0; V_4].$$

(3.106)

The probability that $s_1$, $s_2$ and $s_3$ are all positive gives the probability that an ordinate is a four-point summit. Thus, again using the nomenclature of Cheng [15],

$$V_4 = \begin{pmatrix} 1 & d & d & d \\ d & 1 & a & a \\ d & a & 1 & a \\ d & a & a & 1 \end{pmatrix}$$

(3.107)

with

$$d = \left(\tfrac{1}{2} - \tfrac{1}{2}\rho_1\right)^{1/2} \quad \text{and} \quad a = \tfrac{1}{2}\left(1 - 2\rho_1 + \rho\sqrt{3}\right) \big/ \left(1 - \rho_1\right)$$

(3.108)

where $\rho_1 = \rho(h)$ and $\rho\sqrt{3} = \rho(\sqrt{3}h)$.

If $Y_4$ is the event $(s_1 > 0, s_2 > 0, s_3 > 0)$ then the distribution of four-point summit height is the conditional distribution of $z_0$ given that $Y_4$ has occurred. This can be obtained using the results of Whitehouse and Phillips [4] with $m = 3$:

$$\mathbf{X} = (2 - 2\rho_1)^{-1/2}(s_1, s_2, s_3) \tag{3.109}$$

$$d = \left(\tfrac{1}{2} - \tfrac{1}{2}\rho_1\right)^{1/2} \tag{3.110}$$

and

$$\mathbf{V} = \mathbf{V}_3 = \begin{pmatrix} 1 & a & a \\ a & 1 & a \\ a & a & 1 \end{pmatrix} \tag{3.111}$$

so that

$$\lambda = 1 + 2\alpha. \tag{3.112}$$

Then the probability density function of the four-point summit height distribution is given by

$$p(z_0 \mid Y_4) = \frac{\Phi^{(3)}(z_0[(1 - \rho_1)/(1 + \rho_1)]^{1/2}; \mathbf{V}_c)\Phi(z_0)}{\Phi^{(3)}(0; \mathbf{V}_3)} \tag{3.113}$$

in exactly the same way as for the tetragonal case where

$$\mathbf{V}_c = \begin{pmatrix} 1 & a_c & a_c \\ a_c & 1 & a_c \\ a_c & a_c & 1 \end{pmatrix} \tag{3.114}$$

and

$$a_c = (\rho\sqrt{3 - \rho_1^2})/(1 - \rho_1^2) \tag{3.115}$$

$\Phi^{(n)}(z'; \mathbf{V})$ is the cumulative distribution function at $z'$ of the $n$-dimensional multivariate Gaussian distribution with zero expectation and variance-covariance matrix $\mathbf{V}$. $z$ is used for $(z, z_2, z_3, \ldots) = z'$. $\Phi(x)$ is the probability density function of the univariate standard Gaussian distribution.

The denominator of equation (3.112) is the orthant probability which gives the probability that an ordinate is a four-point summit. Hence

$$pr(Y_4) = \Phi^{(3)}(0; \mathbf{V}_3) = \tfrac{1}{2} - \tfrac{3}{4}(\pi)^{-1}\cos^{-1}a \tag{3.116}$$

The expected (mean) four-point summit height is given by

$$E[z_0 \mid Y_4] = \frac{3(1 - \rho_1)^{1/2}\Phi^{(2)}(0; \mathbf{B}_3)}{2\sqrt{\pi}\Phi^{(3)}(0; \mathbf{V}_3)} \tag{3.117}$$

where

$$B_3 = \begin{bmatrix} 1 - a^2 & a(1-a) \\ a(1-a) & 1 - a^2 \end{bmatrix}. \tag{3.118}$$

Hence

$$E[z_0 \mid Y_4] = 3(1 - \rho_1)^{1/2} \left[ \frac{1}{4} + \frac{1}{2\pi} \sin^{-1} \left( \frac{1 - 2\rho_1 + \rho_{\sqrt{3}}}{3 - 4\rho_1 + \rho_{\sqrt{3}}} \right) \right] \bigg/ \left\{ 2\sqrt{\pi} \left[ \frac{1}{2} - \frac{3}{4\pi} \cos^{-1} \left( \frac{1 - 2\rho_1 + \rho_{\sqrt{3}}}{2 - 2\rho_1} \right) \right] \right\}. \tag{3.119}$$

The distribution of the height $z_0$ of a four-point summit conditional on curvature $C$ is Gaussian with an expectation given by

$$E[z_0 \mid C = c, \ Y_4] = \frac{3h^2(1 - \rho_1)c}{4(2 - 3\rho_1 + \rho_{\sqrt{3}})} \tag{3.120}$$

and variance given by

$$\text{var}[z_0 \mid Y_4] = \frac{1 - 3\rho_1 + 2\rho_{\sqrt{3}}}{2(2 - 3\rho_1 + \rho_{\sqrt{3}})}. \tag{3.121}$$

This is the same as the distribution of the height $z_0$ of an ordinate conditional on the four-point curvature but not conditional on the ordinate being a summit, and is a result which holds for the three values of $k = 2, 3$ and $4$. This is because $V_k$ is of the form

$$V_k = \begin{pmatrix} 1 & dl' \\ dl & \overline{V} \end{pmatrix} \tag{3.122}$$

where $V$ is a correlation matrix with a constant row (column) sum ensured by cofactors $dl$ and $dl'$. This result enables the expected $(k+1)$-point peak (or summit) curvature to be obtained from the expected $(k+1)$-point peak (or summit) height.

Hence the expected four-point summit curvature is given by

$$
\begin{aligned}
E[C \mid Y_4] &= \frac{4(2 - 3\rho_1 + \rho_{\sqrt{3}})}{3h^2(1 - \rho_1)} E[Z_0 \mid T_4] \\
&= 2(2 - 3\rho_1 + \rho_{\sqrt{3}}) \left[ \frac{1}{4} + \frac{1}{2\pi} \sin^{-1} \left( \frac{1 - 2\rho_1 + \rho_{\sqrt{3}}}{3 - 4\rho_1 + \rho_{\sqrt{3}}} \right) \right] \bigg/ \\
&\quad \left\{ h^2 [\pi(1 - \rho_1)]^{1/2} \left[ \frac{1}{2} - \frac{3}{4\pi} \cos^{-1} \left( \frac{1 - 2\rho_1 + \rho_{\sqrt{3}}}{2 - 2\rho_1} \right) \right] \right\}.
\end{aligned}
\tag{3.123}
$$

It is also possible to obtain the following simple connection between the variances of the $(k+1)$-point peak (or summit) height and curvature:

$$\text{var}[Z_0 \mid T_{k+1}] = \text{var}(Z_0 \mid C, Y_{k+1}) + \left[ \frac{E[Z_0 \mid C = C, \ T_{k+1}]}{C} \right]^2 \text{var}(C, T_{k+1}). \tag{3.124}$$

This relation was shown by Whitehouse and Phillips [4] for $k = 1$ (with $z_0 \equiv y_0$) and by Greenwood [5] for $k = 4$.

So, by the application of the theory of Gaussian truncated random variables, it has been possible to obtain connections between the expectations and variances of four-point summit and curvature. Also, it is straightforward to show, without getting involved in the calculation, that the probability that an ordinate is a summit for a hexagonal sampling model is with

$$pr(Y_7) = \Phi^{(6)}(0; \mathsf{V}_6) \tag{3.125}$$

$$\mathsf{V}_6 = \begin{pmatrix}
1 & \frac{1}{2} & b & c & b & \frac{1}{2} \\
\frac{1}{2} & 1 & \frac{1}{2} & b & c & b \\
b & \frac{1}{2} & 1 & \frac{1}{2} & b & c \\
c & b & \frac{1}{2} & 1 & \frac{1}{2} & b \\
b & c & b & \frac{1}{2} & 1 & \frac{1}{2} \\
\frac{1}{2} & b & c & b & \frac{1}{2} & 1
\end{pmatrix} \tag{3.126}$$

where

$$b = (1 - 2\rho_1 + \rho_{\sqrt{3}})/(2 - 2\rho_1) \qquad c = (1 - 2\rho_1 + \rho_2)/(2 - 2\rho_1) \qquad \rho_T = \rho(t\Delta h) \tag{3.127}$$

and $\Delta = 2^{1/2}/3^{1/4}$, giving

$$E[z \mid Y_7] = \frac{3\sqrt{1 - \rho_1}\pi\varphi^5(0, \mathsf{B}_6)}{pr(Y_7)} \tag{3.128}$$

where $\mathsf{B}_6$ is the variance-covariance matrix of the conditional distribution of the differences $s_1, s_2, ..., s_5$ given $s_6$, from which

$$E[z \mid Y_7] = \frac{3\sqrt{(1 - \rho_1/\pi}}{pr(Y_7)} \ \frac{[6\pi - 12\tan^{-1}(1/\sqrt{2}) - \sqrt{2}\pi]}{8\sqrt{2}\pi^2} \Delta \left( \frac{\mathsf{D}_4}{-\mathsf{D}_2} \right)^{1/2} h. \tag{3.129}$$

From this the limiting values as $h \to 0$ can be found and inserted in table 3.2.

**Table 3.2** Expected summit (peak) density.

| Model | $k$ | $pr(T_{k+1})$ Limit as $h \to 0$ | Expected density | | |
|---|---|---|---|---|---|
| | | | Limit as $h \to 0$ | | Limit as $h \to \infty$ |
| Three points | 2 | $\dfrac{1}{2\pi}\left(\dfrac{D_4}{-D_2}\right)^{1/2} h$ | $\dfrac{1}{2\pi}\left(\dfrac{D_4}{-D_2}\right)^{1/2} = D_{\text{peak}}$ | | $\dfrac{1}{3h} = \dfrac{0.333}{h}$ |
| Four points | 3 | $\dfrac{1}{6\pi}\left(\dfrac{D_4}{-D_2}\right)h^2 = 0.0531\left(\dfrac{D_4}{-D_2}\right)h^2$ | $\sqrt{3D_{\text{sum}}} = 1.732\,D_{\text{sum}}$ | | $\dfrac{1}{4h^2} = \dfrac{0.25}{h^2}$ |
| Five points | 4 | $\dfrac{[\pi + 2\sin^{-1}(\frac{1}{3}) + 4\sqrt{2}]}{24\pi^2}\left(\dfrac{D_4}{-D_2}\right)h^2$ | $\dfrac{\sqrt{3}[\pi + 2\sin^{-1}(\frac{1}{3}) + 4\sqrt{2}]}{4\pi}D_{\text{sum}}$ | | $\dfrac{1}{5h^2} = \dfrac{0.2}{h^2}$ |
| | | $= 0.400\left(\dfrac{D_4}{-D_2}\right)h^2$ | $= 1.306\,D_{\text{sum}}$ | | |
| Seven points | 6 | $\dfrac{[\pi + 6(\sqrt{3} - 1)]}{24\pi^2}\left(\dfrac{D_4}{-D_2}\right)h^2$ | $\dfrac{\sqrt{3}[\pi + 6(\sqrt{3} - 1)]}{4\pi}D_{\text{sum}}$ | | $\dfrac{1}{7h^2} = \dfrac{0.143}{h^2}$ |
| | | $= 0.0318\left(\dfrac{D_4}{-D_2}\right)h^2$ | $= 1.038\,D_{\text{sum}}$ | | |

It can be seen from figure 3.20 that the values of the summit properties approach that of the continuous case as $k$ increases—a not unexpected result. It has also been shown for the hexagonal case $k = 6$ that the data storage is 13% less than for the rectangular case and in many instances the processing is quicker (see [52]). Also, if faults are present in the surface they are more easily detected. This result is shown by Whitehouse and Phillips [18].

The most important point to emerge is that the best results are those where the sampling pattern follows most closely the areal bandwidth pattern of the surface. For example, the hexagonal case $k = 6$ is most suitable for isotropic surfaces which have circular symmetry about the origin in frequency (and wavelength). In the case of anisotropic surfaces a suitably scaled sample interval in both directions with a rectangular grid for the tetragonal case will probably be best.

So simply applying a sampling procedure to a surface is not good enough. The sampling pattern should, whenever possible, image the surface properties. There is obviously another good case here for looking at the surface before trying to evaluate it.

### 3.2.8.3 *The effect of sampling interval and limiting results [15] on sample patterns*

It is important to investigate the variation of parameters with $h$ because it is due to the large number of possible differences in sampling interval that the scatter of measured values of parameters occurs between investigators.

The distributions of four-point summit height and curvature have been derived in terms of correlation coefficients between ordinates. These two correlation coefficients are $\rho_1$ for ordinates a distance $h$ apart, and $\rho_{\sqrt{3}}$ for ordinates a distance $\sqrt{3}h$ apart. If the surface is isotropic and the autocorrelation function is $\rho(x)$, then

$\rho_1 = \rho(h)$ and $\rho_{\sqrt{3}} = \rho(\sqrt{3}h)$. So $\rho_1$ and $\rho_3$ will vary as $h$ varies, depending on the shape of the autocorrelation function of the surface.

Results for the summit height have been obtained for the continuous surface, so it is possible to compare these results with those obtained earlier for the discrete results as the sampling interval $h$ converges to zero.

To do this it is necessary, as before, to make assumptions about the behaviour of the autocorrelation function $\rho(h)$ near the origin. It will be assumed as before that

$$\rho(h) = 1 + D_2 h^2/2! + D_4 h^4/4! + O(h^4) \tag{3.130}$$

where

$$D_2 < 0 \quad D_4 > 0 \quad \text{and} \quad \eta = -D_2(D_4)^{-1/2} < \sqrt{\tfrac{5}{6}}. \tag{3.131}$$

$D_2$ and $D_4$ are the second and fourth derivatives of the autocorrelation function at the origin.

Comparison will be made for the estimates of parameters measuring peak and summit properties of the surface. This will be done for the four cases of three-, four-, five- and seven-point estimates.

The first parameter that will be considered is the density of peaks or summits. These parameters are known for a continuous random Gaussian surface and were given for peaks as

$$D_{\text{peak}} = \frac{1}{2\pi}\left(\frac{D_4}{-D_2}\right)^{1/2} \tag{3.132}$$

by Rice [5] and for summits as

$$D_{\text{sum}} = \frac{1}{6\pi\sqrt{3}}\left(\frac{D_4}{-D_2}\right) \tag{3.133}$$

by Nayak [11].

The density of the peaks or summits is the number of peaks per unit length or summits per unit area, using the $(k+1)$-point definition of peak for $k = 2$ and summit for $k = 3$, $4$ and $6$. The expected density of peaks or summits is given by the product of $pr(Y_{k+1})$ and the density of ordinates, where $T_{k+1}$ is the event $(s_1 > 0,..., s_k > 0)$ and $s_1$ to $s_k$ are the differences between the central ordinate and the $k$ adjacent ordinates at a distance $h$.

The limiting behaviour *of $pr(Y_k+_1)$ as $h$* tends to zero, the density of ordinates and the limit of the expected density of peaks (or summits) are given in table 3.2. The limits are given in terms of the limiting results for a continuous surface given by [4] and [11]. It can be seen that the density of peaks (when $k = 2$) converges to the continuous limit. This is not the case for summits (when $k = 3$, $4$ and $6$). In these cases the density would be overestimated by 73%, 31% and 4% respectively (see figure 3.20).

The second parameter that will be considered is the average peak (or summit) height. The results are known for a continuous random Gaussian surface and were given for peaks as

$$E[Z \mid \text{continuous peak}] = \sqrt{\frac{\pi}{2}}\,\eta \tag{3.134}$$

by Rice [5] and Whitehouse and Phillips [4] and for summits as

$$E[Z \mid \text{continuous summit}] = 1.801\left(\frac{\pi}{2}\right)^{1/2}\eta \qquad (3.135)$$

by Nayak [11]. Hence the average summit height is 80% higher than the average peak height.

Again the expected height of peaks (when $k = 2$) converges to the continuous limit for peaks on a profile, *but this is not the case* for summits (when $k = 3$ and 6) as is seen in table 3.3. In these cases the expected

**Table 3.3** Expected summit (peak) height $E[Z_0 1/2 Y_k + 1]$.

| Model | $k$ | Expected summit (peak) height | |
|---|---|---|---|
| | | Limit as $h \to 0$ | Limit as $h \to \infty$ |
| Three points | 2 | $\left(\dfrac{\pi}{2}\right)^{1/2}\eta = 0.555\dfrac{4}{\sqrt{\pi}}\eta$ | 0.846 |
| Four points | 3 | $2\left(\dfrac{3}{\pi}\right)^{1/2}\eta = 1.559\left(\dfrac{\pi}{2}\right)^{1/2}\eta$ <br><br> $= 0.866\dfrac{4}{\sqrt{\pi}}\eta$ | 1.029 |
| Five points | 4 | $\dfrac{8\sqrt{2\pi}}{\pi + 2\sin^{-1}(\frac{1}{3}) + 4\sqrt{2}}\eta = 1.688\left(\dfrac{\pi}{2}\right)^{1/2}\eta$ <br><br> $= 0.938\dfrac{4}{\sqrt{\pi}}\eta$ | 1.163 |
| Seven points | 6 | $\dfrac{3\sqrt{\pi}[6\pi - 12\tan^{-1}(1\sqrt{2}) - \sqrt{2\pi}]}{\sqrt{\pi}[\pi + 6(\sqrt{3}-1)]}\eta = 1.779\left(\dfrac{\pi}{2}\right)^{1/2}\eta$ <br><br> $= 0.988\dfrac{4}{\sqrt{\pi}}\eta$ | 1.352 |

summit height is underestimated by 13% for the four-point case, by 6% for the five-point case and *by only* 1 % for the hexagonal case.

Because the conditional distribution of height given curvature is Gaussian with a mean which is a linear function of curvature, for all values of $k$, the expected summit curvature will converge in the same manner as the expected summit height.

To study the qualitative effect of the change of the sampling interval $h$ on the digital measurements of an isotropic surface it is necessary to specify a model for the autocorrelation function of the surface. Note that *any* autocorrelation could be used. For the model to fit in with observed autocorrelation functions of surfaces it would be desirable to have a negative exponential function with a multiplicative periodic function. Whitehouse and Phillips [4] 'smoothed' the exponential cosine function to give a function that was smooth at the origin. This alternative approach replaced the negative exponential function by another function which is smooth at the origin but behaves like the negative exponential function for large lag values. Lukyanov [21] has used such models extensively. Both of these are close to the autocorrelation functions of many typical practical surfaces. Note, however, that *the results are quite general* and that the specific models are only used to give an idea of actual scale. This model autocorrelation function is given by

$$\rho(h) = \mathrm{sech}\left(\tfrac{1}{2}\pi h A(\theta)\right)\ \cos\left(2\pi\theta h A(\theta)\right) \tag{3.136}$$

where

$$A(\theta) = \mathrm{sech}\left(2\pi\theta\right) + 2 \sum_{r=0}^{\infty} \frac{(-1)^r \theta}{2\pi\{\theta^2 + [1/4(2r+1)]^2\}\sinh(\pi(2r+1)/8\theta)} \tag{3.137}$$

where $h$ is the sampling interval. The values of $\theta$ used are 0 and $\tfrac{1}{2}$. For this autocorrelation function

$$D_2 = -\left(1/2\,\pi\right)^2 [1+(4\theta)^2][A_2(\theta)]^2 \tag{3.138}$$

and

$$D_4 = -\left(1/2\,\pi\right)^4 [5+6(4\theta)^2+(4\theta)^4][A_2(\theta)]^4. \tag{3.139}$$

A point needs to be made here. The rather complicated correlation function (3.136) is not the only function for which the analysis will work. It will work for any correlation function having $\rho_1, \rho_2$, etc, as values at the given spacings. This is wide enough to cover most, if not all, reasonable surfaces. The reason for this complex correlation function is simply to ensure that it has near perfect properties at the origin and elsewhere just to forestall any criticism. For practical cases the exponential, Gaussian, Lorentzian or whatever correlation function could have been picked to give an idea of the quantitative effects of sampling on the tribological properties of random surfaces.

What these results show is that by measuring the correlation values on a surface at spacings $h$, $2h$, $\sqrt{3}h$, etc, the tribological summit parameters can be found simply by inserting the values of $\rho_1$, $\rho_2$, etc, into the derived formulae. There is no need to measure the surface itself to get curvatures etc, providing that it is reasonably Gaussian (to within a skew of $\pm 1$, see Staufert [22]). The very tricky parameters can be simply calculated by knowing the correlation values.

For theoretical purposes a model of the surface can be devised as above and the values of $\rho_1, \rho_2$, etc, calculated. The tribological parameters can then be similarly evaluated.

The expected density of summits is given in figures 3.22 and 3.23 and the expected height of peaks or summits is given in figures 3.24 and 3.25 for the autocorrelation function for $\theta = 0$ and $\tfrac{1}{2}$.

The expected four-, five- and seven-point summits differ little as the sampling interval $h$ exceeds one correlation length. For smaller sampling intervals the four-point expected density of summits exceeds that for the five- and seven-point expectations.

### 3.2.8.4 Discussion

The technique of using discrete measurements has application in fields where it is expensive or time consuming to obtain large amounts of data. The reason for the analysis into sampling schemes is to try and see

**Figure 3.22** Expected density of summits (four and five points) of peaks (three points), $\theta = 0$.



**Figure 3.23** Expected density of summits (four and five points) of peaks (three points), $\theta = \frac{1}{2}$

whether taking measurements using a non-conventional sampling scheme would produce any advantages to outweigh the disadvantage of complexity. The advantages are less information to collect, easier analytical derivation of theoretical results and simpler numerical methods.

The sampling schemes that were considered all had the property that the information could be collected by sampling along parallel straight lines with a fixed sampling interval. (It might be necessary, however, to have a variable starting point, though this would follow a regular pattern.)

This ensured that if a measurement (ordinate) was chosen when using a particular scheme it would always have the same number of adjacent ordinates at a distance $h$ (the chosen sampling interval), provided the chosen ordinate is not on the boundary.

**Figure 3.24** Expected height of summits (four and five points) of peaks (three points), $\theta = 0$



**Figure 3.25** Expected height of summits (four and five points) of peaks (three points), $\theta = \frac{1}{2}$

From the point of view of simplicity of sampling mechanism the square grid ($k = 4$) in the tetragonal case is the best. In this case the spacing between the lines is constant and equal to the sampling interval $h$ along the line. Also the starting points for the sampling all lie along a straight line. However, other schemes do have advantages to offset their complexity.

The trigonal ($k = 3$) case has the advantage that measurements of slope can be taken in three directions as opposed to two for the tetragonal ($k = 4$) case. Though the theoretical results have been restricted to the consideration of isotropic surfaces it may still be of practical value to be able to check the assumption of isotropicity in more than two directions.

The trigonal ($k = 3$) case can be obtained by an alternative sampling method but this involves alternating the sampling interval from $h$ to $2h$. This alternative method is equivalent to rotating the grid through $\pi/6$. For $k = 6$ the sampling is very straightforward, as can be seen from figure 3.21.

From the point of view of collecting digital information, the trigonal ($k = 3$) case is preferable as 'less' information is collected. The density of ordinates is $(4/3\sqrt{3})h^2$ ($= 0.77/h2$) compared with $1/h2$ for the square grid ($k = 4$), so in the same area 23% fewer ordinates would be needed. The advantage of this would need to be weighed against the disadvantages.

Another advantage of the trigonal ($k = 3$) case is that fewer ordinates are used when defining the properties of the extremities. To check the definition of a four-point summit only three conditions have to be obeyed, as opposed to four conditions for the five-point summit and six for the seven point. It should also be noted that some properties of the discrete defined random variables, such as the limiting values of the ($k + 1$)-point summit (or peak) height as the sampling interval tends to infinity, are simply a function of the numerical definition and are independent of the surface being measured. The problem is that the surface has to be measured to get any values for the properties.

Any discrete measurement of a surface must lose information compared with a complete 'map' of the surface. This is inevitable! However, ideally, any discrete measurement should produce results that converge to the results for the continuous surface as the sampling interval $h$ tends to zero.

For sampling along a straight line ($k = 2$) it has been seen that the discrete results do converge to those for the continuous profile. They do not, however, converge to the results of the two-dimensional surface. For example, $D^2_{peak} = 0.83\, D_{sum}$, so that assuming independent measurements at right angles would produce a limit that is 17% too small.

For two-dimensional measurements when sampling with $k = 3$, 4 or 6, the limiting results for expected summit density and expected summit height do not converge to the results for the continuous surface. In the case of expected summit density the limit is 73% too large for $k = 3$, 31% too large for $k = 4$ and 4% for $k = 6$. Again for expected summit height, the case $k = 3$ is worse than for $k = 4$ and $k = 6$ but the differences are not so large. This suggests that some surface parameters may be estimated by discrete methods fairly well but others may not. For the case of average profile slope three sampling schemes agree (for $k = 2$, 3 and 4) but this is, of course, an essentially one-dimensional parameter.

In order to consider the merits of sampling schemes it is necessary to study their theoretical properties. By doing this it is possible to obtain new insights into the general problem, but only by using models which lead to tractable mathematics. The three sampling schemes with $k = 2$, 3 and 4 considered here have been chosen because they have a common property that enables them to be investigated using the analytical results of theoretical statistics. Using the trigonal ($k = 3$) symmetry case leads to a simpler mathematical model than for the tetragonal ($k = 4$) symmetry case, as this reduces the dimension by one. However, taken as a whole it may be that a hexagonal sampling plan where $k = 6$ offers the maximum benefit in terms of the three criteria mentioned above. One message which has emerged is that the conventional grid pattern method of sampling is not necessarily the best. The implications of this in general pattern recognition and image analysis scanning systems could be significant.

The result given here has been concerned primarily with the effect of sampling patterns on the values of parameters obtained from measured surfaces. It has not therefore been aimed at investigating the actual nature of surfaces in general. Well-behaved correlation functions have been assumed and certain specific examples have been used to give the researcher an idea of the value of parameter changes that might be expected to occur on typical measured surfaces. This has been justified by the fact that to some extent all instruments used for obtaining the data have a short-wavelength filter incorporated, whether it be a stylus or a light spot, which tends to force the correlation at the origin to be smooth. However, there can be no denying that the very nature of random manufacture encourages the presence of exponential and other misbehaved characteristics in the correlation function. The effect of sampling patterns on such fractal (multiscale) surfaces will be the subject of further comment in chapter 7.

In the above sections much emphasis has been placed on statistical parameters, distributions and functions. Some effort has to be expended to make sure that the values obtained are realistic. The next section therefore will revert back to a simple discussion of the processing of these parameters. However, because of its central importance in random process theory, the role and evaluation of Fourier transforms will be described here.

### 3.3 Fourier transform and the fast Fourier transform

#### 3.3.1 General properties of the Fourier transform

The Fourier transform may be of use whenever frequency, either in reciprocal time or wavenumber, is of importance. This is often the case in engineering, in a temporal sense in manufacturing systems and spatially in metrology. Examples of its use in correlation and filtering have already been mentioned as well as its use in random process theory.

The general form of the equation relating a time function $f(t)$ to a frequency function $F(\omega)$ is

$$F(\omega) = \int_{-\infty}^{\infty} f(t) \exp(-j\omega t)dt. \tag{3.140}$$

Alternatively this can be written with the factor $1/2\pi$ in front of the integral. The corresponding equation connecting $f(t)$ with $F(\omega)$ is

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} f(\omega) \exp(-j\omega t)d\omega. \tag{3.141}$$

As long as there is a $2\pi$ relationship between these two equations it does not matter in which domain it occurs. Some authors use a factor of $1/\sqrt{2\pi}$ outside both integrals to provide some symmetry. The first equation is referred to as the Fourier integral. Before embarking upon the fast Fourier transform algorithm there will be a short listing of some of the essential properties of the Fourier integral as follows:

1. If $f(t)$ is real, as is usually the case in metrology, then the real and imaginary parts of the $F(\omega)$, namely $R(\omega)$ and $X(\omega)$, are given by

$$R(\omega) = \int_{-\infty}^{\infty} f(t) \cos\omega t\ dt \tag{3.142}$$

$$X(\omega) = - \int_{-\infty}^{\infty} f(t) \sin\omega t\ dt \tag{3.143}$$

that is $R(\omega)$ is an even function and $X(\omega)$ is odd. Thus $R(-\omega) = R(\omega)$ and $X(-\omega) = -X(\omega)$.

2. If $f(t)$ is even, that is $f(-t) = f(t)$,

$$R(\omega) = 2 \int_{-\infty}^{\infty} f(t) \cos\omega t\ dt \qquad\qquad X(\omega) = 0$$

which is particularly relevant when $f(t)$ takes the form of an autocorrelation function.

3. If $f(t)$ is odd, that is $f(-t) = -f(t)$,

$$R(\omega) = 0 \qquad\qquad X(\omega) = -2 \int_{-\infty}^{\infty} f(t) \sin\omega t\ dt. \tag{3.144}$$

4. Linearity: if $F_1(\omega)$ and $F_2(\omega)$ are the Fourier integrals of time functions $f_1(t)$ and $f_2(t)$ then

$$f_1(t) + f_2(t) \leftrightarrow F_1(\omega) + F_2(\omega) \tag{3.145}$$

where the symbol $\leftrightarrow$ indicates Fourier transformation.

5. Scaling:

$$f(Kt) = \frac{1}{K} F\left(\frac{\omega}{K}\right). \tag{3.146}$$

6. Shifting in time by $t_0$:

$$f(t - t_0) \leftrightarrow \exp(-\mathrm{j}t_0\omega)\mathrm{F}(\omega)$$

a property which is useful when applying the theory of the phase-corrected (linear phase) filter. A similar relationship occurs in frequency. Thus

$$\mathrm{F}(\omega - \omega_0) \leftrightarrow \exp(\mathrm{j}\omega_0 t)f(t). \tag{3.147}$$

7. Differentiation:
   (a) Time:

$$\frac{\mathrm{d}^n f}{\mathrm{d}t^n} \leftrightarrow (\mathrm{j}\omega)^n F(\omega) \tag{3.148}$$

a theorem which is useful for evaluating the frequency characteristic of a complicated impulse response. A similar theorem exists which enables the reverse to be done in frequency differentiation.
   (b) Frequency:

$$\frac{\mathrm{d}^n F}{\mathrm{d}\omega^n} \leftrightarrow (-\mathrm{j}t)^n f(t). \tag{3.149}$$

8. Convolution: this has already been used in filtering techniques. In condensed form, if $f_1(t) \leftrightarrow F_1(\omega), f_2(t) \leftrightarrow F_2(\omega)$, then

$$\int_{-\infty}^{\infty} f_1(\tau) f_2(t - \tau) \mathrm{d}\tau \leftrightarrow F_1(\omega) \times F_2(\omega)$$

written $\tag{3.150}$

$$f_1(\tau) * f_2(t) \leftrightarrow F_1(\omega) \times F_2(\omega)$$

where the symbol * denotes convolution. These relationships will be used extensively in chapter 4. The discrete Fourier transform (DFT) may be written

$$F(k) = \frac{1}{N} \sum_{n=0}^{N-1} f(n) \exp\left(\frac{-2\pi\mathrm{j}}{N} nk\right) \quad n = 0, \dots, N-1. \tag{3.151}$$

### 3.3.2 Fast Fourier transform routine

The problem with the DFT as it stands is that it takes quite a long time to compute. This situation has been changed by the rediscovery of the fast Fourier transform (FFT) algorithm. This is not a special transform, it is a generic term enfolding a whole set of algorithms for computing the DFT. Some algorithms are matched to specific types of data etc. Generally, using the FFT, a reduction in the number of operations by a factor of $\log_2 N/N$ is possible, where 'operation' is usually taken to mean a complex multiplication and an addition. Thus the larger the number of data points $N$ the greater the advantage. The FFT is also efficient in terms of storage because intermediate results in the calculation are stored in the original data places so that no extra storage is required for those extras beyond that of the data. Yet another benefit derives directly from the computation reduction: less computation means less numerical errors. Thus the FFT has three remarkable advantages: it is faster, more accurate and requires minimal storage. The disadvantage is that there have to be as many spectral points evaluated as there are data points, which may be more than is needed.

#### 3.3.2.1 Fast Fourier transform, analytic form

This is a fast algorithm of efficiently computing $F(k)$

$$F[k] = \sum_{n=0}^{N-1} f(n) \exp\left(\frac{-\mathrm{j}2\pi kn}{N}\right). \tag{3.152}$$

For $N = 2^r$ there are many variants.

Let $k$ and $n$ be represented in binary form as

$$
\begin{aligned}
k \quad &= k_{r-1}2^{r-1} + k_{r-2}2^{r-2} + \ldots + k_0 \\
n \quad &= n_{r-1}2^{r-1} + \ldots + n_0 \\
W \quad &= \exp(\mathrm{j}2\pi/N)
\end{aligned}
\tag{3.153}
$$

Then

$$
\begin{aligned}
F(k_{r-1}, k_{r-2}, \ldots, k_0) = \sum_{n_0=0}^{1} \sum_{n_1=0}^{1} \cdots \sum_{n_{r-1}=0}^{1} f(n_{r-1}, n_{r-2}, \ldots, n_0) \\
\times \; W^{k_0 n_{r-1} 2^{r-1}} \quad \times \; W^{(k_1 2^1 + k_0) n_{r-2} 2^{r-2}} \\
\times \; \ldots \; \times \; W^{(k_{r-1} 2^{r-1} + k_{r-2} 2^{r-2} + \ldots + k_0) n_0}
\end{aligned}
\tag{3.154}
$$

Performing each of the summations separately and labelling the intermediate results the following is obtained:

$$
\begin{aligned}
F_0[n_{r-1}, n_{r-2}, \ldots, n_0] &= f[n_{r-1}, n_{r-2}, \ldots, n_0] \\
F_1[k_0 n_{r-2}, \ldots, n_0] &= \sum_{n_{r-1}=0}^{1} F_0[n_{r-1}, n_{r-2}, \ldots, n_0] \; W^{k_0 n_{r-1} 2^{r-1}} \\
F_2[k_0, k_1, \ldots, n_0] &= \sum_{n_{r-2}=0}^{1} F_1[k_0, n_{r-2}, \ldots, n_0] \; W^{(k_1 2^1 + k_0) n_{r-2} 2^{r-2}} \\
F_r[k_0, k_1, \ldots, k_{r-1}] &= \sum F_{r-1}[k_0, k_1, \ldots, n_0] \; W^{(k_{r-1} 2^{r-1} + k_{r-2} 2^{r-2} + \ldots + k_0) n_0} \\
F[k_{r-1}, k_{r-2}, \ldots, k_0] &= F_r[k_0, k_1, \ldots, k_{r-1}]
\end{aligned}
\tag{3.155}
$$

This operation is performed in $N \log N$ complex multiplications through the indirect method above. If $F[k]$ *is* evaluated directly as

$$\sum_{n=0}^{N-1} f(n) \exp\left(-j\frac{2\pi n}{N}\right) \tag{3.156}$$

it takes the order of $N^2$ complex multiplications, so the FFT is obviously faster. If $N$ is large, as it usually is for the data set in any surface metrology, the gain in time is very great. The fact that there has to be as much store for the spectrum as there is for the data is no longer a disadvantage, although at one time when storage was at a premium in computers this was.

The only better way of achieving further gains is to use one or other of the clipped transforms, such as the Walsh transform [23]. As far as speed is concerned, some instrumental techniques can dispense altogether with digital methods and use coherent optical methods, as will be seen in the section on Fourier diffraction techniques (section 3.10). In this case the transform is produced at the speed of light. A further advantage is that because the light is shone onto the surface, no digitization of the surface is required so that optical methods are very valid in surface metrology. The possibility of an optical method that is practical for evaluating the Wigner function is also feasible.

In what follows, some conventional digital ways of filtering will be given with particular emphasis on the types of filter used in surface metrology.

Some of the comparative properties of the Fourier transform, ambiguity function and Wigner distribution are given in table 3.4.

As previously mentioned the transform pair involves a pair of complementary transforms, one with positive exponents and one with negative exponents, the pair taken together having a normalization factor of $1/2\pi$ or $1/N$ for the DFT. This choice is arbitrary. The equation with positive exponential is usually called the direct transform and the negative exponential the inverse transform (IDFT):

$$f(t) = \frac{1}{2\pi} \int_{-\infty}^{\infty} F(\omega) \exp(j\omega t) \, d\omega. \tag{3.157}$$

Sometimes both are referred to as the DFT and sometimes the inverse discrete Fourier transform (IDFT) is called the Fourier series, as in most textbooks.

To see how the discrete form of the FFT works consider the discrete form of $f(t)$: $f(n)$

$$f(n) = \sum_{k=0}^{N-1} F(k) \exp\left(\frac{2\pi}{N} jn\right) \qquad n = 0, \ldots, N = 1. \tag{3.158}$$

Equation (3.144) can be written

$$f(n) = \sum_{k=0}^{N-1} F(k) \, W^{nk} \qquad \text{where } W = \exp\left(\frac{2\pi j}{N}\right) \tag{3.159}$$

and $W^{nk} = exp[(2\pi j/N)nk]$ which has a period of $N$.

The basic method used in the FFT algorithm is to make use of the symmetry of the exponential function. $W^{nk}$ is a sinusoid of period $N$ and it also displays odd symmetry about $k = N/2$ for fixed $n$ etc. Thus $W^{nk} = -W^{n(k+N/2)}$.

This cyclic property is of fundamental importance. It can be shown that the efficiency of the FFT routine depends on the data sequence length. In fact it is optimized for the case where $N = 3^m$, where $m$ is an

**Table 3.4** Comparison of the properties of the Fourier transform, the ambiguity function $A_f(\overline{\omega}, x)$ and the Wigner distribution function $W_f(x, \omega)$.

| | Complex valued | Complex valued | Complex valued | Real valued |
|---|---|---|---|---|
| | $f(x)$ | $F(\omega)$ | $A_f(\omega,\chi)$ | $W_f(x, \omega)$ |
| Spatial shifting | $f(x - x_0)$ | $F(\omega)\exp(-j\omega x_0)$ | $A_f(\overline{\omega}, \chi)\exp(-j\overline{\omega}x_0)$ | $W_f(x - x_0, \omega)$ |
| Frequency shifting | $f(x)\exp(j\omega_0 x)$ | $F(\omega - \omega_0)$ | $A_f(\overline{\omega}, \chi)\exp(-j\overline{\omega}_0\chi)$ | $W_f(x, \omega - \omega_0)$ |
| Spatial limited | $[x_a, x_b]$ <br> for $x$ | $[-\infty, \infty]$ <br> for $\omega$ | $[-(x_b - x_a), x_b - x_a]$ <br> for $\chi$ | $[x_a, x_b]$ <br> for $x$ |
| Frequency limited | $[-\infty, \infty]$ <br> for $x$ | $[\omega_a, \omega_b]$ <br> for $\omega$ | $[-(\omega_b - \omega_a)\omega_b - \omega_a]$ <br> for $\overline{\omega}$ | $[\omega_a, \omega_b]$ <br> for $\omega$ |
| Convolution | $\displaystyle\int_{-\infty}^{\infty} f(\chi)h(x - \chi)\mathrm{d}\chi$ | $F(\omega)H(\omega)$ | $\displaystyle\int_{-\infty}^{\infty} A_f(\overline{\omega}, x)A_h(\overline{\omega}, \chi - x)\mathrm{d}x$ | $\displaystyle\int_{-\infty}^{\infty} W_f(\chi, \omega)W_h(x - \chi, \omega)\mathrm{d}\chi$ |
| Modulation | $f(x)m(x)$ | $\dfrac{1}{2\pi}\displaystyle\int_{-\infty}^{\infty} F(\overline{\omega})M(\omega - \overline{\omega})\mathrm{d}\overline{\omega}$ | $\dfrac{1}{2\pi}\displaystyle\int_{-\infty}^{\infty} A_f(\omega, \chi)A_m(\overline{\omega} - \omega, \chi)\mathrm{d}\omega$ | $\dfrac{1}{2\pi}\displaystyle\int_{-\infty}^{\infty} W_f(x, \overline{\omega})W_m(x, \omega - \overline{\omega})\mathrm{d}\overline{\omega}$ |
| Typical stationary | $A \exp\left(j\omega_0 x\right)$ | $A2\pi\delta\left(\omega - \omega_0\right)$ | $\lvert A \rvert^2 \; \exp\left(-j\overline{\omega}_0 x\right)2\pi\delta\left(\overline{\omega}\right)$ | $\lvert A \rvert^2 \; 2\pi\delta\left(\omega - \omega_0\right)$ |
| Typical non-stationary | $A \exp\left(j\dfrac{\alpha}{2} x^2\right)$ | $A\dfrac{2\pi}{\alpha}\exp\left(j\dfrac{\pi}{4}\right)\exp\left(-\left(\dfrac{j\omega^2}{2\alpha}\right)\right)$ | $\lvert A \rvert^2 \; 2\pi\delta\left(\overline{\omega} - \alpha\chi\right)$ | $\lvert A \rvert^2 \; 2\pi\delta\left(\omega - \alpha x\right)$ |

integer. However, for digital computer operation it is more convenient to use $N = 2^m$ or $4^m$ without significant loss in efficiency. Obviously base 2 is more suitable for computational work than base 3.

To see how it is possible to derive a base 2 FFT without recourse to the general theory, consider equation (3.158). If $N$ is divisible by 2 the summation can be split into two smaller sums involving odd and even indices. Thus

$$f(n) = \sum_{k=0}^{N/2-1} F(2k)\, W^{2nk} + \sum_{k=0}^{N/2-1} F(2k+L)W^{(2k+1)}$$

(3.160)

so that the first sum uses values of $F(0)$ to $F(N-2)$ and the second sum uses values of $F(l)$ to $F(N-1)$. Hence $f(n)$ may be written in terms of odd and even indices. Thus

$$f(n) = E(n) + W^n O(n) \qquad \text{where } W^n = \exp\!\left(\frac{2\pi \mathrm{j} n}{N}\right)$$

(3.161)

and

$$E(n) = \sum_{k=0}^{N/2-1} F(2k)\, W^{2nk}$$

$$O(n) = \sum_{k=0}^{N/2-1} F(2k+1)\, W^{2nk}.$$

(3.162)

$W^n$ is a constant for constant $n$.

The two series in equation (3.161) will be recognized as discrete transforms of the odd and even-indexed terms of the original series in equation (3.158). Note that the 2 in the index of $W$ is appropriate because these two smaller series only have a length of $N/2$ and not $N$. It is this factor that is decisive in cutting the number of operations. To see this consider equation (3.158). To evaluate one $f(n)$ requires $N$ operations and to evaluate $N$ values of $f(n)$ therefore requires $N^2$ operations where operation refers to a complex multiplication and addition. However, the situation has changed in equation (3.161) because there are only $N/2$ terms in the even series and $N/2$ terms in the odd series. To evaluate all the even components of all the $F(n)$ therefore requires $(N/2)^2$ operations and to work out all the odd components of the $f(n)$ also requires $(N/2)^2$. The two, together with the $N$ multiplications of $Wn$ and $N$ additions of the $E(n)$ and $O(n)$, give $(N^2/2 + N)$ operations. For $N$ large this tends to $N^2/2$, a gain over the direct method of 50% in computation. Similarly, if $N/2$ is also divisible by 2 a gain may be made by breaking down the two series. This reduction process can be continued $M$ times until finally no further decomposition can take place. The overall number of operations carried out in the entire operation is $N \log_2 N$ as distinct from $N^2$ in the direct approach, the ratio being $\log_2 N/N$. Thus for $N = 256$ the computational time is only 3% of that used in the direct method.

This iterative process is shown in figure 3.26, which illustrates how an eight-point transform can be produced by three iterations. Notice that the input points are not arranged in their natural order. The reason for this becomes clear from the figure; at each stage a pair of points depends only on the pair previously residing at the same locations. Thus the results of each iterative step can overwrite the operands without subsequently affecting future calculations. The computation can proceed 'in place'; no extra storage is needed for intermediate results. The order into which the original sequence has to be arranged is successively determined by separating out the odd from the even indices from the progressively shorter sequences. It can be shown that each element of the sequence should be placed at an address given by reversing the order of the bits in the binary representation of its original address.

Because of the increasing use of the FFT routine some consideration will be given here actually to implementing the routine on a small computer, and one or two special tricks associated with metrology will be indicated. This will be by no means an exhaustive account but it should serve to illustrate some of the features of the technique. It will be followed by some general applications to metrology problems on surfaces.

### 3.3.3 A practical realization of the fast Fourier transform

The term practical realization will be taken here to mean a form of the algorithm suitable for use with a general purpose digital computer which can be efficiently applied to typical transform problems. It is widely recognized that the optimum computational efficiency for a specific problem is obtained only by tailoring the program to its requirements. It follows from this that a general purpose routine will need to compromise between the sometimes conflicting needs of different problems. An obvious example is that if large amounts of data are to be handled, computer storage may be at a premium and it may be preferable to make the FFT



**Figure 3.26** (a) Eight-point fast Fourier transform decimation-in-frequency branch transmissions in powers of $W_8$, output shuffled; (b) eight-point fast Fourier transform decimation-in-tune branch transmissions in powers of $W_8$, input shuffled.

routine small at the expense of reducing its speed. For a 'real-time' problem, the opposite approach of using more storage to increase the computational speed would probably be necessary.

The first compromise is to use a base 2 FFT algorithm. The effect of this is to get good computational efficiency and very little 'housekeeping' so that the actual FFT program may be both fast and physically small. To offset this the data must be equal to $2^m$ points, where *m is a* positive integer.

Having decided on the base 2 FFT, the actual form of the derivation used is of little consequence. For the further discussion here the approach described earlier will be used, that is that based on equations (3.160) and (3.161).

Some general considerations will now be given to the practical realization. In the first place the general situation, which involves complex data, will be considered. This will be followed by some details of how these can be recovered in the face of the less restrictive data often found in surface metrology.

### 3.3.4  General considerations

Remembering that all sequences in the discrete Fourier transform are complex, and that, therefore, each data point will occupy more than one store address, three distinct sources of potential inefficiency can be identified. These are the complex multiply and add operations at the heart of each calculation, the addressing of the data points at each stage of the iteration and the generation of the complex exponents.

There is little than can be done about the actual operation of multiplication without the use of special hardware, particularly if floating point operation is used. In fixed point operations, where scaling between operations is used, some optimization is possible since the nature of the multiplier (the exponential term) is known in advance. Prechecking each operand for zero and unity value can save time by eliminating unnecessary multiplications, but in general the overall time of the complete evaluation may not be reduced since extra time is required before every multiplication in order to perform this test.

Each iteration of the FFT requires at least one access to each point in the data array for reading and one for writing. The order in which the points are accessed will change for each iteration, since the point pairs move their relative positions, so that the data array must be treated as random access rather than cyclic. The accessing of array elements can introduce significant overheads in terms of the speed of operation and so should be kept to a minimum. While at least two accesses are needed to each pair of locations, these always occur as a group before the locations are addressed. Thus the addresses need only be calculated once per iteration, being stored for use during accessing, without significantly increasing the total store requirements. If the exponentials are stored in a look-up table, the amount of addressing to that can be minimized by ordering the calculations so that those using a particular exponent occur consecutively during each iteration.

The generation of complex exponentials is almost universally performed using de Moivre's theorem:

$$\exp(j\theta) = \cos\theta + j\sin\theta.$$

(3.163)

By this method it is only necessary to calculate either the sine or cosine, since they are simply related and the calculation time is relatively small. However, compared with the other operations within the FFT the time to calculate $\sin\theta$ is large. Since the calculation would also require quite a lot of programming, it is more usual to use some form of look-up table for the values. Only the first two quadrants need to be stored. The form which a look-up table takes depends largely on the compromise between speed and storage. The table may be made progressively shorter by increasing the amount of calculation required to access a particular value. The fastest access is by simply storing the $N/2$ complex exponent values. The store may be reduced with little overhead by having two separate arrays of $N/2$ points for sine and cosine, which may, of course, be overlapped by one quadrant, and combining them into a complex number at access time. An interesting method of retaining fairly fast access with only small storage is to use the approximate relationship that, for small *B,*

$$\sin(A + B) = \sin A + \cos A \sin B$$

(3.164)

Here only a short table containing widely spaced $A$ values and another short table of $B$ values are needed.

Also, during the second iteration $W$ may take only the values $+1, -1, j, -j$. By programming the first two iterations separately outside the iterative loop, no multiplications are needed during them.

Because of the ready availability of FFT subroutines, they will not be described here.

### 3.3.4.1 The Fourier series of real data

By taking into account symmetries and identities readily identifiable it is possible to take great advantage of the FFT routine. Furthermore, in the case of evaluating the autocorrelation function from the former spectral density and vice versa, additional benefits can be achieved because the data is not only real but also even (i.e. symmetrical about the origin). Under these circumstances use can be made of the cosine transform rather than the full Fourier transform.

Instead of grouping the data in such a way as to appear complex it is grouped to appear as a conjugate, even series. Thus $z(n)$ may be described as

$$
\begin{aligned}
f_1(n) &= z(2n) \\
f_2(n) &= z(2n+1) - z(2n-1) \\
f(n) &= f_1(n) + j f_2(n) \qquad \text{for } n = 0, 1, \ldots, N-1.
\end{aligned}
\tag{3.165}
$$

Since $z(n)$ is real and even and $f(n)$ is conjugate even, the sequence $f(n)$ need only be calculated for $n = 0, \ldots, N/2$; only the $z(n)$ values of $n = 0, 1, \ldots, N$ are needed for this. Therefore only $N/2 + 1$ complex storage locations are needed for processing $2N$ real, even data.

Using this technique only $N + 1$ real or $N/2 + 1$ complex numbers need to be stored, giving a saving of four times in store over the $2N$ complex data series storage requirement.

Similar statements can be made about the sine transform. It is obvious that the FFT routine has many uses. Some will be outlined here.

### 3.3.5 Applications of Fourier transforms with particular reference to the FFT

Main areas of use are:

(1) digital filtering
(2) power spectral analysis
(3) correlation
(4) other convolutions
(5) interpolation
(6) other analysis, for example Cepstrum
(7) roundness and other shapes—least-squares references.

### 3.3.5.1 Use of Fourier transform for non-recursive filtering

This is basically a convolution filter using the impulse response of the filter $\omega$ as a window for the data. Thus the output $g$ of a filter with input $f$ is $g_i = w * f_i$. Using the transform it is preferable to work in the frequency domain; thus

$$
G_i = W \times F_1
\tag{3.166}
$$

since convolution becomes multiplication in the transform domain. So filtering becomes:

1. transform data
2. multiply transformed data by frequency weighting function.
3 invert transformed and modified data.

Note that it sometimes beneficial to change the spectrum obtained from 1 above to an analytical funtion. This involves making $\omega < 0 = 0$. This operation involves taking the Hilbert transform (See 3.5.3.1.3).

Whether this method is quicker to compute will depend on the lengths of the data sequence and weighting function. Recently this technique has become progressively more attractive as the data lengths have increased. Because of the reciprocal nature of the transform domains, an added advantage of the transform technique may occur: a long, spatial weighting function will produce a short-frequency weighting function so that even less operations are required.

In the frequency domain, the amplitude and phase characteristics of the filter are separated into the real and imaginary parts of the weighting function. Thus it is easy to construct unusual filter characteristics. In particular, the phase-corrected filter is constructed by having real terms describing the required attenuation and by having zeros for the imaginary terms so that no phase modification takes place. Also, the defined amplitude characteristic of the phase-corrected filter is very short in the frequency domain. In the Whitehouse phase-corrected filter it is even simpler because the frequency characteristic is mostly either zero or one and is efficient to compute.

### 3.3.5.2  *Power spectral analysis*

The measurement of power spectral density is the most obvious application of the Fourier transform. Put simply, a data sequence from the surface is Fourier transformed and a power spectrum obtained by squaring each term independently. This gives the periodogram, which is an estimate of the power spectrum.

The periodogram is a rather unstable estimate and various smoothing methods may be applied.

### 3.3.5.3  *Correlation*

The calculation of auto and cross-correlation functions may be undertaken by transform methods since the correlation function and power spectrum form a transform pair.

Autocorrelation is performed by transforming the surface data, squaring term by term and then retransforming back to the spatial domain. For cross-correlation both data sequences must be transformed and a cross-power spectrum obtained so that more computational effort is needed. However, for real data, the system of performing two transforms at once using the FFT can reduce the total time taken.

The transform technique of obtaining correlation functions always gives all possible lag values up to the total data points, whereas, particularly with autocorrelation, often only lags up to a small proportion of the data sequence length are needed. Under these circumstances it may be computationally quicker to use traditional lag calculation for the autocorrelation function, especially if the total data length is fairly small.

The discrete Fourier transform operates by imposing a periodicity of wavelength equal to the data sequence length in that the Fourier integral is defined for an infinite waveform made up of repeats of the actual data. In forming the correlation function it is these constructed infinite sequences that are, in effect, lagged so that points lagged beyond the end of the sequence do not 'drop off' but match up with points on the start of the first repeat. This looks like a circular correlation where the data is placed in a complete, closed circle and the lag produced by rotating this.

This effect can be useful in roundness measurement, for instance when comparing two randomly oriented profiles. However, it can lead to erroneous values in linear situations. To overcome this a sequence of zero-value points of length up to at least the maximum desired lag should be appended to the data sequence. The action of these is made clear in figure 3.27.

### 3.3.5.4 Other convolutions

The main use of convolution is covered under filter and correlation. The heading here just stresses that any convolution in space may be handled as multiplication in frequency. Beware of the circular convolution of the DFT 'described' above.



**Figure 3.27** Circular correlation.

### 3.3.5.5 Interpolation

It is easily possible to produce an interpolation of a data sequence in which the harmonic quality of the original data is retained in the expanded sequence.

The discrete transform pair always has the same number of points in the sequence in either domain. Thus a short data sequence may be transferred and the frequency domain sequence expanded by adding zero-value elements at the high-frequency position (the centre, since the transform is symmetrical). Inverse transformation will produce a longer data sequence with harmonic interpolation, since no non-zero frequency components have been added.

### 3.3.5.6 Other analysis

Various other signal processing techniques use Fourier transforms. Discrete transforms can be used to produce a form of discrete Laplace transform analysis.

Another possibility is Cepstrum analysis. The Cepstrum is defined as the transform of the logarithm of the power spectrum. As such it highlights periodicities of the power spectrum, which are, of course, harmonic sequences in the original data. Should the power spectrum of the data consist of two or more harmonic sequences multiplied together in the Cepstrum (log spectrum) they will be added. The power spectrum of this will clearly separate these different added harmonic sequences.

### 3.3.5.7  Roundness analysis

The Fourier transform will yield roundness reference lines from full circles since the first term of the Fourier series is the amplitude of the least-squares limaçon. Similarly approximation to reference ellipses, and so on, can be obtained. Because of the computation involved, it is unlikely that the DFT would be used just to determine the reference circle.

Direct evaluation of the harmonic series (since the circle is periodic) also seems useful for assessing the performance of bearings.

Perhaps more could be done using the DFT to produce circular cross-correlations, for instance to allow the comparison of cam-forms with stored masters without the need for accurate orientation. Specific applications of Fourier transforms will be dealt with in many places within this book.

## 3.4  Statistical parameters in digital form

### 3.4.1  Amplitude probability density function

This contains the height information in the profile. It has been used extensively in one form or another in surface analysis. Various terms are used to describe the function, its derivatives and its integral.

In figure 3.28 the value for the amplitude probability density function (APDF) at a given height is called $p(z)$ and is the number of profile ordinates having a numerical value between $z$ and $z + \delta z$ divided by the total number of profile ordinates. There are two basic ways of evaluating the APDF digitally. The first and most obvious way is to select a height interval, say between $z$ and $z + \delta z$, and to scan through the whole profile data, counting how many ordinates lie within this height range. The height interval is changed and the operation is repeated. This is carried on until all the vertical range of the profile has been covered. In this method the height interval is selected before the counting is started. The other method involves examining every profile ordinate in sequence, finding to which height interval it belongs and registering a count in the store position corresponding to this level. The advantage of the last method is that the whole of the amplitude distribution (APDF) is found after just one traverse and the only store requirement is that of the number of height intervals to cover the height range; the profile itself need not be stored. Measurement of the central moments, skew and kurtosis can be made from the APDF directly.



**Figure 3.28** Statistical height distributions.

For instance, the RMS value $R_q$ is given by

$$R_q^2 = \sum_{i=1}^{n} (i\Delta z)^2 p_i$$

(3.167)

where $\Delta z$ is the quantization interval of the pockets and where the mean

$$\bar{z} = \sum_{i=1}^{n} (i\Delta z)p_i$$

(3.168)

and there are $N$ levels of $p_i$, in the APDF. Notice that the $R_q$ value is subject to an error of $\Delta z/\sqrt{12}$ which is a penalty for being in discrete rather than continuous form. These moments can be measured more accurately without recourse to increasing the store as running summations in terms of each other. It is not necessary to determine $p_i$ as such, the natural quantization interval of the incoming data $q$ being the only limitation in accuracy. Thus

$$R_q = \left| \frac{1}{m} \sum_{i=1}^{m} z_i^2 - \left( \frac{1}{m} \sum_{i=1}^{m} z_i \right)^2 \right|^{1/2}$$

(3.169)

where the number of ordinates is $m$ and the $z$ values are measured from an arbitrary datum.

## 3.4.2 Statistical moments of the APDF

$$\text{skew} = \left[ \frac{1}{m} \sum_{i=1}^{m} z_i^3 - \frac{3}{m^2} \sum z_i \sum z_i^2 + \frac{2}{m^3} \left( \sum z_i \right)^3 \right] \bigg/ \left[ \frac{1}{m} \sum z_i^2 - \left( \frac{1}{m} \sum z_i \right)^2 \right]^{3/2}$$

(3.170)

$$\text{kurtosis} = \left[ \frac{1}{m} \sum z_i^4 - \frac{4}{m^2} \sum z_i \sum z_i^3 + \frac{6}{m^3} \left( \sum z_i \right)^2 \sum z_i^2 - \frac{3}{m^4} \left( \sum z_i \right)^4 - 3 \right] \bigg/ \left[ \frac{1}{m} \sum z_i^2 - \left( \frac{1}{m} \sum z_i \right)^2 \right]^2.$$

(3.171)

In other words the central moments can all be expressed in terms of the moments about an arbitrary level. It is the central moments which contain the information in surface metrology.

All the important moments can be found directly from the APDF, even the $R_a$ value. It is given by

$$R_a = 2 \left( \bar{z} \sum_{\substack{i=1 \\ z_i < \bar{z}}}^{N} p_i - \sum_{\substack{i=1 \\ z_i < \bar{z}}} \Delta z_i . p_i \right)$$

(3.172)

so that, providing the APDF is measured as the input signal enters, all the moments, peak values (subject to a $\Delta z$ limit on accuracy) and $R_a$ can be measured without actually storing the data!

The distribution function is merely the cumulative sum of $p_i$, up to a given level. Thus

$$P_j = \sum_{i=1}^{j} p_i$$

(3.173)

The bearing ratio, material ratio or Abbott-Firestone curve often used in surface metrology are all $1-P_j$.

One of the problems in assessing the profile parameters without evaluating $p_i$, is that the number of ordinates can be so great that the summations may overload the computational wordlength.

These techniques are not restricted to the measurement of a surface profile; any waveform can be analysed in the same way whether it is derived from the profile itself, the slope of the profile or whatever. There is a problem, however, in isolating the different components of, say, texture. This is because, although these quick and economic techniques for statistical moment evaluation do not require knowledge of the mean level, they do assume that the signal follows the general direction of the surface. In other words, it is not possible to eliminate the need for a reference line but it is possible to relax the stringent positioning of it relative to the profile. Fortunately this is not a problem in measuring slope or curvature; the general trend of the waveform will be zero because of the differentiation that takes place.

### 3.4.3 Autocorrelation function (ACF)

There are a number of ways in which the ACF can be evaluated; which one is chosen depends on time and space factors.

The most obvious method is to store all the initial filtered data, say $m$ points $z_1, z_2, ..., z_m$, and then to evaluate, step by step, the correlation coefficient corresponding to each lag point, making sure that the average values are taken. Thus for a lag $j\Delta x$ the correlation array

$$A(j\Delta x) = \frac{1}{m-j+1} \sum_{i=1}^{m-j+1} z_i z_{i+j} \left/ \frac{1}{m} \sum_{i=1}^{m} z_i^2 \right. .$$

(3.174)

This method requires that all the data is stored, which may well not be convenient. Under these circumstances an alternative way is possible which, although taking longer, requires considerably less storage. The data is sifted and operated on as it becomes available from the input; the correlation array is built up in parallel rather than serial as in the first method.

Organizationally the operation is as shown in figure 3.29. After the $m$ ordinates have been operated on the array $A(1(r)j)$ is normalized with respect to the variance of $z_i$ with the mean value removed and the result in $A(1(r)j)$ *is* the autocorrelation function. Only $j$ data storage locations have been required as opposed to $m$ for the previous method. Because $j$ is of the order of $m/10$ a considerable gain in space can be achieved. This method is used when storage is at a premium. Fortunately, with the spectacular increase in storage capability schemes like this are becoming redundant.



**Figure 3.29** Autocorrelation computation.

By far the most often used method is that involving the fast Fourier transform simply because of the time saving. The method is different from the others because it involves getting the power spectral density (PSD) first and then deriving the autocorrelation from it.

### 3.4.4 Autocorrelation measurement using FFT

Schematically the calculation follows the steps below:

Profile $z(x) \rightarrow z(k\Delta x)$ (sampled profile) fast Fourier transform
↓
DFT discrete Fourier transform using fast Fourier transform
↓
$P(h\Delta f)$ power spectral density (PSD)
↓
DIFT discrete inverse Fourier transform using FFT
↓
$A(i\Delta \tau)$ autocorrelation function (ACF).

In this technique the basic FFT operation is performed twice, once to get the power spectrum and once to get the autocorrelation function from it. This is the reverse of the old method in which the power spectrum is obtained from the ACF.

Apart from the obvious time advantage there is a storage advantage. For autocorrelation this involves a realization that two steps are required and, furthermore, because of the nature both of the data and of the transformation some economy can be made. In this case the original data—the profile, say—is real. This means that the DFT will be conjugate even (the amplitude terms have even symmetry about the DC origin whilst the phases have odd) and consequently only half the transform needs to be evaluated—the other half can be deduced and the storage allocated to these points can be used elsewhere. Remember that the original data points are completely replaced in the store by an equal number of transformed points in the FFT. Furthermore, in undergoing the inverse discrete Fourier transform a further gain can be effected because, although the data (in this case the Fourier transform of the original profile) is conjugate even, only the cosine transformation needs to be used—the phase is lost in the autocorrelation function. Hence space which would normally be used to house the imaginary part of the complex numbers (resulting from the FFT) can be used for other purposes.

All these methods use the time average version of the autocorrelation function for evaluation rather than the ensemble average. Thus,

$$A(\tau) = E[z(x)z(x + \tau)]. \tag{3.175}$$

This choice is for purely practical reasons. It is easier to evaluate $A(\tau)$ from one record than to measure it from a number of records in parallel. It does assume ergodicity, however.

A similar function known as the structure function $S(\tau)$ can be obtained by measuring the expected variance between ordinates $z_1$ and $z_2$ separated by $\tau$. Thus, $y\sigma = R_p$,

$$S(\tau) = E[(z_1 - z_2)^2] = 2\sigma^2[1 - A(\tau)] \tag{3.176}$$

The structure function is a form of cumulative autocorrelation function. It has the advantage that any slope on the waveform is eliminated and the value $z_1—z_2$ can be evaluated directly by means of a skid-stylus combination without the need to have two independent styluses or without having to store any data (providing that $\sigma$ is known). If $\sigma$ is not known, a digital filter can be used to preprocess the data in real time using

very little storage. From the point of view of classification it is not necessary that $\sigma$ be known because it is a constant—the essential variation of the structure with the lag $\tau$ is the same as the autocorrelation function, that is

$$\frac{\mathrm{d}S(\tau)}{\mathrm{d}\tau} = -\frac{k\mathrm{d}A(\tau)}{\mathrm{d}\tau}. \tag{3.177}$$

### 3.4.5 Measurement of power spectral density ( PSD )

The power spectral density (PSD) can be evaluated from the data either via the autocorrelation function or the Fourier transform of the data. The former method is possible because of the Wiener-Khinchine relationship linking the two together. Thus in continuous form

$$P(\omega) = 2\int_{0}^{\infty} A(\tau) \cos \omega\tau \; \mathrm{d}\tau. \tag{3.178}$$

In practice the record from which the autocorrelation function has been measured is finite, of length L. Thus

$$A(\tau) = \frac{1}{L-\tau}\int_{0}^{L-\tau} f(x)f(x+\tau) \; \mathrm{d}x. \tag{3.179}$$

This limits the available length of autocorrelation that can be used to get a PSD. This truncation of the ACF causes some problems in the frequency domain. These problems are twofold: one is that the statistical reliability of the information is limited, and the other is that the shape of the truncation causes spurious side effects in the spectrum. Truncation in the extent of the record is equivalent to multiplying the data waveform by a box function. This box function when transposed into frequency is highly resonant; it produces 'ringing' around simple peaks in the spectrum, for instance frequencies corresponding to the feed of a tool. The transform of the box function is a sine function, which has a considerable lobing up to 25% of the peak value. To reduce this in practice a weighting function can be applied to the ACF prior to transformation. A criterion for the shape of this is that ringing does not occur, or at least to only a few per cent.

The box function is an example of what is in effect a 'lag window'. This lag window would normally be a box function, but it can be shaped in the correlation domain to have a minimum effect in the frequency domain. The frequency equivalent of the lag window is called a spectral window. The criterion for a well-behaved spectral window is that it should have a highly concentrated central lobe with side lobes as small and rapidly decaying as possible [24].

The most used lag window $W_L(\tau)$ is that due to Hanning, which has the formula

$$W_L(\tau) = 0.5 + 0.5 \cos(\pi\tau/\tau_{\mathrm{max}}) \tag{3.180}$$

where $\tau_{\mathrm{max}}$ is the maximum autocorrelation lag allowed by reliability considerations. An alternative is also due to Hamming

$$W_L(\tau) = 0.54 + 0.46 \cos(\pi\tau/\tau_{\mathrm{max}}). \tag{3.181}$$

Some other examples are shown in figure 3.30, which also illustrates the well-used Gaussian window. This is unique in that the spectral window is the same shape as the lag window.

Assuming that the Hanning window is chosen the digital equivalent formula for the power spectrum is given by

$$P(\omega) = 2\Delta\tau \sum_{k=0}^{N} W(k\Delta\tau)A(k\Delta\tau) \cos(\omega k\Delta\tau) \qquad (3.182)$$

For the power spectrum the sample points are taken $\Delta\tau$ apart, usually equal to $\Delta x$, the spacing of the ordinates or a multiple of it.

Another way to measure the power spectrum is directly from the data using the FFT routine described earlier. First the periodgram $|F(\omega)|^2$ is obtained by transforming the real data. This will yield $N$ transform points corresponding to $N$ real data points. To get the PSD from this it is necessary to apply the spectral window corresponding to the Hanning (or other) lag window. Now this is operated on the frequency data by means of convolution. Thus

$$P(\omega_0) = 0.25(F(\omega_{-1}))^2 + 0.5(F(\omega_0))^2 + 0.25(F(\omega_{+1}))^2 \qquad (3.183)$$

for the Hanning window, where $\omega_{-1}$, $\omega_0$ and $\omega_{+1}$ are adjacent frequencies in the array. For the special case where $\omega = 0$

$$P(0) = 0.5(F(0))^2 + 0.5(F(\omega_{+1}))^2 \qquad (3.184)$$

and for the Hamming window

$$P(\omega_0) = 0.23(F(\omega_{-1}))^2 + 0.54(F(\omega_0))^2 + 0.23(F(\omega_{+1}))^2. \qquad (3.185)$$



**Figure 3.30** Windows for the correlation function.

## 3.5 Properties and implementation of the ambiguity function and Wigner distribution function

### 3.5.1 General

As seen earlier in chapter 2 the ambiguity function and the Wigner distribution are closely allied to the Fourier transform and the autocorrelation and can be applied in a similar way.

Some properties will be given for comparison with the Fourier transform without proof. Other properties and discussion are given in section 3.8 rather than the digital properties described here.

### 3.5.2 Ambiguity function

This is

$$A_f(\chi, \overline{\omega}) = \int_{-\infty}^{\infty} f\left(x + \frac{\chi}{2}\right) f^*\left(x - \frac{\chi}{2}\right) \exp(-j\overline{\omega}x)\,dx \tag{3.186}$$

and

$$A_f(\chi, \overline{\omega}) = \frac{1}{2\pi} \int_{-\infty}^{\infty} F\left(\omega + \frac{\overline{\omega}}{2}\right) F^*\left(\omega - \frac{\overline{\omega}}{2}\right) \exp(+j\overline{\omega}\chi)\,d\omega \tag{3.187}$$

and is shown in figure 3.31.



**Figure 3.31**

#### 3.5.2.1 Spatial shift

$$f(x - x_0) \xleftrightarrow{\ \ A\ \ } A_f(\overline{\omega}, \chi)\exp(-j\overline{\omega}x_0) \tag{3.188}$$

where $\xleftrightarrow{\ \ A\ \ }$ denotes 'ambiguity function of'.

#### 3.5.2.2 Frequency shift

$$f(x)\exp(j\omega_0 x) \xleftrightarrow{\ \ A\ \ } A_f(\overline{\omega}, \chi)\exp(j\omega_0 x). \tag{3.189}$$

#### 3.5.2.3 Spatial-limited signals

If $f(x)$ is restricted to $[x_a, x_b]$, then $A_f(\overline{\omega}\ \chi)$ is restricted to $[-(x_b - x_a), (x_b - x_a)]$ with respect to $\chi$.

#### 3.5.2.4 Frequency-limited signals

If $f(a)$ is band limited to $[\omega_a, \omega_b]$, then $A_f(\overline{\omega}\ \chi)$ is limited to $[-(\omega_b - \omega_a), (\omega_b - \omega_a)]$ in terms of $\overline{\omega}$.

### 3.5.2.5 Concentration of energy

$$| A_f(\overline{\omega}, \chi) | \le \int_{-\infty}^{\infty} | f(x) |^2 \, \mathrm{d}x = A_f(00).$$

(3.190)

### 3.5.2.6 Total energy

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} | A_f(\overline{\omega}, \chi) |^2 \, \mathrm{d}\overline{\omega}\mathrm{d}x = \| f(x) \|^4 .$$

(3.191)

### 3.5.2.7 Convolution

If $g(x) = f(x){*}h(x)$, then

$$A_g(\overline{\omega}, \chi) = \int_{-\infty}^{\infty} A_f(\overline{\omega}, \chi)A_h(\overline{\omega}, \chi - x)\mathrm{d}x.$$

(3.192)

### 3.5.2.8 Modulation

If $g(x) = f(x)m(x)$, then

$$A_g(\overline{\omega}, x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} A_f(\omega, x)A_m(\overline{\omega} - \omega, x)\mathrm{d}\omega.$$

(3.193)

### 3.5.2.9 Discrete ambiguity function ( DAF )

For the discrete form of a complex signal $f(n)$, $n = ... -2, -1, 0, 1$, the discrete ambiguity function (DAF) $A_f(v, \eta)$ is

$$A_f(v, \eta) = \sum_{n=-\infty}^{\infty} 2f[n + \eta]f^*[n - \eta]\exp(-\mathrm{j}2v\eta)$$

(3.194)

or

$$A_f(v, \eta) = \frac{1}{\pi} \int_{-n}^{n} F[\theta + v]F^*[\theta - v]\exp(\mathrm{j}2\theta\zeta)\mathrm{d}\theta.$$

(3.195)

Note that the function is discrete in the spatial domain and continuous in the frequency domain. Note also the frequency doubling.

As in the continuous case

$$f[n] \xleftarrow{\quad A \quad} A_f[v, n]$$

(3.196)

where $A$ is the ambiguity transformation.

### 3.5.2.10 Computation of DAF ( where n is restricted to 0, N—1, and N is even )

1. Compute $r[\eta, n] = 2f[n + \eta]f^*[n-\eta]$ for $n = 0, 1,..., N-1$, $\eta = -M,..., 0,..., M$ where $M = -N/2-1$.

2. Using DFT compute

$$\sum_{n=0}^{N-1} r[\eta, n] \exp\left(-j\frac{k\pi n}{N}\right)$$

(3.197)

for $\eta = -M, ..., 0, ..., M$.

3. $A_f(\bar{\omega}, \eta)$ is then obtained for $\bar{\omega} = k\pi/N$ where $k = 0, 1, ..., N-1$ and $\eta = -M, ..., 0, ..., M$. Because $A_f(\bar{\omega}, \eta)$ is restricted in $\eta$ and periodic in $\bar{\omega}$, so $A_f[fk, m] = A_f[k\pi/N, \eta]$ is obtained for $k, \eta = -2, -1, 0, 1, 2$; for odd $N$ it is similar provided that it can be represented by $N = N_1, ..., N_2, ..., N_k.\Delta$

The point to notice about the DAF is that the variable is the spatial $x$ and spectral $\omega$, shown in figure 3.32. It covers an area in $x, \omega$ space of area $\chi\bar{\omega}$; for this reason features of the order $X\bar{\omega}$ will be highlighted in the space-frequency domain—hence the comment that the ambiguity function is good for identifying the presence or absence of features ($\chi \bar{\omega}$). However, it is not good for positioning absolutely in space or frequency; in other words, it is not very useful for finding the non-stationarity of a signal.



**Figure 3.32** Movement of variables in ambiguity space.

The box $\chi\bar{\omega}$ scanning the space-frequency domain acts as a correlation with $\chi, \bar{\omega}$ being fixed for all $x$, and $\omega$ and acting as space-frequency lags. The instantaneous values extend by $x' \pm \chi/2$ and $\omega' \pm \bar{\omega}/2$, where $x'\omega'$ are the positions at the centre of the domain in $x', \omega$ of extent $\chi$ and $\bar{\omega}$, as seen in figure 3.32.

### 3.5.3 The Wigner distribution function

This

$$W_f(x, \omega) = \int_{-\infty}^{\infty} f\left(x + \frac{\chi}{2}\right).f^*\left(x - \frac{\chi}{2}\right)\exp(-j\omega\chi)d\chi$$

$$= \frac{1}{2\pi}\int_{-\infty}^{\infty} F\left(\omega + \frac{\bar{\omega}}{2}\right).F^*\left(\omega - \frac{\bar{\omega}}{2}\right)\exp(j\bar{\omega}x)d\bar{\omega}.$$

#### 3.5.3.1 Properties

$$f(x) \xleftarrow{\quad W \quad} W_f(x, \omega).$$

(3.198)

$W$ is Wigner transformation.

### 3.5.3.2  Symmetry

$$W_f(x,\omega) = W_f(x,-\omega). \tag{3.199}$$

### 3.5.3.3  Realness

For any complex-valued signal the Wigner distribution is always real:

$$W_f(x,\omega) = (W_f(x,\omega))^*. \tag{3.200}$$

### 3.5.3.4  Spatial shift

$$f(x + x_0) \xleftrightarrow{\ W\ } W_f(x + x_0, \omega). \tag{3.201}$$

### 3.5.3.5  Frequency shift

$$f(x)\exp(j\omega_0 x) \xleftrightarrow{\ W\ } W_f(x, \omega - \omega_0). \tag{3.202}$$

### 3.5.3.6  Spatial-limited signal

If $f(x)$ is restricted to $[x_a, x_b]$, so also is $W_f(x, \omega)$.

### 3.5.3.7  Frequency limiting

If $f(x)$ is restricted to $[\omega_a, \omega_b]$, so is $W_f(x, \omega)$.

### 3.5.3.8  Spatial energy

$$\frac{1}{2\pi}\int_{-\infty}^{\infty} W_f(x,\omega)\mathrm{d}\omega = |\,f(x)\,|^2. \tag{3.203}$$

### 3.5.3.9  Frequency energy

$$\int_{-\pi}^{\pi} W_f(x,\omega)\mathrm{d}x = |\,F(\omega)\,|^2. \tag{3.204}$$

### 3.5.3.10  Total energy

$$\frac{1}{2\pi}\int_{-\infty}^{\infty}\int_{-\infty}^{\infty} W_f(x,\omega)\mathrm{d}x\,\mathrm{d}\omega = \|\,f(x)\,\|^2. \tag{3.205}$$

The integral of the Wigner distribution over the whole plane $(x\omega)$ is the total energy of the signal $f(x)$.

### 3.5.3.11 Convolution

If $g(x) = f(x)*h(x)$, then

$$W_g(x,\omega) \int_{-\infty}^{\infty} W_f(\chi,\omega)W_h(x-\chi,\omega)\mathrm{d}\chi. \tag{3.206}$$

### 3.5.3.12 Modulation

If $g(x) = f(x),(x)$, then

$$W_g(x,\omega) = \frac{1}{2\pi} \int_{-\infty}^{\infty} W_f(x,\overline{\omega})W_m(x,\omega-\overline{\omega})\mathrm{d}\overline{\omega}. \tag{3.207}$$

### 3.5.3.13 Analytic signals

Note

For real-valued signals mainly found in surface texture, in order for the Wigner distribution to be applied the signal should be converted to the analytic form (whose frequency content $= 0$, $\omega \le 0$). Thus, taking the Hilbert transform to obtain the analytic signal yields

$$\hat{f}(x) = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{f(\chi)}{x-\chi}\mathrm{d}\chi \tag{3.208}$$

so that $f_a = f(x) + \mathrm{j}f(x)$.

In discrete form

$$W_f(n,\theta) = \sum_{\eta=-\infty}^{\infty} 2f(n+\eta)f^*(n-\eta)\exp(-\mathrm{j}2\theta\eta)$$

$$W_f(n,\theta) = \frac{1}{\pi} \int_{-\pi}^{\pi} \sum_{\eta=-\infty}^{\infty} F(\theta+v),F^*(\theta-v)\exp(\mathrm{j}2vn)\mathrm{d}v \tag{3.209}$$

where $F(v)$ is the DFT of $f(n)$ and

$$f_a(n) = f(n) + \mathrm{j}\hat{f}(n)$$

where $\hat{f}(n)$ is the discrete Hilbert transform of $f[n]$, $H_d f(n)$:

$$\text{discrete } \hat{f}(n) = (H_d f)[n] = \sum_{\eta=-\infty}^{\infty} f[\eta]\frac{2}{\pi}\frac{\sin^2(\pi(n-\eta)/2)}{(n-\eta)}. \tag{3.210}$$

### 3.5.3.14 Moments

The discrete moments are important in surface analysis:

$$\text{zeroth in frequency } P_f(n) = \frac{1}{2\pi} \int_{-\pi/2}^{\pi/2} W_f(n,\theta)\mathrm{d}\theta$$

$$P_f(n) = |f(n)|^2. \tag{3.211}$$

Because the argument in the Wigner distribution has two variables, moments exist in both space and frequency as well as global moments involving both.

The first order, $\theta_f(n)$, is

$$\theta_f(n) = \frac{1}{2} \arg\left( \int_{-\pi/2}^{\pi/2} \exp(j2\theta) W_f(n,\theta) d\theta \right).$$

(3.212)

Then

$$\theta_f(n) = \frac{1}{2} \arg\left( f(n+1) f^*(n-1) \right).$$

(3.213)

If $f(n) = v[n] \exp(j\varphi(n))$, then

$$\theta_f(n) = \frac{\varphi(n+1) - \varphi(n-1)}{2} \mod \pi$$

(3.214)

that is, the first-order moment in frequency represents the instantaneous frequency, being the differential of the phase $\varphi$ over the region.

The second-order moment in frequency, $M_f(n)$, is

$$M_f[n] = \left( P_f(n) - \left| \frac{1}{2\pi} \int_{-\pi/2}^{\pi/2} \exp(j2\theta) W_f(n,\theta) d\theta \right| \right) \bigg/ \left( P_f(n) - \left| \frac{1}{2\pi} \int_{-\pi/2}^{\pi/2} \exp(j2\theta) W_f(n,\theta) d\theta \right| \right).$$

(3.215)

$M_f(n)$ can be expressed as

$$M_f(n) = \frac{\lfloor f(n) \rfloor^2 - \mid f(n+1) f^*(n-1) \mid}{\mid f(n) \mid^2 + \mid f(n+1) f^*(n-1) \mid}.$$

(3.216)

The zeroth-order moment in space, $P_f(\theta)$, is

$$P_f(\theta) = \sum_{n=-\infty}^{\infty} W_f(n,\theta) = \mid F(\theta) \mid^2 + \mid F(\theta + \pi) \mid^2.$$

(3.217)

The first-order moment $X_f(\theta)$ is

$$X_f(\theta) = \sum_{n=-\infty}^{\infty} W_f(n,\theta) n / P_f(\theta).$$

(3.218)

For analytic signals $X_f$ can be expressed as

$$X_f(\theta) = -\text{Im} \frac{(d \ln F(\theta))}{d\theta} = -\text{Im}\left( \frac{F'(\theta)}{F(\theta)} \right)$$

(3.219)

and similarly the second order as

$$M_f(\theta) = -\frac{1}{2} \text{Re} \frac{d}{d\theta} \left( \frac{F'(\theta)}{F(\theta)} \right)$$

(3.220)

(see figure 3.33).

**Figure 3.33** Movement of variables in Wigner space.

The Wigner distribution (figure 3.33) centres on specific $x$ and $\omega$ and scans via $\chi$ and $\bar{\omega}$ over all the area. It is in effect a two-dimensional convolution with $\chi$ and $\bar{\omega}$ as dummy variables. The two-dimensional integral gives complete information about the spot $x'$, $\omega'$.

It was shown in chapter 2 how the moments can be used to extract information about amplitude modulation, frequency modulation and phase modulation effects as well as chirp. Also, the spatial moments pick out the position and width of pulses. The dual character of this function enables the extraction of the most relevant information from the correlation technique and the spectral technique with the benefit of having to use only one function.

### 3.5.4  Some examples of Wigner distribution: application to signals—waviness

It was suggested earlier that the Wigner function might be useful in characterizing waviness. This is because the moments of the function clearly reveal the type of modulation that might be causing the waviness even if it is non-linear, where ordinary spectral decomposition breaks down.

According to the theory given earlier, the local moments in frequency are

$$M_{f0} = |z(x)|^2 = |F(\omega)|^2 = |a|^2 \qquad \text{the instantaneous power}$$

$$M_{f1} = \mathrm{Im}\left\{\frac{z'(x)}{z(x)}\right\} = \varphi'(x) \qquad \text{the frequency of the signal}$$

$$M_{f2} = -\tfrac{1}{2}\frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{a'(x)}{a(x)}\right) = \qquad \text{the envelope of the signal}$$

where the right-hand sides of the equations are for a signal $z(x) = a(x)\exp(j\varphi(x))$.

For a frequency-modulated type of surface or a non-linear chirp signal it seems evident that the first moment in frequency would reveal the most information because this moment is sensitive to the phase, whereas the second moment responds to the envelope. Hence the amplitude of the envelope is most likely to be seen here, so the figures in figure 3.34 would be easily separated. As an example, consider a signal with a quadrate change in frequency with distance—commonly called a chirp signal, $z(x) = a \exp(j(\alpha/2)x^2)$

$$\text{zero order} = a^2$$
$$\text{first order} = \alpha x$$
$$\text{second order} = 0.$$

**Figure 3.34** The zeroth and first moments of the Wigner distribution for a chirp signal and a frequency-modulated signal.

For a frequency-modulated signal $z(x) = a \exp[j(\omega_{0x} + \varphi_0 + b\sin(\omega_m x + \varphi_m))]$ the moments are:

$$
\begin{aligned}
\text{zero order} \quad &= a^2 \\
\text{first order} \quad &= \omega_0 + b\omega_m \cos(\omega_m x + \varphi_m) \\
\text{second order} &= 0
\end{aligned}
$$

which are obviously very distinguishable from each other.

For a signal with amplitude modulation $z(x) = a(x)\exp(j\phi(x)x)$

$$
\text{the zeroth moment} \quad = a^2(x)
$$

$$
\text{the first moment} \quad = 0
$$

$$
\text{the second moment} = -\tfrac{1}{2}\frac{\mathrm{d}}{\mathrm{d}x}\left[\frac{a'(x)}{a(x)}\right].
$$

Hence it can be seen that use of the local frequency moments can isolate the various forms that envelopes can take. This is most likely to be important in functional cases.

### 3.5.5 *Comparison of the Fourier transform, the ambiguity function, the Wigner distribution function and wavelet transform*

A comparison of the properties of three transforms has been given in table 3.4. All three are based on the Fourier kernel and so can be seen to have similar properties. However, whereas the Fourier transform is in either the frequency domain or the space domain the ambiguity and Wigner functions are in both.

It is useful to think of the space-frequency domain as a 'music score'. The only difference is that the space-frequency 'score' shows the intensity as well as the pitch and spacing of the notes.

The ambiguity function as seen in figure 3.32 could be considered to be a way of scanning across the whole score for a spectral pattern or a time/space pattern within a given size space-frequency frame. The Wigner function on the other hand tends to give a space-frequency 'feel' for the content of the score centred on a specific $x$ and $\omega$ as seen in figure 3.33. The common feature of both is the 'score' itself. The Fourier transform only gives the axes!

The wavelet transform is another alternative (ref. 58) see later in Transformations.

## 3.6 Digital estimation of reference lines for surface metrology

### 3.6.1 *Numerical filtering methods for establishing mean lines for roughness and waviness profiles*

Filtering methods are the natural way to isolate specific bands of information of the surface. Obviously, originally the most usual way to apply filtering techniques was to pass the signal through a passive or active electrical network so that the breakdown of the signal occurred in frequency. This is not the natural domain of the surface, which is spatial, so the recent introduction of digital methods has been an important step in surface analysis. Each digital measurement can now be referred more easily to a point on the surface. Fundamentally the great advantage of filters, as mentioned in the previous chapter, is that they take the signal as it is and do not assume any particular waveform. They take the waveform 'as received' and operate on it, unlike best-fit polynomial curves which can completely distort the remnant signal if the order of the polynomial is wrong when compared with the general shape of the surface.

### 3.6.2 Convolution filtering

The first step is to work out the impulse response of the filter and express it in digital form or sequence. Thus, let $h(n)$ be the digital sequence representing the impulse response and $z(n)$ be the profile signal. Then the output from the filter $g(n)$ becomes

$$g(n) = \sum_{m=0}^{n} h(m)z(n-m)$$

(3.221)

where $h(m)$ and $z(n-m)$ are zero outside the limits.

The relationship between the weighting function and the impulse response has been explained in section 2.2. The working formula is

$$g = k(a_0 b_0 + a_1 b_1 + \dots + a_m b_m).$$

(3.222)

Usually $a_1 a_2$ are equally spaced ordinates, but not necessarily; $K$ is a constant equal to

$$\frac{1}{\text{number of ordinates per sample length}}.$$

(3.223)

The $b$ terms are the digitized values of the weighting function (reversed impulse response of the filter).

### 3.6.3 Standard filter

The required spacing, and the associated value of $K$, are shown in table 3.5. These values are shown to agree with the old standard filter.

How the weighting function appears digitally with respect to the profile (shown with the stylus) is illustrated in figure 3.35. Here the weighting function is that of the $2CR$ filter.

**Table 3.5**

| No of ordinates per cut-off length | $K$ | Ordinate spacing (in) on surface for: | | |
|---|---|---|---|---|
| | | 0.01 in cut-off | 0.03 in cut-off | 0.1 in cut-off |
| 50 | 0.02 | 0.0002 | 0.0006 | 0.002 |
| 150 | 0.0067 | 0.0000 667 | 0.0002 | 0.000 667 |

| No of ordinates per cut-off length | $K$ | Ordinate spacing (mm) on surface for: | | |
|---|---|---|---|---|
| | | 0.25 mm cut-off | 0.08 mm cut-off | 2.5mm cut-off |
| 50 | 0.02 | 0.005 | 0.016 | 0.05 |
| 150 | 0.0067 | 0.001 67 | 0.003 33 | 0.0167 |

**Figure 3.35** Application of weighting function.

For the sake of clarity, the procedure will be illustrated by the example in figure 3.36, where the work is set out for manual calculation, even though it would take far too much time and space to work it like this in practice, unless far fewer ordinates were used merely as an exercise.

Assume that the horizontal magnification is 100 and that the mean line is required for 0.025mm cut-off. From table 3.5, the required ordinate spacing on the graph is seen to be 0.15mm. The profile ordinates $a_0$, $a_1$, $a_2$, etc, are therefore measured at this spacing throughout the length of the graph.

Starting, for example, from the right-hand side, ordinates are taken in groups of 100, say, and each ordinate in the group is multiplied by the appropriate weighting factor, the process being continued until there are 100 products. The products are then added algebraically, noting that the first 28 are always positive and the remainder negative, and the sum (imagine it to be 35.5 in) is multiplied by the factor $K$, in this case 0.02, that is 1/50. Then 35.5in $\times$ 0.02 = 0.71 in is the value of the mean line ordinate $z'_0$ (or 0.71 $\times$ 2.54 mm).

The procedure is now repeated for the next required mean line ordinate. For drawing the mean line through the profile graph, it will generally be sufficient to calculate five per cut-off length, so that the ordinate $a_{10}$ will become the second starting point. These products, taken up to $a_{10} \times b_{99}$, will provide the mean line ordinate $z'_{10}$. The third starting point will be the profile ordinate $a_{20} \times b_0$, the products being calculated up to $a_{119} \times b_{99}$, to provide the mean line ordinate $y'_{20}$. The procedure is continued for as far along the profile as desired.

### 3.6.4 *Phase-corrected ( linear phase ) filters, other filters and filtering issues*

An exactly similar procedure can be carried out for any other weighting function, including the phase-corrected (or linear phase) filter described earlier.

A pictorial demonstration of how the impulse response approach is constructed can be seen in figures 3.37 and 3.38. For continuity the 2$CR$ filter is shown but the principle is exactly the same for all other weighting functions. It should be noticed, however, that whereas the weighting function for the standard type of filter is the *inverse of the impulse* response in the case of phase-corrected filters, the weighting function is the same as the impulse response because the function impulse response is an even function in $x$.

**Figure 3.36** Application of weighting function to get mean line of surface texture.

In the figure, the labels at top: $Z_{40}$, $Z_{30}$, $Z_{20}$, $Z_{10}$, $Z_0$

$$Z_{10} = K(a_0 b_0 + \dots + a_{99} b_{99})$$

$$Z_{10} = K(a_{10} b_0 + \dots + a_{109} b_{99})$$

$$Z_{20} = K(a_{20} b_0 + \dots + a_{119} b_{99})$$

Cut-off length   Cut-off length



**Figure 3.37** (*a*) Profile, (*b*) unit impulse response of 2*CR* filter, (*c*) inverted impulse response (without impulse), (*d*) tabulated version of weighting function.

To get the mean line value multiply each weighting factor by the profile ordinate in line with it, add the products and divide by the normalizing factor, which is $1/k$. This gives the height of the mean line above the same datum and in the same units used for the profile ordinates.

Each mean line ordinate that has been calculated in this way refers to *that profile ordinate which is directly in line with the maximum value of the weighting function or centre of symmetry*. The weighting factor corresponding to the maximum $b_0$ of the weighting function corresponds to the first term in causal systems $b_0$ for example in the standard 2*CR* filter.

A comparison of figures 3.35 and 3.39 immediately illustrates the difference in the shape of the two weighting functions and, more importantly, it shows the fact that the mean line position acts in a different place relative to the weighting function for a phase-corrected weighting function than it does for a weighting function which is not phase-corrected.

Subtraction of the mean line value from the profile ordinate to which it refers gives the filtered profile at that point. Taking the average of the differences (without taking account of sign) over the number of cut-off

**Figure 3.38** Pictorial demonstration of mean line procedure: (*e*) peak line weighting function $1/RC_1$ exp $(\tau/RC)(2+\tau/RC)$; (*f*) mean line weighting function shifted to $t_1$; (*g*) weighted profile at $t_1$; (*h*) mean line ordinate at $\tau_1$ plotted on profile; (*j*) mean line ordinate at $t_1$, and $t_2$ plotted on profile; (*k*) mean line and profile.

lengths (sampling lengths) that make up the traversing length (e.g. five lengths for a 0.03 inch or 0.8 mm cut-off) gives the $R_a$ value as normally determined. Other parameters can be measured in a similar way.

### 3.6.5 *Gaussian filter*

An interesting property of the Gaussian function can be used in one application in surface metrology, and this is the use of a Gaussian filter to exclude waviness. The property is that if a number of windows are convoluted together they will always eventually produce an effect which is Gaussian, irrespective of the shape of the window. A good example is the box function (figure 3.42). Three convolutions of the box function will produce an equivalent weighting function window which is already very close to Gaussian (figure 3.40).

**Figure 3.39** Application of linear phase weighting function.



**Figure 3.40**

Some instrument manufacturers advocate the Gaussian filter simply because it has got this property. It means that if a running-average procedure is repeated three times, the end result is as if the profile signal had been subjected to a low-pass Gaussian filter. The very simplicity of the technique makes it fast and inexpensive.

In this form the weighting function has a low-pass characteristic, that is the waviness profile. To get the surface roughness value of the surface data corresponding to the mid-point of the weighting function should be selected and the waviness line, as found by this Gaussian weighting function taken from it (figure 3.41).



**Figure 3.41**

The Gaussian shape has some advantages. Probably the most important is that it is recognizable in industry, being the basis for the acceptance curve used in statistical process control. This means that production engineers are comfortable working with it. Also the curve is always positive and it falls off rather faster than the $2CR$ transmission curve. Large components of waviness therefore do not impinge on the roughness.

The weighting function is given by

$$h(x) = \frac{1}{\alpha \lambda_c} \exp\left[-\pi\left(\frac{x}{\lambda_c \alpha}\right)^2\right]$$

(3.224)

and its transform giving the transmission curve of (3.229) where $x$ is measured from the axis of symmetry.

Thus

$$H\left(\frac{1}{x}\right) \text{is} \exp\left[-\pi\left(\frac{\lambda_c \alpha}{x}\right)^2\right]$$

where

$$\alpha = \sqrt{\frac{\ln}{\pi}} = 0.4697.$$

Another more practical reason for using Gaussian type filters is that they minimize the RMS duration of the product of $h(x)$ and $H(w)$ where $w = \frac{2\pi}{x}$.

Bodschwinna [26] proposed a modified form of the Gaussian filter incorporating a second order polynomial to reduce the influence of end effects caused by the finite duration of the weighting function. The combination of least square polynomial as well as filtering is somewhat messy. However, he suggests a form $h(x)$ given by

$$h(x) = \frac{1}{\lambda_c}\sqrt{\frac{\pi}{C}}\left[1.5 - \frac{\pi^2}{C\lambda_c^2}x^2\right]\exp\left(-\frac{\pi^2}{C\lambda_0^2}x^2\right)$$

(3.225)

which has a transmission.

$$H\left(\frac{1}{x}\right) = \left[1 + C\left(\frac{\lambda_c}{\lambda}\right)^2\right] \exp\left(-C\left(\frac{\lambda_c}{\lambda}\right)^2\right)$$

(3.226)

Krystek [26] uses a spline to get rid of form and waviness and he produces a complicated form for the filter. It is questionable whether the extra complication is needed.

Wavelets have also been proposed as a means of filtering. The ability to vary the resolution and range – usually in octave bands [25] makes them useful in fractal analysis [27]. The term 'mathematical zoom lens' has been used for wavelet analysis. Apart from the standard filtering use of wavelets, which has a questionable credibility it seems that defect detection using 'raised wavelets' has some benefits [28].

### 3.6.6 Box functions

Whilst on the subject of box functions, it should be recognized that, when convoluted with the profile signal, a box function is a simple running average, the extent of the average being the length of the box. This average is taken to be at the mid-point of the box (figure 3.42).



**Figure 3.42**

The process of using the averaging procedure was advocated by Reason [29] as a means of producing a mean line of the surface texture. It was then called the 'mid-point locus' line. Unfortunately it has a poor frequency response (figure 3.43). That is,

$$B\left(\frac{x}{\lambda_c}\right) \Leftrightarrow \frac{\sin[2\pi|(\lambda_c/x|]}{2\pi|\lambda_c/x|}.$$

(3.227)

However, the 'averaging' procedure produced by a running box function is a convolution and, therefore, the whole process acts as a filter. Furthermore, it is in effect a phase-corrected filter providing that the output is taken to act at the centre of the box. It therefore represents the first attempt at phase-corrected (or linear phase) working in surface metrology.

It is obvious that in order to get no distortion of the signal but just an attenuation at certain frequencies a phase-corrected characteristic should be used. Another way to get phase-corrected characteristics, which is simple, is to use a 'double-pass' method. This is useful when only analogue methods are available, but the technique has also been used digitally. To see how this technique works imagine a filter has an impulse response $h(t)$ whose Fourier transform is $H(\omega)$. Then if the signal is put through this filter and the output *reversed* and put through the same filter again the final output is

$$\left|H(\omega)\right|^2 \exp(-j\omega T).$$

(3.228)

**Figure 3.43** Mid-point focus characteristic.

This is achieved because turning the output round is equivalent to conjugation, that is

$$h(T - t) \Leftrightarrow \exp(-j\omega T)h^*(\omega) \tag{3.229}$$

where $T$ is the time taken to reverse the output signal and enter it into the filter again. The second pass (in the reversed mode) is simply

$$\exp(-j\omega T)H^*(\omega) \text{ yielding overall } \left|H(\omega)\right|^2 \exp(-j\omega T). \tag{3.230}$$

This equation has the linear phase term $\exp(-j\omega T)$ in it. This is phase corrected about the same time $T$.

Physically this means that to enable the signal (suitably filtered) to take on phase-corrected characteristics it has to be delayed by time $T$. It is impossible to get phase-corrected properties without the delay because an impulse response cannot be an even function about $t = 0$ and all realizable systems are causal (i.e. do not exist before $t = 0$). This argument is the same in spatial terms.

To use this technique it is essential that the square root of the intended final transmission can be taken, that is $H(\omega)$ must exist. This means that any characteristic, if capable of being expressed as $H(\omega)^2$, can be made into a phase-corrected filter by using this reversal technique. Once this is done the calculation is very fast and the storage small.

### 3.6.7   Truncation

The finite length of the impulse response may call for comment. Many impulse responses have infinite extent (IIR) (infinite impulse response). Practically, in communication theory they are difficult to use because of the need to store masses of data points. What happens is that the impulse response is curtailed as in the case when only 100 ordinates are used. This must result in errors. These can be quantified easily enough because, instead of the frequency output being $G(\omega) = H(\omega)Z(\omega)$, it becomes

$$G'(\omega) = H(\omega) * B(\omega)Z(\omega), \tag{3.231}$$

where $B(\omega)$ is the frequency spectrum of the truncating function, usually a box function. This shows the other operation performed by the box function, a multiplication rather than a convolution in averaging.

The effect of truncation is shown in figure 3.44.

**Figure 3.44** Truncated function (left) and normalized Fourier transform (right) for: (*a*) no truncation, (*b*) and (*c*) truncating function.

The amount of truncation allowable is usually decided by the percentage error produced on one of the averaging parameters, such as the $R_a$ value. It has been found that making the weighting function between two and three cut-offs long suffices for an error of less than 5%. Even this can be greatly reduced by increasing (or decreasing) the weighting factor ordinates so that, if the truncation length is $a$, then

$$\int_{-a/2}^{a/2} h(x)\mathrm{d}x = 1 \tag{3.232}$$

rather than the ideal

$$\int_{-\infty}^{+\infty} h(x)\mathrm{d}x = 1. \tag{3.233}$$

If this crude compensation is used then at least the *average* frequencies are correct.

Later in this chapter recursive filters will be discussed. Here it will be found that they have less of a truncation problem than the convolution methods because the equivalent weighting function builds up as the number of evaluated points increases.

Errors in parameters resulting from truncated weighting functions are easy to evaluate for $R_q$ and, hence, assuming a correction factor, for $R_a$. But for peak parameters the problem is much more complicated.

However, it can be argued that communication criteria should not be dominant in surface metrology. Two questions have to be asked. The first is whether an infinite impulse response means anything functionally when two surfaces make contact. The second is whether the requirements for filtering should be different for the functional applications of surfaces and for the control of manufacture.

Referring to the first question, it could well be that the effective impulse response could use, as a weighting function, the pressure ellipse of the Hertzian contact zone [30] (figure 3.45). The term functional filtering has been used [31,32]. The advantage of such an approach is that it more nearly represents what goes on in a functional situation. The disadvantage is that the instrument maker cannot provide every possibility. The obvious answer is to let the customer choose the filter nearest to his/her requirement.



**Figure 3.45** Use of weighting function in functional simulation.

The advantage of the convolution method, using a weighting function representing a functional condition and not the impulse response of an electrical filter, has many advantages. One is that it is versatile. The filter characteristics can be chosen to fit the function, within reason; the function can effectively be shaped to fit. The disadvantage is that the arithmetic operation is expensive on storage and takes a long time to perform. Simple tricks like the equal-weight method outlined below can be used. There are, however, other methods which preserve the equal-spacing criteria of ordinates.

### 3.6.8 Alternative methods of computation

The first technique is called the overlap-add method. If the profile is long ($N_1$) and the weighting function $h(n)$ short ($N_2$), the profile is split up into samples (of $N_3$ ordinates). Convenient values of $N_2$ would correspond to the sampling length

$$z(n) = \sum_{i=0}^{\infty} z_i(n).$$

(3.234)

The convolution becomes

$$g(n) = \sum h(m) \sum_{i=0}^{\infty} z_i(n-m)$$

$$= \sum_{i=0}^{\infty} h(n) * z_i(n) = \sum_{i=0}^{\infty} g_i(n).$$

(3.235)

The duration of each of the convolutions of equation (3.240) is $(N_3 + N_2 - 1)$ samples—so there is a region of $(N_2 - 1)$ samples over which the $i$th convolution overlaps the $(k + 1)$ convolution and the outputs from each therefore overlap and so have to be added; hence the name overlap-add when referred to the outputs.

The second method to be used, which achieves an output reasonably quickly, is called the overlap-save method. This differs from the previous one in that it involves overlapping input sections rather than output sections and is used most often in the filtering of periodic signals, as in the case of roundness. This technique makes use of circular convolution.

In this method the input is sectioned into overlapping sections of length $L$ and overlap $M - 1$. Then each of the input sections is convoluted with the filter and the resultant outputs are added together, but the first $M - 1$ outputs from each section are discarded because these will have been the last $M - 1$ values of the previous section [22].

These methods are valuable in metrology when using conventional computers because they do allow some output to be obtained very quickly rather than to have to wait until all the input data is stored. This time factor is not essential for post-process inspection but can be serious when in-process measurement is envisaged.

Different methods can also be used to speed up the processing time, one of which follows.

The evaluation of filtered results in the frequency domain can obviously be done by taking the data, doing a fast Fourier transform on it, and then multiplying the spectrum by the desired filter characteristic. To get the filtered profile all that is needed is to carry out an inverse FFT. This is much quicker because all convolution operations are replaced by multiplications. However, if functional filtering is the objective, for example in contact situations, the spatial convolution method is more suitable because the weighting function can be shaped (as, for example, to the shape of a pressure distribution) to run across the surface (figure 3.45).

### 3.6.9 Equal-weight techniques

In the general convolution form for filtering shown earlier, the output at time $t$ is given by $g(t)$ where

$$g(t) = \text{constant} \sum_{i=0}^{n} a_i b_i. \tag{3.236}$$

The constant here is a constant of proportionality dependent on the ordinate spacing $h$. In this equation the $b_i$ represent weighting function ordinates derived from the impulse response of the filter. Letting the constant be $h'$ (equal to the reciprocal of the density of ordinates in a sampling length), the equation can be written in a different form which makes better use of the data. The equation shows that each ordinate $a_i$ is multiplied by a weighting function: $b_i \times h'$ represents a strength or area of weighting function to be associated with ordinate $a_i$ as in figure 3.46. Each ordinate of the profile coming within the weighting function has a multiplier associated with it of a different value (i.e. different area in the numerical integration). It is possible to simplify the calculation by arranging that some profile ordinates have equal-value multipliers instead. The filtering is achieved by the selection of ordinates. Thus, in figure 3.46, only those ordinates opposite crosses would be operated on. The areas shown are all of equal value equal to $V$, say. Then the output of the filter becomes

$$V \times \sum_{i=0}^{k} a_i. \tag{3.237}$$

This method has a number of advantages. First $k$ is usually much smaller than $N$, perhaps 15 instead of 200. Second the $a$ values of the profile are simply added and not multiplied, therefore it is much less prone to freak values than the *equal-spacing* weighting function. The only disadvantage is that the $a_i$ values are not equally spaced: their location is inversely proportional to the height of the weighting function. Their location addresses have to be stored in much the same way that the weighting function ordinates in the conventional form have to be stored. However, the look-up process is much faster than the multiplication process and gains in speed of an order of magnitude are possible. This numerical procedure for equal weight can also be applied to a graphical method—which needs to be possible so as to enable calibration routines to be established for people without computers.

**Figure 3.46** (*a*) Numerical application of weighting function (i) equal interval, (ii) equal interval, phase corrected; (iii) equal weight. (*b*) Graphical filtering using equal-weight rule.

So, instead of multiplying equally spaced ordinates each by a different weighting factor, a useful approximation can be arrived at by multiplying unequally spaced ordinates by a constant factor, which reduces to adding up the graph ordinates at these spacings and multiplying by the factor (or dividing by its reciprocal if that is more convenient). The use of 26 ordinates was found to give reasonable accuracy and a convenient divisor of 20. A 26-ordinate template designed for 0.8mm cut-off and a horizontal magnification of 100 is shown over a profile graph in figure 3.46.

In use, having laid the template over the graph with the zero lines of graph and template in coincidence, the successive ordinates are read off, added together and divided by the stated constant to give the mean line ordinate at the front edge of the rule. The last three ordinates must always be taken negatively. Thus

$$y' = \frac{a_0 + a_1 + \ldots + a_{21} + a_{22} - a_{23} - a_{24} - a_{25}}{20}. \tag{3.238}$$

Because of the discontinuous nature of the integration, the values of the mean line ordinates determined with the template may be a little larger or smaller than their true values, but in the planimetric determination of the $R_a$ that will finally be made, these fluctuations will tend to cancel out.

### 3.6.10 Recursive filters

The convolution method requires a lot of storage and computation. The output depends on previous and present inputs. Since a filter effectively contains a memory of past inputs, it might be expected that considerable savings in total effort could be obtained by calculating the past and present inputs and past outputs.

This technique is analogous to feedback in a linear control system and for this reason is called recursive. With conventional filtering methods such as that involved in the $2CR$ filter, the savings in time and storage over the convolution method are substantial. For instance, the example given earlier, which involved 100 weighting factors requires 100 multiplications and additions for each output point. The recursive approach reduces this to the order of only four each.

To appreciate this method (which is similar in every way to the ARMA models used for surface characterization examined in chapter 2) an example will be given.

### 3.6.11 The discrete transfer function

A convenient notation for the discussion of sampled systems is the $z$ transform in which the term $z^{-1}$ is the unit-delay operator. Then a sequence of sampled values can be represented as the coefficients of a power series in $z^{-1}$. For instance, the infinite sequence $X = 1, B, B^2, B^3,...$ can be written as

$$X(z) = 1 + Bz^{-1} + B^2 z^{-2} + \ldots$$
$$= \frac{1}{1 - Bz^{-1}} \quad \text{for } B < 1. \tag{3.239}$$

The unit-delay operator is related to the $z$ transform operator as

$$z^{-k} @ \delta(n - k) \tag{3.240}$$

where $T$ is the sampling interval so that although digital filters can be directly realized, it is more usual to approach the problem from a linear analysis. The Laplace operator is $p$.

There are several ways of converting the continuous filter into a digital one. The standard method is that known as the impulse-invariant technique in which the discrete impulse response is identical to the sampled impulse response of the continuous filter. In this the continuous function is expanded into partial fractions and then transformed by

$$\frac{1}{p + a} \rightarrow \frac{1}{1 - \exp(-aT)z^{-1}}. \tag{3.241}$$

Tables exist for transformations of this type. The difficulty of this approach is that the gain of the digital filter is proportional to the sampling frequency, which may cause problems if this frequency is not always the same value. Also, it may be noted that the continuous transfer function may not easily resolve into a form suitable for transformation.

A way of avoiding the need to expand the transfer function is to map the $p$ plane onto the $z$ plane.

A suitable way of doing this is to use the bilinear substitution:

$$p \rightarrow \frac{2}{T}\left(\frac{1 - z^{-1}}{1 + z^{-1}}\right). \tag{3.242}$$

One of the main difficulties of this method may be shown by considering the unit circle in the 2 plane. Now

$$z = \exp(-j\omega_d T)$$

$$p = \frac{z-1}{z+1}$$

becomes

$$p = j\tan\left(\frac{\omega_d T}{2}\right) \tag{3.243}$$

so that, as required, the imaginary axis of the $p$ plane is mapped onto the unit circle of the $z$ plane. However, it also follows that

$$\omega_a = \tan\left(\frac{\omega_d T}{2}\right). \tag{3.244}$$

Thus the method causes a non-linear warping of the frequency scale. It would, therefore, generally be necessary to precompensate the continuous function.

A third method uses an approximation of the convolution integral on a continuous transfer function that has been either factorized or resolved into partial fractions. This method, of which the impulse-variant approach is a special case, has a gain that is independent of sampling rate and is more readily applied to a high-pass filter than the impulse-invariant transformation.

As computational requirements suggested the use of a cascade approach and the $2CR$ filter causes no factoring problems, the third method of transformation can be used and this will be discussed in more detail.

For convenience in the following analysis, the single-stage filter will be considered throughout. The two-stage filter is derived directly from it by squaring the transfer function. The basic form of filter section is the low pass and such sections in cascade have their transfer functions multiplied. So the general low-pass filter has the form

$$H(p) = k\prod_{i=1}^{N}\frac{1}{p+a_i}. \tag{3.245}$$

Considering a single stage

$$H(p) = \frac{Z(p)}{X(p)} = \frac{1}{p+a} \tag{3.246}$$

and expressing its time domain response by the convolution integral, assuming the filter to be initially relaxed gives, for the output $z(t)$ and input $x(t)$,

$$z(t) = \int_0^t \exp[-a(t-\tau)]x(\tau)d\tau$$

$$z(t) = \exp(-at)\int_0^t \exp(a\tau)x(\tau)d\tau. \tag{3.247}$$

Now, in a digital system, time is necessarily discrete (it is no longer 'real time') and equation (3.247) may be set to successive values $nT$ and $nT-\tau$, where $n$ is an integer:

$$z(nT) = \exp(-anT)\int_0^{nT} \exp(a\tau)x(\tau)\mathrm{d}\tau$$

(3.248)

$$z(nT - T) = \exp[-a(nT - T)]\int_0^{nT-T} \exp(a\tau)x(\tau)\mathrm{d}\tau.$$

(3.249)

Multiplying equation (3.253) by $\exp(-aT)$ and subtracting from equation (3.252) gives

$$z(nT) = \exp(-aT)z(nT - T) + \exp(-anT)\int_{nT-T}^{nT} \exp(a\tau)x(\tau)\mathrm{d}\tau.$$

(3.250)

This integral requires that an analytic expression for $x(\tau)$ exists within the interval $(nT - T)$ to $nT$. An approximate solution can be obtained by applying the constraint, which is in any case normally imposed by the analogue-to-digital conversion, that $x(\tau)$ is constant in the interval $nT - T \le \tau < nT$. Then equation (3.250) becomes

$$z(nT) = \exp(-aT)z(nT - T) + \exp(-anT)x(nT - T)\frac{\exp(anT)}{a}[1 - \exp(-aT)]$$

$$z(nT) = \exp(-aT)z(nT - T) + \frac{[1 - \exp(aT)]}{a}x(nT - T).$$

(3.251)

Since time $(nT - T)$ is one unit delay from time $nT$ the discrete transfer function can be found directly to be

$$H(z^{-1}) = \frac{Z(z^{-1})}{X(z^{-1})} = \frac{1}{a}\frac{[1 - \exp(-aT)]z^{-1}}{(1 - \exp(-aT)z^{-1})}.$$

(3.252)

Having established the transfer function for the low-pass section, the other elementary sections can be derived. For instance, the bandpass can be expressed as the difference of two low-pass sections. Of particular interest is the high-pass section, which has a transfer function

$$H(p) = \frac{p}{p + a} = 1 - \frac{a}{p + a}.$$

(3.253)

Thus the high-pass is realized as a direct link of the profile in parallel with and opposite to a low-pass section. However, in the discrete system, if this approach were used, the fact that time is not continuous means that the signal applied to the low-pass section could not affect the output until a finite time later and so would not combine properly with the direct-link signal. The filter response would thus deteriorate, so it is necessary, for best performance, to introduce a compensating delay in the direct-link path. To find the required delay, it is necessary to estimate the phase angle of the discrete low-pass section. Using equations (3.251) and (3.253) and grouping the constants gives

$$H(p) = \frac{A\exp(-pT)}{1 - B\exp(-pT)}.$$

(3.254)

Substituting $p = j\omega$ and applying de Moivre's theorem gives

$$H(\omega T) = \frac{A(\cos\omega T - \mathrm{j}\sin\omega T)}{1 - B\cos\omega T + \mathrm{j}B\sin\omega T}$$

(3.255)

so that the phase angle is

$$\varphi = -\omega T - \beta - \tan^{-1}\left(\frac{B \sin \omega T}{1 - B \cos \omega T}\right) \tag{3.256}$$

where $\beta$ is a constant.

For $B$ close to unity, which is equivalent to the product $|aT|$ being much less than unity, equation (3.256) simplifies to

$$\varphi = -\omega T - \beta - \tan^{-1}\left(\cot \frac{\omega T}{2}\right). \tag{3.257}$$

From this, the group delay of the section is

$$\frac{-\mathrm{d}\varphi}{\mathrm{d}\omega} = \frac{T}{2} \tag{3.258}$$

so that a delay of half a sampling period should be introduced into the direct path. The high-pass section now has the form

$$h_{\mathrm{hp}}(z^{-1}) = z^{1/2} - \frac{Az^{-1}}{1 - Bz^{-1}} \tag{3.259}$$

which is not a realizable linear digital filter since its transfer function is not rational in $z^{-1}$. A further transformation is therefore required and a suitable one is $z^{-1} \rightarrow z^{-2}$ which in the continuous filter would be $p \rightarrow 2p$. Thus, to retain the same frequency response under this transformation, the cut-off frequency must be doubled and the fully realized high-pass section is

$$H_{\mathrm{hp}}(z^{-1}) = z^{-1} - \frac{[1 - \exp(-2aT)]z^{-2}}{1 - \exp(-2aT)z^{-2}}. \tag{3.260}$$

Thus the required two-stage high-pass filter has the transfer function [30]

$$H(z^{-1}) = \left| z^{-1} - \frac{[1 - \exp(-2aT)]z^{-2}}{1 - \exp(-2aT)z^{-2}} \right|^2. \tag{3.261}$$

The discrete transfer function of equation (3.260) has a pole near the Nyquist frequency ($1/2T$ Hz) which is caused by the frequency-doubling procedure of the foregoing analysis and so will not be usable near this frequency.

It is the correction for group delay used here that gives the improved high-pass characteristics of the convolution integral approach over that of the impulse-invariant method. If the analysis embodied in equations (3.255)-(3.259) were carried out on the impulse-invariant-derived transfer function, it would be found that the group delay was $T/2$, that is, a half time-period advance would be needed as compensation. Although this could be arranged when processing previously stored data, it is not a physically realizable system since it requires information about the future. It may be noted that this problem is similar to that encountered when a phase-corrected filter is being realized, because a true linear phase method requires both linear and group delay correction.

The single-stage low-pass section of equation (3.252) can be rearranged to give

$$z = AXz^{-1} + BZz^{-1} \tag{3.262}$$

where $X$ and $z$ are the input and output sequences respectively.

The high-pass section described by equation (3.260) can be implemented as it stands or multiplied out to give

$$H_{hp}(z^{-1}) = \frac{z^{-1} - Az^{-2} - Bz^{-3}}{1 - Bz^{-2}}.$$

(3.263)

It will be noted that the parallel representation, equation (3.252), requires less arithmetic than the canonical form. In general, therefore, this approach should be both faster and more accurate.

### 3.6.12 The 2CR filter

A particular application of the recursive 2CR filter is in the fast assessment of profiles. With the limited memory size available at the time, the saving of the need to store a weighting function is significant and, for on-line work, the filter should preferably work in 'real time'. In this sense, 'real time' means that the time taken to process one sample in the filter is less than the smallest sampling period used. For the reasons discussed in the previous sections, the chosen method of implementation could be a cascade implementation of the 'convolution-approximation' section.

The requirements of other programs with which the filter has to be used demands that, initially, there should be 400 samples per cut-off. (This requirement is to ensure that a reasonably well-defined 100:1 bandwidth for average wavelength measurements could be produced.) However, at small cut-offs (sampling lengths) this represents samples considerably less than one stylus width apart and requires high sampling rates. (It is likely that, in general, less severe restrictions will be imposed.) As an example, for the two-stage filter to demonstrate 75% transmission at the cut-off it can be shown that for each stage

$$a = \frac{\omega_c}{\sqrt{3}} = \frac{2\pi f_c}{\sqrt{3}}$$

(3.264)

where $f_c$ is the cut-off frequency. Relating this to wavelength

$$a = \frac{2\pi}{\sqrt{3}\lambda_c}$$

(3.265)

or, if there are $N$ ordinates per cut-off,

$$a = \frac{2\pi}{\sqrt{3}NT}.$$

(3.266)

For 50% transmission at the cut-off, (the new proposition) $a$ would be simply $2\pi/\lambda c$.

Thus for an $N$ of 400, the $\exp(-2aT)$ will be very close to unity and, consequently, $1 - \exp(-2aT)$ will be near zero. The multipliers that are used in the filter are fractional and so cannot be handled by the integer arithmetic. There are methods of pre- and postscaling coefficients and coordinates which allow fractional quantities to be handled but in this case the near zero value is likely to introduce considerable errors. In the interests of accuracy it may be better to use floating point arithmetic although, since this depends on computer software, it is slower in operation.

Recursive methods can only be used in cases where the discrete transfer function can be written in terms of a rational ratio of polynomials such as

$$H(z^{-1}) = \frac{a_0 + a_1 z^{-1} + a_2 z^{-2} \dots}{1 + b_1 z^{-1} + b_2 z^{-2} \dots}.$$

(3.267)

The question arises of whether this can be suitable for linear phase (phase-corrected) filters. Another point is whether the poles and zeros of equation (3.267) are realizable. In practice the latter is decided by the transmission characteristics of the filter and whether or not it suffers discontinuities. From the earlier discussion on optimized filters it is clear that the 3:1 and 2:1 filters are not well behaved in so far as they have discontinuous slope in the transform and so cannot be directly implemented in a recursive mode. However, the fact that these filters cannot be implemented does not mean that linear phase filters as a whole cannot. One simple way of implementing filters of the kind shown above relies on the ability to describe the transfer function in a square root mode.

Remember that, if the impulse response is $h(t)$ whose transform is $H(\omega)$, then reversing the output can be shown to be equivalent to conjugation, with an additional linear phase term; that is, if $h(t) \Leftrightarrow H(\omega)$ then

$$h(T - t) \Leftrightarrow \exp(-j\omega T)H^*(\omega) \tag{3.268}$$

where $H^*(\omega)$ is the complex conjugate of $H(\omega)$. Hence, passing the signal through a filter, storing the output, reversing it and passing it back through the filter yields

$$H(\omega)H^*(\omega)\exp(j\omega t) = |H(\omega)|^2 \exp(-j\omega t). \tag{3.269}$$

If a recursive method is possible then this method is very fast.

### 3.6.13 Use of the FFT in surface metrology filtering—areal case

The FFT routine can be used very effectively in surface metrology because it is very easy to implement any desired filter characteristic. Because the input data (i.e. the surface) is real, the storage is reduced by a factor of 2 relative to conventional complex data. Once transformed the desired amplitude characteristic is simply a multiplication of the transform data by the required characteristic and then this is followed by an inverse FFT to reconstitute the filtered signals. Phase effects are invariant if the linear phase filter is being simulated.

As shown in the earlier section on FFTs there is a very big speed advantage of filtering in this way. However, there is the problem that all the data has to be stored before filtering can take place. This is not the case in the direct convolution method, so depending upon the storage capability either one or the other can be used. It is invariably a compromise between storage, speed and accuracy. From the latter point of view the effective weighting function length is, as in the recursive method, equal to the length of the data.

Whatever is true for profile filtering is even more true of two-dimensional filtering. One method due to Rabiner and Gold [8] allows a 2D filtering to take place in terms of profile filtering methods. If $x, y$ are integer numbers in the $xy$ plane and $v, \omega$ in the frequency domain, then the two-dimensional transform of $z(x,y)$ is $Z(v,\omega)$ where

$$Z(v,\omega) = \sum_{x=0}^{N_1-1} \sum_{y=0}^{N_2-1} z(x,y) \exp\left[-j\left(\frac{2\pi x\omega}{N_1} + \frac{2\pi yv}{N_2}\right)\right]$$
$$= \sum_{x=0}^{N_1-1} \sum_{y=0}^{N_2-1} z(x,y) \exp\left[-j\left(\frac{2\pi x\omega}{N_1}\right)\right] \exp\left[-j\left(\frac{2\pi yv}{N_2}\right)\right] \tag{3.270}$$

or

$$Z(v,\omega) = \sum_{x=0}^{N_1-1} \exp\left[-j\left(\frac{2\pi x\omega}{N_1}\right)\right]\left\{\sum_{y=0}^{N_2-1} z(x,y) \exp\left[-j\left(\frac{2\pi yv}{N_2}\right)\right]\right\}. \tag{3.271}$$

The term in braces in equation (3.271) is a series of $N_1$ one-dimensional discrete Fourier transforms (DFTs) obtained by varying $x$ from 0 to $N_1 - 1$ as shown in figure 3.47.

**Figure 3.47** Areal FFT.

So the transforms for each profile can be taken with $x$ fixed for each profile, yielding $N_1$ transforms. These give the transform for a given $\omega$ given $Z(v,\omega)$. Once $Z(v,\omega)$ is mapped it can be multiplied (because it is in the frequency domain) by whatever 2D filtering characteristic is required, so fast convolution using the FFT is probably the most important means for realizing two-dimensional filters. Note, however, that every one of the profile graphs shown in the figure has to have mechanical fidelity relative to each other. It is no good measuring each profile using a skid instrument because this will lose the reference level for each profile and the two-dimensional spectrum will have little macrogeometry spectrum detail.

The filtered characteristic would become

$$Z_f(v,\omega) = Z(v,\omega)W(v,\omega) \tag{3.272}$$

where $W$ is the amplitude characteristic of the required filter in two dimensions. From this the filtered surface $z_f(x,y)$ emerges as

$$z_f(x,y) = \sum_{v=0}^{N_1-1}\sum_{\omega=0}^{N_2-1} Z(v,\omega)W(v,\omega) \exp\left[j\left(\frac{2\pi xv}{N_1}+\frac{2\pi y\omega}{N_2}\right)\right]. \tag{3.273}$$

This is much more efficient than the direct convolution. A typical value of the number of points in a surface map is 250 000. To give some idea of the advantage in time of this method, although the ratio of data points (assuming 250 per profile) is 250:1 the ratio of the FFT method to the direct method is about 50:1 for a profile and about 14 000:1 for the surface area, so the gain in time is tremendous. The problem for surface metrology is visualization; a true contact acting over an area is easy to imagine and is a functional effect. Unfortunately, because its effect is non-linear it cannot readily be transformed into the frequency domain. It can therefore be advantageous to develop convolution-type operations in the spatial domain in order to simulate contact effects.

### 3.6.14 *Examples of numerical problems in straightness and flatness*

Early on the best-fit line was given as $m$ (equation (2.18)) where

$$m = \frac{\sum x \sum z - N \sum xz}{(\sum x)^2 - N \sum x^2} \tag{3.274}$$

and $m$ is relatively small.

If the spacings of measurement in the $x$ direction are taken to be of equal increments then the slope becomes

$$m = \frac{12\sum_{i=1}^{N} iz_i - 6(n+1)\sum_{i=1}^{N} z_i}{N(N^2-1)} = \sum_{i=1}^{N} z_i \left[ \frac{12i}{N(N^2-1)} - \frac{6(N+1)}{N(N^2-1)} \right] \tag{3.275}$$

or

$$m = \sum_{i=1}^{N} z_i (k_1 i - k_2) \tag{3.276}$$

where

$$k_1 = \frac{12}{N(N^2-1)} \quad \text{and} \quad k_2 = \frac{6}{N(N-1)}. \tag{3.277}$$

So, from the computational point of view, each $z$ can be multiplied by a weighting factor which, for a given fixed $N$, can always be the same.

Similarly for flatness, the slopes become

$$M_1 = \frac{12\sum_{i=-N/2}^{N/2} (i\sum_{j=-M/2}^{j=M/2} z_{ij})}{MN(N+1)(N+2)} \tag{3.278}$$

$$M_2 = \frac{12\sum_{j=-M/2}^{j=M/2} (j\sum_{i=-N/2}^{N/2} z_{ij})}{MN(M+1)(M+2)} \tag{3.279}$$

assuming that the measurements are being taken from the centre of the coordinate system.

Typical of the dilemmas that can arise is the play-off between the quantization and the sampling of a roughness waveform to determine the best-fit least-squares line. Usually one requires a small quantization interval and high sampling to get good results. But if one of these is not available it may not be advantageous to hold to the other; a mutual relaxation may sometimes be necessary. In this example, the least-squares mean line slope $m$ as above is given by

$$m = \frac{12\sum_{i=1}^{N} iz_i - 6(N+1)\sum_{i=1}^{N} z_i}{N(N^2-1)} \tag{3.280}$$

where the $z$ values are profile ordinates taken at $N$ equal unit intervals along the surface. The essential numerical point to notice is that the numerator is made up of the difference between two summations. Two quantization factors influence the value of $m$ that is obtained; first, the resolution of the digital measurement of the analogue signal, and second the resolution in the computer, that is the wordlength. The best way to show these effects is by a simple example.

Suppose that five measurements are taken on a surface using a very high-resolution A/D converter and let these be processed in a computer of a long wordlength. Let the numbers as measured be 10.000, 10.000, 10.000, 10.000, 10.1999. Substituting these into the equation for $m$ yields

$$m = \frac{1811.91 \pm 1807.164}{120} = 0.0398 = 2.2790°.$$

If now the same measurements had been taken with an A/D convertor capable of seeing only three decimal digits, the numbers would be 10.0, 10.0, 10.0, 10.0, 10.2 giving

$$m = \frac{1806 - 1803.6}{120} = 0.0200 = 1.1458°.$$

Again if a high-resolution A/D convertor was used and the computer could only work to four decimal digits (13 bits), then

$$m = \frac{1811 - 1807}{120} = 0.0333 = 1.909°.$$

Finally if the A/D convertor measures three digits and the computer four, then

$$m = \frac{1806 - 1803}{120} = 0.025 = 1.432°.$$

Four different answers have been obtained (2.279°, 1.1458°, 1.909° and 1.432°) using what were intended to be the same data from the surface profile—nearly 100% variation!

Part of the discrepancy is due to the small variations between the numbers being suppressed by the limited resolution of the A/D convertor and part due to the small difference between the two summations in the numerator being suppressed or modified by the limited wordlength. Both have entered into it. Yet another factor has to be taken into account and this is the sampling rate. If there were three times as many numbers, each of the originals counted three times, what effect does this have? If the A/D convertor has three digits and the computer four, the answer should be $0.025 = 1.432°$ in the five-ordinate case, but it is $0.0089 = 0.511°$ for the 15 ordinates 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.0, 10.1, 10.1, 10.1, which contain exactly the same information. The reason for this further discrepancy is that merely adding more numbers into each of the summations of the numerator eventually leads to the least significant numbers (which contain the vital information) being rounded off or truncated. This is due to the limited wordlength. Again, as before, too much sampling aggravates the problems associated with quantization.

What is to be done? A number of possibilities are open. The first is only to use independent samples; do not use redundant information by sampling too fast. Next make sure that the real variations are preserved in the computer. This can be achieved in the example above by removing an estimate of the mean value of the data from each measurement before working out the numerator. Under these circumstances, the data becomes 0, 0, 0, 0, 0.1 and the danger of the summation overflowing the wordlength limitation is considerably reduced. In addition to these measures, it is always possible to use double wordlength arithmetic if the computer allows it, and it is also possible to increase the resolution of the A/D convertor at the expense of extra work and slower data capture rates. In general it is more advantageous to anyone working on the numerical analysis problem than it is on validating the data itself. The latter is of limited use without the former.

### 3.6.15 *Effect of computer word format*

Often spurious resolution can appear to develop during a calculation, especially when converting from one type of word in the computer to another, for instance from integer to floating arithmetic. The number 10 becomes, say, 10.000; the decimal digits to the right of the decimal point are not significant in this case. Suppose the number 10 represents a profile ordinate that the A/D convertor could not resolve better than the unity digit. The signal itself may have been anywhere between 10 and 11. Certainly the probability of it being 10.000 is remote, but this does not mean that all floating point numbers derived from integer numbers are incorrect. As an example, if a profile is made up of a string of integer numbers 583, 621, 718, etc, then the low-pass-filtered profile may be expressed as 592.3 for instance, because the uncertainty in a mean (weighted mean for filtering) is less than that of individual numbers. It still does not mean, however, that all the printout of the floating point numbers is significant. As a rule of thumb, if $q$ is the uncertainty in the individual num-

bers then the mean can be taken significantly to about one decimal digit further than the individual values. Note that if the A/D converter is two digits and the wordlength of the computer is 16 bits (four digits), the last digit is not likely to be significant if a simple convolution type of exercise is being carried out. This has nothing to do with numerical analysis problems—it is purely the accuracy in the values themselves, whether theoretically correct or not!

Another form of number notation called the scientific or E notation appears at first sight to offer both high resolution and range, that is 0.63152 E 23 would mean $0.63152 \times 10^{23}$. The reason for this apparent benefit is that most computers use two words to present a number in the scientific notation. Thus, for a 16 bit word, 32 bits are available, usually 24 bits for the mantissa and 6 bits for the exponent, so an increase in accuracy is only at the expense of store. Again whether this accuracy is real or not depends on the initial data and the arithmetic operations carried out.

Problems of this nature are always arising in curve fitting, matrix inversion, etc. Their individual solution depends on the particular problem and the type of computer. There is always an improvement if care is taken.

Some obvious checks should be made to ensure that sensible values are being obtained. One example of this is in the measurement of slope. Here the model, quantization and sampling again conspire to confuse the issue (see figure 3.48). This figure shows a case where all three are poor: the quantization interval is too big, the sampling too fast and the model (say three-point) is too restrictive. Let the sampling interval be $q / 10$.



**Figure 3.48** Balance needed between quantization; model and sampling example in slope measurement.

Then the slopes measured at points 1, 2, 3, . . . will all be unity until the point E is reached. Then the slope will be $q / 2 \times q / 10 = 5$, giving an angle of 78°! Not only does this look ridiculous, it cannot make physical sense if the device used for obtaining the analogue signal is a stylus with a semi-angle of 45°. Real slopes greater than this cannot be seen. Also, curvature measurement as revealed digitally cannot or should not be greater than that of the stylus itself! Common-sense rules like this often show up numerical errors. A point to note here is that, although some emphasis has been placed upon the problems of digital analysis, they are in no way simplified by reverting to analogue methods. What usually is the case is that the digital method forces attention onto the real problems. Too often the analogue signal (or assessment) is regarded as sacred.

## 3.7  Algorithms

### 3.7.1  *Differences between surface and dimensional metrology algorithms: least-squares evaluation of geometric elements*

In assessing the departure from straightness, roundness, flatness, sphericity, cylindricity, conality and other geometric forms the method of least-squares occupies a dominant place, being the preferred method often included in international and national standards. Although this would appear to be straightforward, it is not so because of the many algorithms that could be used to achieve the assessment. Depending on which is used the result could be obtained quickly and accurately or, sometimes, not at all. Also, another problem that arises is very often the parameter is not linear, in which case, prior to applying least-squares some linearization has to take place, from which an approximate solution is obtained and used in the next iteration.

In the text so far the least-squares method has been used in a number of situations. There is, however, one general approach that might prove to be useful. This is based on the technique developed at NPL by Forbes [33] for use in coordinate-measuring problems. The technique is valid for surface metrology problems providing linearization is taken into account. The algorithms employ a common approach which has stable parameterization of the elements. What has to be watched carefully is the nature of the data obtained from the measuring instrument. Surface metrology instruments are different from dimensional measuring machines. In measuring an arc, for example, both can produce different forms of data; one can usually easily be linearized, the other cannot! Very unstable results can be produced if one sort of algorithm is applied to the wrong sort of data.

The actual shapes of elements are those already considered in chapter 2 and which are specified in BS 7172 [34].

The rationale is as follows. It concerns parameterization, which in turn concerns the way in which the problem of solving for the best fit solution is posed.

As an example, the best-fit plane will be discussed. This plane or any other can be specified by a point in it, say $(x_0, y_0, z_0)$, and the direction cosines of the normal perpendicularly such that $a^2+b^2+c^2=1$. Thus, any plane can be defined by six parameters which are not all independent of each other.

If a plane exists as above and it is required to see how a set of points fit to it the intuitive thing to do is to measure the distance of each point from it. For a point $x, y, z$ the distance $d$ is given by

$$d = a(x - x_0) + b(y - y_0) + c(z - z_0).$$
(3.281)

The $i$th point has a distance $d_i$ :

$$d_i = a(x_i - x_0) + b(y_i - y_0) + c(z_i - z_0).$$
(3.282)

and the one way to estimate the goodness of fit is to look at $S = \sum d_i^2$. This sum $S$ depends on the parameters of the plane $x_0, y_0, z_0$ and $a, b, c$. These parameters have to be chosen to minimize $S$.

The steps for a least-squares fit are therefore as follows.

1. Choose parameters to describe the position, shape and sometimes also the orientation and size of the geometrical part.
2. Derive a formula for the distance of a point from this geometric object.
3. Express the sum of squares $S$ in terms of a set of data points and their distances from the geometric object.
4. Develop a stable algorithm for determining the parameters such as $a, b, c, x_0, y_0, z_0$.

A precursor to evaluating the best-fit least-squares to improve the numerical accuracy is to find the centroid of the data points $\bar{x}, \bar{y}, \bar{z}$ and to remove these from the general data points so that the new $x_i, y_i, z_i$ are equal to the original $x_i, y_i, z_i - \bar{x}, \bar{y}, \bar{z}$.

The case for the best-fit line and best-fit plane reduce to the simple case of solving a matrix equation, that is finding the eigenvectors of matrices.

### 3.7.1.1 Optimization

In general the function $S(u) = \sum_i^M (u)$ has to be minimized with $m$ data points and, say, $n$ parameters to be minimized with respect to $u$ where

$$u = (u_1, u_2, \ldots, u_n)^{\mathrm{T}} \tag{3.283}$$

where T indicates the transpose $u,$ the vector form of $u$.

### 3.7.1.2 Linear least squares

Here $d_i$ is a linear function of the parameters $u$ and there exist constraints $a_{ij}$ and $b_i$ such that

$$d_i = a_{i1}u_1 + \ldots + a_{ij}u_j + \ldots a_{in}u_n - b_i. \tag{3.284}$$

This is a set of linear equations which can be put in the matrix form $\mathsf{A}u = b$ :

$$\begin{pmatrix} a_{11} & a_{12} & a_{in} \\ a_{21} & a_{22} & a_{2n} \\ \vdots & \vdots & \vdots \\ a_{ml} & \ldots & a_{mnl} \end{pmatrix} \begin{pmatrix} u_1 \\ u_2 \\ \vdots \\ u_m \end{pmatrix} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}. \tag{3.285}$$

In general $m > n$ so that all the equations will not be satisfied simultaneously. Reliable and efficient algorithms do exist to solve (3.285) in the least-squares sense:

$$\mathsf{A}^{\mathrm{T}}\mathsf{A}u = \mathsf{A}^{\mathrm{T}}b. \tag{3.286}$$

These are the 'normal equations'.

In chapter 2 it has already been shown that, in surface metrology, this approach is really sufficient, because in the case of a circle and sphere, for example, the limaçon approach basically changes the equation problem from being that of a quadratic to that of being linear. *It is this feature which differentiates the coordination measurement from that of the surface metrology approach.* Rarely is it impossible to linearize the system from the point of view of surface metrology. This is because the 'surface skin', of whatever shape, is in practice so very much smaller than the dimension or position of the geometric object. *Surface metrology instruments are built to see only the former; coordinate-measuring machines have to cater for the latter.*

However, because the two methods are converging, as the scale of size reduces, it is informative to take the coordinate-measuring approach.

### 3.7.1.3 Eigenvectors and singular value decomposition

For linear systems a good algorithm follows.

Given a square matrix $B$, an eigenvector $u$ of $B$ is such that

$$Bu = \lambda u \tag{3.287}$$

for some eigenvalue $\lambda$. The case in question is such that

$$B = A^T A, \tag{3.288}$$

as in (3.287), for some $m \times n$ rectangular matrix $A$ where $m > n$. In this situation a stable numerical solution is obtained by finding a 'singular value decomposition' (SVD) of the matrix $A$. In this $A$ can be written as a product.

$$A = USV^T \tag{3.289}$$

with $U$ and $V$ orthogonal matrices and $S$ a diagonal matrix containing the singular values of $A$. If $B$ is as in (3.288) the squares of the diagonal elements of $S$ are the eigenvalues of $B$ and the columns of $V$ are the corresponding eigenvectors. These are usually produced as standard output from most software implementations. SVD is now standard for many solutions. See, for example [39].

### 3.7.2 Best-fit shapes

### 3.7.2.1 Best-fit plane

The steps are as follows.

(1) Specify point $x_0 , y_0 , z_0$ on a plane.
(2) Evaluate direction cosines $(a, b, c)$ of a normal to a plane; note that any point on a plane satisfies

$$a(x - x_0) + b(y - y_0) + c(z - z_0) = 0.$$

(3) Evaluate distance from plane

$$d_i = a(x_i - x_0) + b(y_i - y_0) + c(z_i - z_0).$$

(4) Describe the algorithm.

The best-fit plane $P$ passes through the centroid $\bar{x}, \bar{y}, \bar{z}$ and this specifies a point in the plane $P$. It is required to find the direction cosines of $P$. For this $(a, b, c)$ is the eigenvector associated with the smallest eigenvalue of

$$B = A^T A \tag{3.290}$$

where $A$ is the $m \times 3$ matrix whose $i$th row is $(x_i - \bar{x}, y_i - \bar{y}, z_i - \bar{z})$; alternatively $(a, b, c)$ is the singular vector associated with the smallest singular value of $A$. Thus, an algorithm to find the best-fit line in 3D is:

1. Calculate the centroid $\bar{x}, y, \bar{z}$.
2. Form matrix $A$ from the data points and $\bar{x}, \bar{y}, \bar{z}$.
3. Find the SVD (singular value decomposition) of $A$ and choose the singular vector $(a, b, c)$ corresponding to the smallest singular value. The best-fit plane is therefore $\bar{x}, y, \bar{z}, a, b, c$.

Similarly for the best-fit line to data in 2D.

### 3.7.2.2 Circles, spheres, etc

These shapes involving an axis of revolution are usually evaluated by linearization of the basic equations mechanically as stated by the process of radius suppression by mechanically shifting the instrument reference to a position near to the surface skin of the geometric element being measured. Failing this an iterative method has to be used.

#### (a) Gauss—Newton method, iterative method

This can be used when the relationship between the distances $d_i$ and the parameters $u_j$ is non-linear. Hence an iterative scheme has to be used. This is similar to the Deming method of . The situation is shown in figure 3.49.



**Figure 3.49** Gauss-Newton method.

One iteration of the Newton algorithm for computing the zero of a function is as follows.
Suppose that there is a first estimate $u_0$ of where the function $u$ crosses the $u$ axis. Then:

1. Evaluate $f(u_0)$.
2. Form a tangent to the graph at $(u_0, f(u_0))$ as shown in figure 3.49.
3. Find $u_1$ where the tangent crosses the $u$ axis.

Then

$$u_1 = u_0 + \frac{f(u_0)}{f'(u_0)} = u_0 + p. \tag{3.291}$$

$u_1$ is now the new estimate of where $f(u)$ crosses the $u$ axis. This is repeated until the result is close enough to $u^*$.

Basically the Gauss-Newton method is as follows.
Suppose there is a first estimate $u$ of $u^*$. Then solve the linear least-squares system

$$\mathsf{J}p = -d \tag{3.292}$$

where $\mathsf{J}$ is the $m \times n$ Jacobian matrix whose $i$th row is the gradient of $d_i$ with respect to $u$, that is

$$J_{ij} = \frac{\partial d_i}{\partial u_j}. \tag{3.293}$$

This is evaluated at $u$ and the $i$th component of $d$ is $d_i(u)$. Finally, the estimate of the solution is

$$u := u + p \quad (:= \text{means update}).$$

These steps are repeated until $u$ is close enough to $u^*$. Ideally, changes in the iteration should be small for this method to be quick in convergence and stable.

For example, for the best-fit circle:

1. Specify circle centre $x_0$, $y_0$, radius $r$. Note that $(x - x_0)^2 + (y - y_0)^2 = r^2$.
2. Obtain distance from the circle point:

$$d_i = r_i - r$$
$$r_i = [(x_i - x_0)^2 + (y_i - y_0)^2]^{1/2}. \tag{3.294}$$

3. The elements of the Jacobian are

$$\frac{\partial d_i}{\partial x_0} = -(x_i - x_0)/r_i$$
$$\frac{\partial d_i}{\partial y_0} = -(y_i - y_0)/r_i$$
$$\frac{\partial d_i}{\partial r} = -1. \tag{3.295}$$

4. Algorithm: knowing $x_0$, $y_0$ and $r$ for the circle centre and radius estimates, use them in a Gauss-Newton iteration. Form $\mathsf{J}\,p = -d$ from the $d$ of equation (3.294) and the $\mathsf{J}$ of equation (3.295).
5. Solve

$$\mathsf{J}\begin{pmatrix} p_{x_0} \\ p_{y_0} \\ p_r \end{pmatrix} = -d \tag{3.296}$$

for $p$.
6. Update the $x_0$, $y_0$, $r$ according to

$$x_0 := x_0 + p_{x0}$$
$$y_0 := y_0 + p_{y0}$$
$$r := r + r_r. \tag{3.297}$$

Carry on until successful and the algorithm has converged.

(b)  *Linear best-fit circle*

This is an approximation described earlier used by Scott (chapter 2, reference [109]). In this model

$$S = \sum f_i^2 \tag{3.298}$$

is minimized, where $f_i = r_i^2 - r^2$ rather than $r_i - r$ as in the linear case—the trick is to make the $f_i^2$ linear. By changing the parameters $f$ can be made into a linear function of $x_0$, $y_0$ and $\rho = x_0^2 + y_0^2 - r^2$

$$f_i = (x_i - x_0)^2 + (y_0 - y_i)^2 - r^2$$
$$= -2x_i x_0 - 2y_i y_0 - (x_0^2 + y_0^2 - r^2) + (x_i^2 + y_i^2). \tag{3.299}$$

Thus

$$\mathsf{A}\begin{pmatrix} x_0 \\ y_0 \\ \rho \end{pmatrix} = \boldsymbol{b} \tag{3.300}$$

(from which $x_0$, $y_0$ and $\rho$ are found) where the elements of the $i$th row of $\mathsf{A}$ are the coefficients $(2x_i{}^1 2y_i{}^1 -1)$ and the $i$th element of $\boldsymbol{b}$ is $x_i^2 + y_i^2$.

An estimate of $r$ is

$$\sqrt{x_0^2 + y_0^2 - \rho}. \tag{3.301}$$

This can be used to get a first estimate of the parameter for the non-linear method if required.
Both the linear and non-linear methods described above can be used for spheres.

### 3.7.2.3  *Cylinders and cones*

It has been suggested [33] that a modified Gauss-Newton iterative routine should be used in the case of cylinders because one of the parameters is the direction of a line, that is the axis of the cylinder. Such a line $x_0, y_0, z_0, a, b, c$, in 3D can be specified by four parameters together with two rules to allow the other two to be obtained:

Rule 1: represent a direction $(a, b, 1)$.
Rule 2: given the direction above, ensure $z_0 = -ax_0 - by_0$.

For nearly vertical lines these two rules give stable parameterization for $a$, $b$, $x_0$ and $y_0$. The problem of finding the distance of a data point to an axis is quite complicated. The following strategy is therefore followed based on the fact that for axes which are vertical and pass through the origin, $a=b=x_0=y_0=0$ and all expressions become simple.

The strategy is as follows.

1. Iterate as usual but, at the beginning of each iteration, translate and rotate the data so that the trial best-fit cylinder (corresponding to the current estimates of the parameters) has a vertical axis passing through the origin
2. This means that when it is time to evaluate the Jacobian matrix the special orientation can be used to simplify the calculations. At the end of the iteration use the inverse rotation and translation to update the parameterizing vectors $x_0, y_0, z_0, a, b, c$, and thereby determine the new positions and orientation of the axis.

Note
To rotate a point $(x, y, z)$ apply a $3 \times 3$ matrix $\mathsf{U}$ to the vector $(x, y, z)^{\mathrm{T}}$; the inverse rotation can be achieved by using the transpose of $\mathsf{U}$:

$$\begin{pmatrix} u \\ v \\ w \end{pmatrix} = \mathsf{U}\begin{pmatrix} x \\ y \\ z \end{pmatrix}.$$
$$\begin{pmatrix} x \\ y \\ z \end{pmatrix} = \mathsf{U}^{\mathrm{T}}\begin{pmatrix} u \\ v \\ w \end{pmatrix}. \tag{3.302}$$

A simple way to construct a rotation matrix $\mathsf{U}$ to rotate a point so that it lies on the $z$ axis is to have $\mathsf{U}$ of the form

$$\mathsf{U} = \begin{pmatrix} C_2 & 0 & S_2 \\ 0 & 1 & 0 \\ -S_2 & 0 & C_2 \end{pmatrix} \begin{pmatrix} 1 & 0 & 0 \\ 0 & C_1 & S_1 \\ 0 & S_1 & C_1 \end{pmatrix} \tag{3.303}$$

where $C_i = \cos\theta$ and $S_i = \sin\theta_i$, $i = 1, 2$. So if it is required to rotate $(a,b,c)$ to a point on the $z$ axis, choose $\theta_1$ so that $bC_1 + cS_1 = 0$ and $\theta_2 = aC_2 + (cC_1 - bS_1)S_2 = 0$.

These notes are only suggestions. There are other methods that can be used but these are most relevant to geometric parts like cylinders, which in surface metrology can usually be oriented to be in reasonable positions, that is for a cylinder nearly vertical.

Care should be taken to make sure that the algorithm is still stable if reasonable positions for the part cannot be guaranteed.

*(a) Cylinder*

1. Specify a point $x_0, y_0, z_0$ on its origin, a vector $a$, $b$, $c$ pointing along the axis and radius $r$.
2. Choose a point on the axis. For nearly vertical cylinders

$$z_0 = -ax_0 - by_0 \qquad c = 1. \tag{3.304}$$

3. Distance of the chosen point to cylinder is

$$\begin{aligned} d_i &= r_i - r \\ r_i &= \frac{(u_i^2 + v_i^2 + w_i^2)^{1/2}}{(a^2 + b^2 + c^2)^{1/2}} \end{aligned} \tag{3.305}$$

where

$$\begin{aligned} u_i &= c(y_i - y_0) - b(z_i - z_0) \\ v_i &= a(z_i - z_0) - c(x_i - x_0) \\ w_i &= b(x_i - x_0) - a(y_i - y_0). \end{aligned} \tag{3.306}$$

To implement the Gauss-Newton algorithm to minimize the sum of the square distances the partial deviation needs to be obtained with the five parameters $x_0$, $y_0$, a, b, r (the five independent variables for a cylinder). These are complicated unless

$$x_0 = y_0 = a = b = 0 \quad \text{in which case} \quad r_i = \sqrt{x_i^2 + y_i^2} \tag{3.307}$$

$$\begin{aligned} \frac{\partial d_i}{\partial x_0} &= -x_i/r_i \\ \frac{\partial d_i}{\partial y_0} &= -y_i/r_i \\ \frac{\partial d_i}{\partial a} &= -x_i z_i/r_i \\ \frac{\partial d_i}{\partial b} &= -y_i z_i/r_i \\ \frac{\partial d_i}{\partial r} &= -1. \end{aligned} \tag{3.308}$$

*(b)    Algorithm operation*

1. Translate data so that the point on the axis lies at the origin:

$$(x_i, y_i, z_i) := (x_i, y_i, z_i) - (x_0, y_0, z_0).$$

2. Transform the data by a rotation matrix $\mathsf{U}$ which rotates *a, b, c* to a point on the *z* axis:

$$\begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix} := \mathsf{U} \begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix}. \qquad (3.309)$$

3. Form the right-hand side vector *d* and Jacobian according to expressions (3.305), (3.307) and (3.308).
4. Solve the linear least-squares system for $P_{x0}$ etc:

$$\mathsf{J} \begin{pmatrix} P_{x_0} \\ P_{y_0} \\ P_a \\ P_b \\ P_r \end{pmatrix} = -d. \qquad (3.310)$$

5. Update the parameter estimates to

$$\begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} := \begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} + \mathsf{U}^{\mathsf{T}} \begin{pmatrix} P \\ P_{y_0} \\ -P_{x_0} P_a - P_{y_0} P_b \end{pmatrix}$$

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} := \mathsf{U}^{\mathsf{T}} \begin{pmatrix} P_a \\ P_b \\ 1 \end{pmatrix}$$

$$r := r + P_r. \qquad (3.311)$$

These steps are repeated until the algorithm has converged. In step 1 always start with (a copy of) the original data set rather than a transformed set from the previous iteration.

If it is required to have $(x_0, y_0, z_0)$ representing the point on the line nearest to the origin, then one further step is put in:

$$\begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} := \begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} - \left( \frac{a_{x_0} + b_{y_0} + c_{z_0}}{a^2 + b^2 + c^2} \right) \begin{pmatrix} a \\ b \\ c \end{pmatrix}. \qquad (3.312)$$

(See Forbes [33] for those situations where no estimates are available.) Luckily, in surface metrology, these iterative routines are rarely needed. Also it is not yet clear what will be the proportion of workpieces in the miniature domain that will have axes of centrosymmetry. Until now, in microdynamics the roughness of rotors and stators has precluded the measurement of shape. The criterion at present is one of fitting an estimation of size and not shape. However, there is no doubt that shape will soon be a factor and then calculations such as the one above will be necessary.

## (c) Cones

These can be tackled in the same way except that there are now six independent parameters from $(x_0, y_0, z_0)$, $(a, b, c)$, $\varphi$ and $t$, where $t$ is shown in figure 3.50. $z_0$ and $c$ can be obtained dependent on the other parameters.



**Figure 3.50** Coordinate arrangement for cone.

Specify the cone, a point $x_0, y_0, z_0$ on its axis and a vector $(a, b, c)$ along its axis, and the angle $\varphi$ at the apex giving information about where on the axis the cone is to be positioned. Parameterization requires a systematic way to decide which point on the axis to choose, along with a constraint on $(a, b, c)$. For this

$$c = 1 \quad \text{and} \quad z_0 = S_0 - ax_0 - by_0 \tag{3.313}$$

for some constraint $S_0$, which is position sensitive and the choice of which has to be determined or specified with care depending on whether the cone is narrow angle or wide angle.

The distance of the point from the cone, $d_i$, is given by

$$d_i = e_i \cos(\varphi/2) + f_i \sin(\varphi/2) - t \tag{3.314}$$

where $e_i$ is the distance from $x_i, y_i, z_i$ to the line specified by $(x_0, y_0, z_0)$ and $(a, b, c)$.

Again, as for the cylinder, making $x_0 = y_0 = a = b = 0$:

$$\begin{aligned} e_i = r_i &= \sqrt{x_i^2 + y_i^2} \\ f_i &= z_i - S_0 \end{aligned} \tag{3.315}$$

and

$$\frac{\partial d_i}{\partial x_0} = -x_i \cos(\varphi/2)/r_i$$

$$\frac{\partial d_i}{\partial y_0} = -y_i \cos(\varphi/2)/r_i$$

$$\frac{\partial d_i}{\partial a} = -x_i w_i/r_i$$

$$\frac{\partial d_i}{\partial b} = -y_i w_i/r_i$$

$$\frac{\partial d_i}{\partial \varphi} = +w_i/2 \tag{3.316}$$

$$\frac{\partial d_i}{\partial t} = -1$$

with

$$w_i = (z_i - S_0)\cos(\varphi/2) - r_i \sin(\varphi/2). \tag{3.317}$$

For most cones $S_0$ is chosen such that

$$S_0 = \bar{z} - \bar{r}\tan(\varphi/2). \tag{3.318}$$

(*i*) *Algorithm description*
For cones with moderate apex angle ($< 0.9\pi$) let $S_0 = 0$:

1. Translate data so that the point on the axis lies at the origin:

$$(x_i, y_i, z_i) := (x_i, y_i, z_{i,}) - (x_0, y_0, z_{0,}).$$

2. Transform the data by a rotation matrix $\mathsf{U}$ that rotates ($a$, $b$, $c$) to a point on the $z$ axis:

$$\begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix} := \mathsf{U} \begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix}. \tag{3.319}$$

3. Form RHS vector $\boldsymbol{d}$ and $\mathsf{J}$ according to the above with $S_0 = 0$.
4. Solve the linear least-squares system

$$\mathsf{J} \begin{pmatrix} P_{x_0} \\ P_{y_0} \\ P_a \\ P_b \\ P_{\mathsf{g}} \\ P_t \end{pmatrix} = -d. \tag{3.320}$$

5. Update the parameter estimate

$$\begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} := \begin{pmatrix} x_0 \\ y_0 \\ z_0 \end{pmatrix} + \mathsf{U}^{\mathsf{T}} \begin{pmatrix} P \\ P_{y_0} \\ -P_{x_0}P_a - P_{y_0}P_b \end{pmatrix} \tag{3.321}$$

$$\begin{pmatrix} a \\ b \\ c \end{pmatrix} := \mathsf{U}^{\mathsf{T}} \begin{pmatrix} P_a \\ P_b \\ 1 \end{pmatrix} \tag{3.322}$$

$$\begin{aligned} \varphi &:= \varphi + P_\varphi \\ t &= t + P_t. \end{aligned} \tag{3.323}$$

6. This step is the same as for the cylinder.

For the special case when the periphery of the component is incomplete, the best-fit are has to be generated in order to provide a reference. This has been developed earlier [36]. It suffices here to repeat the formula. The formula used is the simplified version where the angular reference is taken to be the bisector of the arc angle $2\theta$.

Thus the centre of curvature and radius are given by

$$
\bar{x}\left(\int_{-\theta}^{\theta} r\cos\theta\,d\theta - \frac{\sin\theta}{\theta}\int_{-\theta}^{\theta} r\,d\theta\right)\bigg/\left(\theta_3 + \frac{\sin2\theta}{2}\frac{1}{\theta_3} + \frac{\cos2\theta}{\theta}\right)
$$

$$
\bar{y} = \left(\int_{-\theta}^{\theta} r\sin\theta\,d\theta\right)\frac{1}{(\theta_3 - \frac{1}{2}\sin2\theta)}
$$

$$
R = \frac{1}{2\theta}\left(\int_{-\theta}^{\theta} r\,d\theta - 2\bar{x}\sin\theta\right).
$$

(3.324)

These results have already been given in chapter 2 in the roundness section. They can obviously be extended to the 3D case. They are included here for completeness.

Equations (3.324) give the best-fit conditions for a partial arc which can enclose any amount of the full circle. Often it is necesssary to find the unique best-fit centre to a concentric pair of circular arcs. This involves minimizing the total sum of squares of the deviations. Arc 1 has sum of squares $S_1$:

$$
S_1 = \sum_{i=1}^{M}\left(r_{1i} - R_1 - \bar{x}\cos\theta_i - \bar{y}\sin\theta_i\right)^2
$$

and arc 2

(3.325)

$$
S_2 = \sum_{j=1}^{N}\left(r_{2j} - R_2 - \bar{x}\cos\theta_j - \bar{y}\sin\theta_j\right)^2.
$$

Minimizing $S_1 + S_2$ and differentiating these polar equations with respect to $\bar{x}$, $\bar{y}$, $R_1$ and $R_2$ gives

$$
\begin{pmatrix}
\sum\cos^2\theta_i + \sum\cos^2\theta_j & \sum\sin\theta_i\cos\theta_i + \sum\sin\theta_j\cos\theta_j & \sum\cos\theta_i & \sum\cos\theta_j \\
\sum\sin\theta_i\cos\theta_i + \sum\sin\theta_j\cos\theta_j & \sum\sin^2\theta_i + \sum\sin^2\theta_j & \sum\sin\theta_i & \sum\sin\theta_j \\
\sum\cos\theta_i & \sum\sin\theta_i & M & 0 \\
\sum\cos\theta_j & \sum\sin\theta_j & 0 & N
\end{pmatrix}
\begin{pmatrix}\bar{x}\\ \bar{y}\\ R_1\\ R_2\end{pmatrix}
$$

$$
=
\begin{pmatrix}
\sum r_{1i}\cos\theta_i + \sum r_{2j}\cos\theta_j \\
\sum r_{1i}\sin\theta_i + \sum r_{2j}\sin\theta_j \\
\sum r_{1i} \\
\sum r_{2j}
\end{pmatrix}.
$$

(3.326)

These equations are useful when data is available in the polar form. But when data is available in the Cartesian form the other criterion, namely minimizing the deviation from the property of the conic, is useful as described below. In this case, the equations of the arcs are written as

$$
x^2 + y^2 - ux - vy - D_1 = 0
$$
$$
x^2 + y^2 - ux - vy - D_2 = 0
$$

(3.327)

and the total sum of the squares of the deviation from the property of the arc/conic are defined as

$$
E_s = \sum(x_i^2 + y_i^2 - ux_i - vy_i - D_1)^2 + \sum(x_j^2 + y_j^2 - ux_j - vy_j - D_2)^2
$$

(3.328)

where $D$s are dummy radii. Differentiating partially with respect to $u$, $v$, $D_1$, $D_2$, the equations in matrix form to find the solution of $u$, $v$, $D_1$, $D_2$ are given by

$$
\begin{pmatrix}
\sum x_i^2 + \sum x_j^2 & \sum x_i y_i + \sum x_j y_j & \sum x_i & \sum x_j \\
\sum x_i y_i + \sum x_j y_j & \sum y_i^2 + \sum y_j^2 & \sum y_i & \sum y_j \\
\sum x_i & \sum y_i & M & 0 \\
\sum x_j & \sum y_j & 0 & N
\end{pmatrix}
\begin{pmatrix}
u \\
v \\
D_1 \\
D_2
\end{pmatrix}
=
\begin{pmatrix}
\sum (x_i^2 + y_i^2)x_i + \sum (x_j^2 + y_j^2)x_j \\
\sum (x_i^2 + y_i^2)y_i + \sum (x_j^2 + y_j^2)y_j \\
\sum x_i^2 + \sum y_i^2 \\
\sum x_j^2 + \sum y_j^2
\end{pmatrix}.
$$

(3.329)

Then

$$
\bar{x} = u/2 \quad \bar{y} = v/2
$$
$$
R_1 = \sqrt{D_1 + (u^2 + v^2)/4}
$$
$$
R_2 = \sqrt{D_2 + (u^2 + v^2)/4}.
$$

(3.330)

Obviously the key to solving these sorts of problems is how to make the equations linear enough for simple solution. This is usually done automatically by the choice of instrument used to obtain the data. The fact that a roundness instrument has been used means that the centre $a$, $b$ is not far from the axis of rotation—which allows the limaçon approximation to be valid. If a CMM had been used this would not be the case unless the centre portions were carefully arranged.

The best-fit methods above have hinged on the best-fit limaçon technique because this is the natural way in which a roundness instrument sees the signal. Should the data be obtained with the radius not suppressed it can be treated as a circle. The equation for minimization then becomes

$$
k(x^2 + y^2) + ux + vy - 1 = 0
$$

and the sum of errors $S$ is

$$
S = \sum [k(x_i^2 + y_i^2) + ux_i + vy_i - 1]^2.
$$

Differentiating partially with respect to $k$, $u$ and $v$ to minimize $S$, the sum of squared residuals gives the matrix form

$$
\begin{pmatrix}
\sum (x_i^2 + y_i^2)^2 & \sum x_i(x_i^2 + y_i^2) & \sum y_i(x_i^2 + y_i^2) \\
\sum x_i(x_i^2 + y_i^2) & \sum x_i & \sum x_i y_i \\
\sum y_i(x_i^2 + y_i^2) & \sum x_i y_i & \sum y_i
\end{pmatrix}
\begin{pmatrix}
k \\
u \\
v
\end{pmatrix}
=
\begin{pmatrix}
\sum (x_i^2 + y_i^2) \\
\sum x_i \\
\sum y_i
\end{pmatrix}.
$$

(3.331)

Then the unknowns $\bar{x}$, $\bar{y}$ and $R$ are given by

$$
\bar{x} = u/\left(2k\right) \quad \bar{y} = v/\left(2k\right) \quad R = \sqrt{(\bar{x}^2 + \bar{y}^2 + 1)/k}.
$$

An alternative least-squares approach will be given at the end of the chapter for all shapes.

### 3.7.3  Other methods

#### 3.7.3.1  Minimum zone method

The other basic method adopted in most standards is the minimum zone method. This is found by an iterative method based on the simplex method.

This method is a search routine (of which the Steifel exchange is just a subset) which is designed to climb mathematical maxima or minima. The simplex figure is obtained from the geometrical figure used in the search process. In 2D it is an equilateral triangle and in 3D it is a triangular pyramid or tetrahedron.

The basic principle is as follows:

1. Express the problem of minimization in mathematical terms:

$$R_i = [(x_i - \bar{x})^2 + (y_i - \bar{y})^2].$$

The objective or object function is to get a function, say *B,* minimized, that is

$$B = \min[\max R_1 - \min R_i].$$

It is required that *a* and *b* satisfy this.

2. Select a set of feasible values for the independent variables and use this as a starting point. This starting point is invariably chosen as the best-fit centre, that is

$$\bar{x} = 2\sum x_i/N \qquad \bar{y} = 2\sum y_i/N.$$

3. Evaluate the objective function from this centre, that is the *B* value from $\bar{x}$ and $\bar{y}$.
4. Choose, by means of the simplex figure, a second location near to $\bar{x}, \bar{y}$.
5. Compare the objective. If it is smaller move to this new point; if not try again.

Rules for minimization:

1. Reject highest point of *B*.
2. Do not return to original point.
3. Terminate after a given number of iterations.

This is one of many hill-climbing techniques but it is probably the simplest.

### 3.7.4 Minimax methods—constrained optimization

The criteria for parameters encountered in surface metrology fall into two categories: those based on least squares and those based on peaks and valleys or extrema.

Consider, for example, roundness. Typical references are the best-fit circle, the minimum circumscribing circle, the maximum inscribed circle and the minimum zone. The zonal methods, unlike the best-fit methods, rely for their basis on a limited number of high peaks and valleys rather than all of the data, as is the case with the best-fit methods. The problem in the extreme cases is, fundamentally, one of some sort of maximization or minimization, subject to a number of constraints. Thus, for example, the requirement for finding the plug gauge circle is to maximize the radius of a circle subject to the constraint that all data points lie on or outside the circle. In other words, the requirement for finding the plug gauge circle (or limaçon) is to find the maximum 'radius' for which a limaçon may be constructed such that the limaçon lies completely inside the data representing the data being measured.

A number of methods are available for finding the optimum radius or zone subject to such constraints. Originally the assessment was based on a trial-and-error estimate. For roundness measurement, this involved drawing many circles on a chart in order to get the centre which maximized, minimized or optimized the circle or zone.

### 3.7.5 Simplex methods

This technique, although crude, could give answers accurate to a few per cent providing that the workpiece was reasonably centred—in which case it is acceptable to use compasses and to draw circles on the chart.

There is a simple method which gives a good approximation to the true method. As an example of this consider the roundness profile shown in figure 3.51. Imagine that the minimum circumscribing circle (ring gauge method) is being used.

In all problems like this there are a number of ways of tackling the situation. Often there is a problem between technological relevance and mathematical elegance. In metrology the former has to take precedence. It has been pointed out many times by Lotze [37] that agreement on terminology and methodology should pre-empt any agreement on which algorithms to use for the calculation of any of the parameters in dimensional as well as in surface metrology. This is very true because as yet there is no agreement as to the procedures, let alone algorithms, for use in minimax routes (for example [38]). This has caused a lot of confusion between surface metrologists and coordinate-measuring machine metrologists. To illustrate this a simple approach will be given to show that it is easily possible to get to within a reasonable range of the true value in a search routine. This approach will be followed by demonstrating that it is one application of linear programming. The basic elements of linear programming in metrology will be reviewed together with an indication of the concept of the elegant dual linear programming methods for use in metrology. Finally, an example of an algorithmic approach will be given following Dhanish and Shunmugam [39] which could, in principle, be applied to dimensional and surface metrology, with the proviso that great care in setting up the components is taken.

Select a point in figure 3.51 at $O_1$. This can be arbitrary or it can be based on a least-squares centre [36]. Draw a circle from it with a minimum radius to touch just one point $P_1$. Move in the direction of $OP'$ at each step sweeping a radius around and reducing it until two points are at the same minimum radius. Then move the centre along the bisector again reducing the radius until a further point $P_3$ is touched. The centre position $O_3$ is then very often the exact centre frame on which the ring gauge circle is centred. This is because only three points are needed to define a circle, the centre coordinates and the radius, so three contacting points are required. The same is true for the maximum inscribed circle. Four points are needed for the minimum zone circle: two for the centre and two for the two radii. A sphere needs four, a cylinder five, etc.

The method described above is not necessarily unique, even for the circumscribing circle, which should, in principle, give a unique centre. The plug gauge method does not, neither does the minimum zone.

It should be pointed out here that in the case of roundness measurement it should be limaçons and not circles that are being described so that the evaluation should be of the least circumscribing limaçon etc.

Although the method described above works in almost all cases it is not formal. The best method for constrained optimization is to use linear programming techniques. The simple method outlined above contains the basic elements of the technique, so for completeness the basic method will be outlined below and then followed by some examples. Although the basic idea for this optimizing technique has been used in the past quite independently [40], the formal derivation has been achieved by Chetwynd [3] and his terminology will be followed here. However, other attempts have also been used. It will become obvious that the subject is quite complicated and that a number of different approaches to get a workable, cheap and fast solution are possible.



**Figure 3.51** Simple algorithm for minimum circumscribing circle.

**Figure 3.52** The 180° rule in simplex iteration.

Linear programming implies constrained optimization involving either minimization or maximizing a function (called an objective function) while keeping other relationships within predefined bounds. If these relationships are linear, then they can be expressed as a set of linear parameters and the optimization becomes linear. This so-called linear programming is fundamental to the understanding and operation of the 'exchange'-type algorithms which are now used extensively in metrology. See, for example, figure 3.52.

Take the measurement of straightness, for example. The criteria, expressed in instrument coordinates, are, given a sequence of Cartesian datum points $(x_i, y_i)$, and the minimum zone value $Z$

$$\text{minimize } Z = h \quad \text{subject to } mx_i + c + h \geqslant y_i$$
$$mx_i + c - h \leqslant y_i \tag{3.332}$$

for all $(x_i, y_i)$ simultaneously. This illustrates a convenient parameterization, namely a single line (slope $m$, intercept $c$) together with a zone of acceptability of width $2h$ set symmetrically about it. Equation (3.332) is a linear programme in $(m, c, h)$. (It is also a simple form of the minimax polynomial fit for which the so-called Steifel exchange algorithm offers an efficient solution. This may be derived from, and owes its efficiency to the properties of, the associated linear programme, figure 3.53.)

Standards present a method, originally for calculating the least-squares parameters, which has been extensively studied and is known as the 'limaçon approximation' for roundness measurement. Referring to the notation adopted, the eccentric circle is reduced, providing $e \ll R$, to

$$\rho \approx a\cos\theta + b\sin\theta + R. \tag{3.333}$$



**Figure 3.53** Steifel exchange mechanism.

This is a linearization of the parameters about the origin, which is produced mechanically by the instrument. Linearization about any other point involves considerable extra complexity of the coefficients. Whenever the limaçon approximation is valid in surface metrology the calculation of limiting reference circles becomes a linear programme. For example, the minimum circumscribing figure to a set of data points $(r_i, \theta_i)$ is expressible as

$$\text{minimize } Z = e \quad \text{subject to } a \cos \theta_i + b \sin \theta_i + R \geqslant r_i \tag{3.334}$$

for all $i$. Others may be expressed in the form of either equation (3.332) or (3.334).

Before proceeding to develop algorithms from these formulations, it is useful to establish a practical context and a mathematical notation by first illustrating the earlier work on reference circles and reviewing, extremely briefly, the main points of linear programming theory following Chetwynd [41].

### 3.7.6  Basic concepts in linear programming

### 3.7.6.1  General

A linear programme is an optimization in which the objective function (e.g. minimizing a zone or maximizing a radius) and all the constraints are linear in the parameters. Using vector notation for brevity, it can be expressed as

$$\text{maximize } Z = c^{\mathrm{T}} \boldsymbol{x} \quad \text{subject to } \boldsymbol{A} \boldsymbol{x} \leqslant \boldsymbol{b} \tag{3.335}$$

where, for $m$ positive parameters, $\boldsymbol{x}$, and $n$ constraints, $\boldsymbol{c}$ is an $m$ vector, $\boldsymbol{b}$ an $n$ vector and $\boldsymbol{A}$ an $m \times n$ matrix.

It is known (there is extensive literature on this subject) that the optimum solution occurs when each of the constraints ($c$) is satisfied to its limit by one of the parameters. Hence only certain combinations of parameter values need be examined. An orderly search through these is obtained by using the simplex method in which iterations involve only elementary row operations on the matrix-vector representation. Simplex organizes these vectors as a partitioned matrix (a tableau). This has the form

$$\begin{array}{c|c} \mathsf{K} & b \\ \hline c^{\mathrm{T}} & Z \end{array} \tag{3.336}$$

where $\mathsf{K}$ is $\mathsf{A}$ augmented by an $n \times n$ identity matrix and $c$ is correspondingly extended by $n$ zero elements. This adds $n$ 'slack variables' to the original parameters. If the $i$th parameter is limiting a particular constraint, the column $\boldsymbol{K}_i$ in $\mathsf{K}$ will have value $+1$ in the row corresponding to that constraint and zero in all other elements. The set of defining parameters so identified form the 'basis'. Initially the basis is the $n$ slack variables. Iterations attempt to match parameters to constraints in such a way that $Z$ is rapidly maximized. It is usual always to maintain the feasibility of the current iteration by ensuring that no constraint is ever violated, that is that no element of $\boldsymbol{b}'$ becomes negative. This is one of the problems not easily addressed in the simplex method. The prime indicates the vector that currently occupies the position originally occupied by $\boldsymbol{b}$. At each iteration the largest positive element of $c'^{\mathrm{T}}$ is chosen and its column brought actively into the solution (this is the strategy of 'steepest descent'). When no positive elements remain in $c'^{\mathrm{T}}$, optimality has been achieved and the solution values are readily interpreted from the tableau.

At any iteration, the columns which originally consisted of the identity matrix carry a complete and interpretable record of the row transformations carried out on the tableau. Likewise, the columns of the current basis carry the same information in the inverse of their original form. The computationally efficient method of revised simplex does not update the full tableau but merely notes what would have been done at each iteration. The work required relates to that of inverting $n \times n$ matrices. It may, therefore, be

advantageous to use a dual programme. For any $m \times n$ linear programme (termed the primal), we may define an $n \times m$ dual as

$$
\begin{array}{c|c}
\mathsf{K} & c \\
\hline
b^{\mathrm{T}} & Z
\end{array}
$$

(3.337)

where $\mathsf{K}$ is now the augmented form of $\mathsf{A}^{\mathrm{T}}$ (compare (3.336)) and the optimization has changed from minimization to maximization or vice versa. It contains exactly the same information as before, subject to the correct relative interpretation of specific elements.

### 3.7.6.2 Dual linear programmes in surface metrology

Straightness, flatness and all routine roundness measurements involve reference fittings which appear naturally as linear programmes. For more complex geometries, the errors inherent in parameter linearization may be judged acceptable when weighed against the computational efficiency of simplex. All the resulting formulations essentially have features in common indicating that the dual programme will offer the most efficient solutions.

The sign of the parameters required for simplex cannot be guaranteed with metrological data and so each parameter is replaced by a pair having equal magnitude but opposite sign. Even then the number of constraints usually dominates the number of parameters. Thus, a circumscribing limaçon fit involves six parameters and the minimum zone seven, but typical measurements involve several hundred profile points each generating a constraint; each generates two in the case of the minimum zone because it may contribute to the inner or outer circles of the zone. The sources of the difficulties encountered with early attempts at circle fittings are now apparent. They did not exploit the simplex method of searching only certain basic solutions and, furthermore, they worked with a primal or conventional formulation involving, say, six parameters and 500 constraints or points, rather than dual which, while having 500 parameters, has only six constraints. This makes the computing long-winded, so, in moving from the primal to the dual, the roles of vectors $b$ and $c$ are interchanged. If at any iteration the dual is maintained in a feasible condition (all elements of $c$ positive), the corresponding primal would be interpreted as being in an optimal, but generally infeasible, condition. The implications of dual feasibility are critical to what is to follow. Consider a physical interpretation for the case of a circumscribing limaçon (or circle). The primal feasibility condition amounts to starting with a figure which is too large but which certainly encloses the profile and then shrinking it to the smallest radius that still closes the profile. Dual feasibility would entail initially choosing a figure which is the smallest to enclose some of the data points and then expanding it as little as possible so as to include all the data—the same problem looked at a different way.

Note:

If a primal has three parameters, the dual has three constraints. The corresponding geometric observation is that a circle is defined by exactly three contacts with the data, which makes physical sense.

### 3.7.6.3 Minimum radius circumscribing limaçon

Transferring the primal, geometrical statement of the minimum radius circumscribing limaçon to the dual, the initial tableau can be written as a minimization:

$$
\begin{array}{ccccc|c}
\cos\theta_1 & \dots & \cos\theta_i & \dots & \cos\theta_n & 0 \\
\sin\theta_1 & \dots & \sin\theta_i & \dots & \sin\theta_n & 0 \\
1 & \dots & 1 & \dots & 1 & 1 \\
\hline
-r_1 & \dots & -r_i & \dots & -r_n &
\end{array}
$$

(3.338)

At any iteration giving a feasible solution, the basis will be formed from three of these columns, so taking three general contact points at $\theta_i$, $\theta_j$ and $\theta_k$ means that the basis $\beta$ is given by

$$\beta^{-1} = \begin{pmatrix} \cos\theta_i & \cos\theta_j & \cos\theta_k \\ \sin\theta & \sin\theta_j & \sin\theta \\ 1 & 1 & 1 \end{pmatrix}. \tag{3.339}$$

No significance (such as $\theta_i < \theta_j$, for example) can be read into this matrix; the relative positioning of columns depends upon the workings of revised simplex in previous iterations. The determinant of $\beta^{-1}$ is given by the sum of the cofactors of its third row, that is by the same cofactors which identify the elements of the third column of $\beta$. The non-negativity of the elements of the third column of $\beta$ thus requires that these cofactors, $\Delta_{ij}$, $\Delta_{jk}$, $\Delta_{ki}$, must have the same sign where

$$\Delta_{ij} = \begin{vmatrix} \cos\theta_i & \cos\theta_j \\ \sin\theta_i & \sin\theta_j \end{vmatrix} = \sin(\theta_j - \theta_i) \tag{3.340}$$

and similarly for the others. Using Cartesian coordinates, the cofactor can be expressed as

$$\Delta_{ij} = \frac{1}{r_i r_j} \begin{vmatrix} x_i & x_j \\ y_i & y_j \end{vmatrix} \tag{3.341}$$

and related to this is a function

$$\Delta_i = \frac{1}{r_i r} \begin{vmatrix} x_i & x \\ y_i & y \end{vmatrix} = 0 \tag{3.342}$$

which (apart from an indeterminacy at the origin, of little importance here) is a straight line passing through $(x_i, y_i)$ and $(0, 0)$ and dividing the $xy$ plane into the two areas where $\Delta_i > 0$ and where $\Delta_i < 0$. The line is also the locus of all points having $\theta_i$ as their argument. Noting the order of indices, dual feasibility requires that $\Delta_{ij}$ and $\Delta_{ik}$ have opposite sign and so lie on opposite sides of the line. An exactly similar argument applies to the other points and $\Delta_j = 0$ or $\Delta_k = 0$. If point $k$ is to lie on the opposite side of $\Delta_{jr} = 0$ from point $j$ and on the opposite side of $\Delta_{jr} = 0$ from point $i$, it can only occupy the sector shown in figure 3.52. As it is only in this geometry that $\Delta_{ik}$ and $\Delta_{jk}$ will have opposite signs, as required for dual feasibility, the following theorem, termed here 'the 180° rule', is proved. (A point known for a long time but proved by Chetwynd.)

*The 180° rule*

A circumscribing limaçon on a given origin to a set of points is the minimim radius circumscribing limaçon to those points if it is in contact with three of them such that no two adjacent contact points subtend an angle at the origin of more than 180°, where the term 'adjacent' implies that the angle to be measured is that of the sector not including the third contact point.

A complete simplex iteration for the minimum radius circumscribing limaçon in the dual consists of selecting any point which violates the reference (conventionally, the point giving the largest violation is chosen) and substituting it for one of the points defining the reference in such a way that dual feasibility is maintained. The 180° rule allows the general iteration to be simplified to the following exchange algorithm:

1. Choose any three data points such that no two adjacent ones subtend an angle at the origin of more than 180°.

2. Construct a reference limaçon through these three points.
3. If no data points lie outside this limaçon the solution is found, otherwise choose the point which violates the reference by the largest amount.
4. Replace one of the reference points by this new point such that the 180° rule is still obeyed and go back to step 2.

The exchange between any new point and the contacts is always unique.

An exchange algorithm depends upon the iterations moving monotonically towards an optimum solution in order to guarantee that cyclical exchanges do not occur. Here this is the case, because, as the exchange is unique at each iteration, it must be identical to the variable change at the simplex iteration of the linear programme and that is known to converge monotonically.

A similar procedure needs to be followed to get the formal linear programming matrix for the minimum zone but essentially the result is as follows.

The conditions for the optimum solution to the minimum radial zone limaçons give rise to the following geometric interpretation. Expressing the zone as a band of width $2h$ placed symmetrically about a single limaçon:

1. All data points must lie not more than a distance $h$, measured radially from the origin, from the limaçon.
2. There must be four data points all lying exactly $h$ from the limaçon such that they lie with increasing angle $\theta_i$ alternately inside and outside the limaçon.

It may be noted that this alternation property is not unique to this problem; it occurs, for instance, in the Steifel exchange algorithm for best-fit polynomials, which may also be derived by a linear programming method [42]. These rules may be used to formulate an exchange algorithm. Thus:

1. Choose arbitrarily four data points.
2. Fit to these a reference limaçon such that they are radially equidistant from it and lie alternately on either side of it with increasing angle.
3. If no other points are further from the reference the solution is found.
4. Otherwise substitute the point which lies furthest from the reference for one of the four defining points such that the new set of points lie alternately on either side of the reference and return to step 2.

### 3.7.6.4  Minimum zone, straight lines and planes

The minimum separation parallel and straight lines belong to the well-documented class of minimax polynomials, that is curves having the smallest possible maximum divergence from the data. The condition for this to occur is that, relative to an $n$th-order polynomial, the data must have $n + 2$ maxima and minima, all of equal magnitude. The solution can be found by the Steifel exchange algorithm, which proceeds by fitting the polynomial according to this condition to $n + 2$ points and then exchanging points further away from it with the defining set, while maintaining the condition. In terms of the minimum zone straight lines there will be three points, two contacting one line and one the other in an alternate sequence, which are iterated by exchanges (see figure 3.53).

The minimum zone planes can be expressed, in instrument coordinates, as

$$\text{minimize } Z - h \quad \text{subject to } ax_i + by_i + c + h \geqslant z_i$$
$$ax_i + by_i + c - h \leqslant z_i \tag{3.343}$$

for all data points $(x_i, y_i, z_i)$. $a$, $b$ and $c$ are sign unrestricted and $h \leq 0$. Noting that $h = 0$ is valid only for the trivial condition that all points are coplanar, then it may be asserted that four points will be represented in the basis of the dual, which can be expressed as

$$\beta^{-1} = \begin{pmatrix} S_i x_i & S_j x_j & S_k x_k & S_l x_l \\ S_i y_i & S_j y_j & S_k y_k & S_l y_l \\ S_i & S_j & S_k & S_l \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

(3.344)

where $S_i$ etc take values of $+1$ or $-1$ according to whether $(x_i, y_i, z_i)$ contacts the upper or lower of the minimum zone planes. As before, dual feasibility is guaranteed if all terms in the final column of $\beta$ are positive, which will be true providing that the cofactors of the final row of $\beta^{-1}$ all have the same sign. These cofactors are

$$-S_j S_k S_l \Delta_{jkl} \qquad -S_i S_j S_l \Delta_{ijl} \qquad S_i S_k S_l \Delta_{ikl} \qquad S_i S_j S_k \Delta_{ijk}.$$

(3.345)

They must all have the same sign. Consider the determinant equation representing the boundary between positive and negative regions of $\Delta_{jkl}$ :

$$\Delta_{jk} = \begin{vmatrix} x_j & x_k & x \\ y_j & y_k & y \\ 1 & 1 & 1 \end{vmatrix} = 0.$$

(3.346)

It is a plane parallel to the $z$ axis (since it is independent of $z$), passing through points $(z_l, y_l)$ and $(x_k, y_k)$. Dual feasibility requires that if $S_i = S_l$ (contacts with the same place) $\Delta_{jkl}$ and $\Delta_{jki}$ must have different signs and *vice versa*. So if the $i$th and $l$th contacts are with the same plane they lie on opposite sides of $\Delta_{jkr} = 0$, but if they contact different planes they both lie on the same side of $\Delta_{jkr} = 0$. A parallel argument shows that the same is true for all pairs of points.

These relationships show the relative positions of contacts that give dual feasibility (figure 3.54). There can be two contacts with each of the minimum zone planes, in which case the plan of lines joining the alternate types must form a convex quadrilateral or a 3:1 split, in which case the single contact must lie in the plan of the triangle formed by the other three contacts.



**Figure 3.54** Contacts which give dual feasibility.

Even with this most simple of three-dimensional zone fits, the advantage of using specific exchange algorithms rather than a general revised simplex solution in an automatic system is becoming questionable—hence the difficulty of standardization.

Note:

It may, at first glance, seem surprising that limaçon fitting rather than the apparently simpler case of flat surfaces has been used as the primary example. This section, being primarily concerned with linear programming, does not report detailed comparisons between algorithms but some comments are needed. The following remarks are based on practical tests made by Chetwynd [41].

A typical roundness 'profile' would have 512 equally spaced radial ordinates each resolved over a 10 or 12 bit working range. Exchange algorithm systems have now been working with data of this type in both industrial and research environments for several years and their robustness has been established. Even with an arbitrary choice of points for the initial basis, the exchange algorithm virtually always solves for the minimum circumscribing limaçon in five or less iterations, while the minimum zone only occasionally needs more than five on real engineering profiles. The earlier (primal-based) algorithms were run with poorly defined end conditions, typically making 32 relatively coarse-stepped iterations and then 32 finer steps, after which the process was terminated with a result assumed to be close to the desired optimum. The dual techniques yield at least a 10:1 saving in the number of iterations as well as giving a fully determined convergence and so better accuracy. With both algorithms the iteration is dominated by the almost identical computation and checking of the updated figure, so the program size and the cycle times are closely similar on similar machines programmed in the same language. A 10:1 speed increase is also obtained.

The direct use of revised simplex on dual programmes representing limaçon fitting has been studied using a specially developed package containing only the subroutines essential for solving this class of problem. Memory requirements are only slightly larger than those of exchange algorithms and execution is typically about 20% slower. This is due to the simple way artificial variables are treated. This difference can be removed at the cost of extra programme length.

The limaçon fits to roundness measurement have simple exchange rules that can be expressed in a few numerical comparisons and logic operations. Thus in a specialized system both a size reduction and a speed increase would be obtained by replacing the direct use of revised simplex on the dual by an exchange algorithm. However, the exchange logic is specific, so if several different references are to be implemented there will be less shared code. With more complex geometries it is of even greater importance that the efficiency of dual-based methods is obtained. Yet, even with the simplest three-dimensional geometry, the exchange rules are becoming quite complicated. Using a conclusion made by Chetwynd, duality theory has shown that the 'obvious' geometrical method is rarely the best approach.

The study of exchange algorithms gives a very clear insight into the geometrical implications of reference fitting. This is important for metrology. Measurement should always be based on engineering relevance rather than a mathematically convenient abstraction. The exchange algorithm also provides a good method for solution by hand should it be necessary. Relatively flexible measurement systems are likely to use a more general implementation of a revised simplex algorithm. This is no cause for concern: both are firmly based on the same theoretical foundation.

Note:

Even these considerations should show that the use of dual and exchange methods requires a considerable knowledge of the subject so that the parameters to be measured and the constraints to be applied can be formulated properly. The gains are massive if the mathematical background is sound, but it should never be forgotten that these methods rely on a relatively small number of geometrical points to satisfy the ultimate optimum solution. This is a weakness because of the fleeting reliability of such extrema. If this is a consideration in one dimension such as in roundness and straightness, it is even more so in the two-dimensional cases of sphericity, cylindricity, etc. In these cases, especially that of cylindricity, the whole problem becomes very complicated because the straightforward requirement of five constraints becomes masked in the data produced by a surface metrology instrument due to tilt and eccentricity effects. For a more detailed discussion on this subject refer to reference [43].

### 3.7.6.5 Minimax problems

#### (a) Example of a general algorithmic approach

Because there is no absolutely correct algorithm or even an agreed algorithm for the minimax problems in surface and dimensional metrology it is informative to give an example here. The algorithm chosen is to

determine the minimax zone value of a number of different forms commonly found in industry, ranging from the flat to the cylindrical in parallel with the best-fit least-squares procedures given earlier.

The reason why the minimum zone has been selected is because, next to least squares, it is often the preferred way of specifying the tolerance on a form. In fact the ISO standards refer to it often but do not attempt to indicate how to obtain it! This is not surprising because nobody has yet agreed how to do it. The optimum obviously depends on a number of things, not the least being simplicity and technological relevance rather than mathematical elegance.

Such a straightforward approach has been suggested by Dhanish and Shunmugam [39]. This approach assumes that the form parameters have already been linearized, which is consistent with the constraints of most surface metrology instruments and potentially with coordinate-measuring machines providing that the part is very well aligned. Problems such as tilt, eccentriticy and direction of measurement have been ignored. The basic problem with tilt and eccentricity has only really been attempted properly by Chetwynd.

### (b) Definition of the problem

Let the function of the surface be $b$, and $\varphi$ the function which best fits the surface. Assuming that the form shape has been linearized in surface metrology

$$\varphi_i = a_{i1}u_1 + \ldots + a_{ij}u_j + \ldots + a_{in}u_n. \tag{3.347}$$

$n$ is the number of variables, $a_{ij}$ denotes the value of the $j$th variable at the $i$th place, and $u_j$ is the coefficient of the $j$th variable.

As in least squares the deviation $e_i = \varphi_i - b_i$ is

$$e_l = a_{i1}u_1 + \ldots + a_{in} - b_i. \tag{3.348}$$

The form error is computed as

$$h_i = \left| e_{max} \right| + \left| e_{min} \right| \tag{3.349}$$

where $e_{max}$ and $e_{min}$ are the maximum and minimum errors.

The general problem for the minimum zone is, given $a_{ij}$ and $b_i$ find $u_j$ of $\varphi$ such that $\max|e_i|$ is a minimum (following Dhanish and Shunmugam's method below). The steps are as follows:

1. Choose a set of $m + 1$ points as a reference set. It is usual to use a set based on the least-squares approximation, usually the points giving maximum deviation from the least squares. Call these points $v_{kj}$ so that $v_{hj} = a_{ij}$; $i$ takes the values in the reference set. The steps from 2 onwards operate only on the reference set.

2. Find $\lambda_k$, $k = 1 \ldots m + 1$, solutions of the equations

$$\sum_{k=1}^{m+1} \lambda_k v_{ij} = 0. \tag{3.350}$$

Since these are $m$ homogeneous equations with $m + 1$ unknowns the solution in determinant form is

$$\lambda_k = (-1)^{k+1} \left| v_{ij} \right| \qquad j = 1 \ldots m, \qquad i = 1 \ldots m + 1(i \neq k). \tag{3.351}$$

3. Calculate the reference deviation $d$ where

$$d = \sum_{k=1}^{m+1} \lambda_k b_k \left/ \left( -\sum_{k=1}^{m+1} \left| \lambda_k \right| \right) \right. . \tag{3.352}$$

If the denominator is zero the reference is degenerate, so choose a new reference. This is usually done by rotating the points to include the next point in order and dropping the first.

4. Find the values of $\varphi_k$ for this reference set:

$$\varphi_k = b_k + \mathrm{sgn}(\lambda_k)d \qquad k = 1 \ldots m+1. \tag{3.353}$$

If $\lambda_k = 0$ then Haar's condition is violated and the situation is ambiguous with respect to sign, as mentioned earlier, so both the positive and negative $d$ have to be tried. For both $+$ and $-$ values continue to step 6 with that sign which gives the lowest value of error. Also this sign shows the direction to go.

5. Find the levelled reference by interpolation:

$$\varphi_k = \sum_{j=1}^{m} u_j v_{kj} \qquad k = 1 \ldots m+1 \tag{3.354}$$

so $m+1$ equations exist for $m$ coefficients $u_1 u_2 \ldots$ so that any $m$ equations will suffice.

6. Find the value $d_i$ of this reference function at all the given points. The value of $d_i$ whose absolute value is the highest is referred to as the supreme error $e^*$ if

$$| e^* | \le d. \tag{3.355}$$

Then the criterion of optimality is satisfied and the search is stopped. Then the function $\phi$ is the Chebychev approximation sought and the form error is calculated based on the equation (3.349) for $h_i$.

Otherwise the point $i^*$ corresponding to the supreme error $e^*$ enters the reference set and the search is continued.

7. Find the point to be discarded from the reference set by the Steifel exchange mechanism used by Chetwynd. Solve for $\mu_k$ in the equations

$$\sum_{k=1}^{m+1} \mu_k v_{kj} + a_{ij} = 0 \qquad j = 1 \ldots m. \tag{3.356}$$

Since there are $m+1$ variables and only $m$ equations a suitable value for one of the $\mu$ could be assumed; then calculate $q_k$ where

$$q_k = \mu_k / \lambda_k \qquad k = 1 \ldots m+1. \tag{3.357}$$

If $\mathrm{sgn}(d)e^* > 0$ the point with minimum $q_k$ should be removed from the reference set, otherwise the point with maximum $q_k$ is discarded. (If any $\lambda_k = 0$ a 'trial exchange' has to be tried. Each point in the reference set is replaced by the point of supreme error one at a time and the trial reference deviation is computed for each set. The combination leading to the largest value of $|d|$ is taken as the next reference step.) With this reference set go back to step 2. Note that the trial reference is not the exchange therein but alternative to it in the special case $\lambda_k = 0$ where the exchange will not work.

Any example taken from reference [39] shows how this algorithm works. There are worked examples given for flatness, roundness, cylindricity and sphericity.

However, even though the parameterization has been linearized there is still quite a number of special situations which have to be taken into account for obtaining a general algorithm. This detail is much greater when the system has to be linearized by iteration and when the effects of misalignment and eccentricity are taken into account. It seems that the more versatile and wide ranging the instrumentation becomes, the more difficult it will be to finalize a simple general algorithm. It may be that the best that will be done is to agree on a procedure whereby the algorithms will be able to be well behaved. For the fine detail of algorithms for working out minimax problems see books on computational geometry (e.g. [44]).

### 3.8 Transformations in surface metrology

#### 3.8.1 General

The application of random process analysis has had a profound influence on the understanding of surface texture from the point of view of generation and of function. Random process analysis is concerned with the Fourier transform and its inverse. However, the success of the Fourier transform has been so great that it has brought on an appetite for more and better transforms. Do they exist?

Are there transforms which have simpler or faster properties? Does another transform tell more about the surface?

These questions have been raised many times in recent years and have produced a plethora of transforms in the literature almost equivalent to the 'parameter rash' reported earlier [45].

Whilst there is an undenied excitement in introducing a new transform into the field, care has to be exercised that it does not produce an over-reaction. What usually happens is that a transform may give a benefit in one area, but not overall.

The transforms introduced fall into two main categories:

1. faster and simpler ways of producing comparable results;
2. transforms giving more information.

In the former class there are the Walsh, Hartley, Hadamard, BIFORE and Haar transforms (all orthogonal transforms) and in the latter case the Wigner, ambiguity, wavelet and Gabor transforms (these being space-frequency transforms).

Another factor that enters into the question of the use of transforms different from the Fourier is that the old criteria of calculation are no longer as pressing as they once were, for example storage size and speed of implementation. Nowadays storage is no problem and the speed of modern, even small, computers is such that real-time operations are possible, so the benefit obtained by using other transforms has to be significant to warrant their use. Also there is the point that the Fourier transform is very well known and understood and is standard in many existing instruments. To displace it would pose many educational and usage problems.

The basic question in surface metrology is this: how relevant to surface behaviour are new transforms that have usually been devised because of new demands in the subjects of communications and coding? Fundamentally the new transforms are devised for temporal properties. They all need to be transformed in one way or another to be useful in surface metrology. A case in point could be the curvature of an areal summit. It is hard to see how this could be important in communications yet it is very important in tribology. Unless or until transforms are developed which deal with the vertical (or top-down) properties of spatial features directly and not via the Fourier transform and the multinormal distribution, it seems that it is probably best to stick to the Fourier approach. However, because of their existence in the literature, mention of them will be made here. Some transforms are very close to the Fourier transform, for example the Hartley transform [46].

#### 3.8.2 Hartley transform

The Hartley transform $H(u, v)$ is the difference of the real and imaginary parts of the Fourier transform $F(u, v)$. Thus if

$$F(u,v) = F_{\text{real}}(u,v) + F_{\text{imag}}(u,v) \tag{3.358}$$

then

$$H(u,v) = F_{\text{real}}(u,v) - F_{\text{imag}}(u,v)$$
$$jF_{\text{imag}}(u,v) = \tfrac{1}{2}[F(u,v) - F(-u,-v)]. \tag{3.359}$$

A coherently illuminated object in the front focal plane of an optical system produces the Fourier transform in the back focal plane. Unfortunately this is not the case in the Hartley transform, but there are some features of the Fourier transform that are inconvenient. For example, the phase of the field contains substantial information but available optical sensing elements respond to intensity and not to phase. A photographic record of a Fourier transform plane abandons phase information which can be very important. On the other hand, the direct recording of the squared modulus of the Hartley transform would provide much more information. In cases where the transform does not go negative $|H(u,v)|^2$ suffices to recover $f(x, y)$ in full. In other cases, where the Hartley transform does go negative, knowledge of $|H(u,v)|^2$ does not by itself determine the sign; however, sign ambiguity is a much less serious defect than the complete absence of phase knowledge and sign knowledge can often be inferred

The Hartley transform does not inherently use complex notation:

$$H(u) = \int_0^L f(x)\mathrm{cas}(kux)\mathrm{d}x \tag{3.360}$$

where the symbol 'cas' means cosine and sine:

$$\mathrm{cas}\,\theta = \cos\,\theta + \sin\,\theta. \tag{3.361}$$

The discrete one-dimensional Hartley transform offers speed gains over the FFT for numerical spectral analysis and therefore it has great potential in communications, but there has been difficulty in carrying the advantage over to more than one dimension. This is because, whereas $\exp[-jk(ux + vy)]$ is easily separable into two factors, the function $\mathrm{cas}[k(ux + vy)]$ is not. Until this becomes straightforward it seems that the Hartley transform will be of use only in profile evaluation [35].

### 3.8.3 Square wave functions—Walsh functions

The most popular of the square wave functions is the Walsh function. It is similar to the Fourier transform except that the sinusoidal function is replaced by an on-off signal, that is a pseudo-square-wave signal.

Because it is not a sinusoidal form the frequency axis has been renamed 'sequency'. This is a measure of 'half the number of zero crossings per unit time or length', so any signal can be represented by a Walsh spectrum as opposed to a Fourier spectrum. Thus

$$\begin{aligned} f(\theta) &= c_0 + c_1\,\mathrm{Wal}(1,\theta) + c_2\,\mathrm{Wal}(2,\theta) \\ &= a_0 + a_1\,\mathrm{Cal}(1,\theta) + a_2\,\mathrm{Cal}(2,\theta) + b_1\,\mathrm{Sal}(1,\theta) + b_2\,\mathrm{Sal}(2,\theta) \end{aligned} \tag{3.362}$$

where Cal = cos Walsh and Sal = sin Walsh.

The Walsh functions, originally defined in 1923, take only the values $+1, 0, -1$ and consequently they are very fast to compute.

If the Walsh function is factorized using the Cooley-Tukey FFT routine then the speed advantage over the FFT is considerable because the usual routine of complex multiplication is replaced by one of multiplying simply by $\pm 1$ (similar to the Stieltjes technique for correlation; see section 3.10.1).

Which of the two methods, Fourier or Walsh, is best? The answer is that it depends on the use. Common-sense metrology would say, if the surface geometry is angular or discontinuous in any way, use the Walsh; if it is continuous and smooth use the Fourier, the argument being that the metrology should follow the function wherever possible. This is illustrated by the fact that a square wave has a Fourier transform of an infinite number of sinusoidal components, and one Walsh component. A sine wave has infinite Walsh coefficients, yet only one Fourier coefficient [47].

It is also beneficial that the surface can be reconstituted from the Walsh spectrum albeit with many coefficients.

The whole of the Walsh philosophy is discontinuous rather than continuous and as such its operation and properties have been examined by Smith and Walmsley [47] using dyatic calculus. In this Nayak's results have been shown to have an equivalent derived via the Walsh spectrum rather than the Fourier. Because the Walsh spectra $P_w$ are most suited to discrete operation they appear in the literature in this form:

$$F_w = (\theta) - \frac{1}{N} \sum_{n=0}^{N-1} f(n) \, \text{Wal}(n\theta) \qquad \theta = 0, 1, 2, ..., N-1$$

and the Walsh power spectrum by

$$P_w(0) = F_w^2(0)$$
$$P_w(\theta) = F_w^2(0)(2\theta - 1) + F_w^2(2\theta)... \qquad 0 = 1, ..., N/2 - 1$$
$$P_w N/2 = F_w^2(N - 1)$$

$$(3.363)$$

Unlike the power spectra of the Fourier transform the Walsh transform is not invariant to circular time shifts. This has led to the development of phase-invariant square wave transforms, in particular the Walsh phase-shift-invariant transform which is obtained from the autocorrelation function by means of a series of translation matrices. This has the effect of summing and averaging all the possible Walsh transforms of the time-shifted versions of the data—not very elegant but effective!

Another possibility is the BIFORE transformation $P_b$, which is obtained from the Walsh transforms (BIFORE = binary Fourier representation):

$$P_b(0) = P_w(0)$$
$$P_b(\theta) = \frac{N}{2^{(\theta+1)}} \sum_{N=1}^{N/2} P_w(2^\theta n - 2^{(\theta-1)}) \qquad 0 = 1, 2, ..., \log N.$$

$$(364)$$

This has only $1 + \log_2 N$ spectral values spaced logarithmically in sequence.

Mulvaney [48] compared some of these transforms and came to the general conclusion that, after all, the Fourier transform is best for surface characterization. Although the orthogonal binary transforms such as the Walsh, Haar, etc, were faster they did not give spectral estimates which converged rapidly as a function of record length. The phase-invariant ones, Fourier and BIFORE, were slower to evaluate yet converged quickly.

A fundamental problem keeps on emerging when considering the use of binary-type transforms for use in surface analysis, and this is the difficulty in identifying strongly periodic continuous characteristics such as may well result from a surface analysis. Transforms like the Walsh do not respond well to such features; other ones are the Haar. Consequently many features that reflect machine tool misbehaviour, such as chatter, would be missed at worst or degenerated at best. This seriously limits the use of such transforms. Other orthogonal transforms such as the Hadamard, like the Hartley, have been used in optics and in general signal processing. The Hadamard, which is another binary-type transform based on the Walsh functions, gives an order of magnitude increase in speed over the Fourier transform yet inherits the difficulties of application of the Walsh for basically continuous functions. In the past few years another variant of the Walsh-like transforms has emerged. This is called the wavelet transform [49]. This, like the others above, is binary and is especially useful for the compression of data containing many edges, such as fingerprints and in some cases TV images. It is not essentially suitable for surface analysis, which is basically continuous, but it may be that with the growing importance of fractal-like characteristics the wavelet transform may find a place.

It should be noted that these classes of transform based on the binary signal were originally devised to reduce the bandwidth of a signal being transmitted from one place to another, so they can hardly be expected to provide substantial advantages in their application to surface metrology.

### 3.8.4 Space—frequency functions

A great deal of thought and effort has gone into the mathematical tools needed to characterize the surface roughness. There is also little doubt that so far the most important of these is probably that of random process analysis in terms of spectral analysis and the correlation functions. Random process analysis is a very powerful tool for extracting the most significant information from a very noisy background. So far it has been shown that it can be used for characterizing both deterministic and random surfaces. It will be shown in chapters 6 and 7 how it can be used with great effect to give information on manufacture and performance. Some examples of this have already been hinted at.

However, there is scope for improvement. Random process analysis, as used here, is best fitted to examine stationary random surfaces and their statistical properties. It tends to bring out average statistics. But in many cases the very thing that causes failure of a component is the non-uniformity of a part, not its uniformity—corrosion patches and fatigue cracks are just two functional cases. Tool wear is one example in manufacture, Neither correlation nor spectra are well placed to look for non-stationarity because they are basically simple averages in space or frequency and, as a result, tend to integrate out changes from position to position or frequency to frequency. Ideally a function should be used that has the equal capability of characterizing random and periodic signals and their non-stationary embodiments. This implies a function which has two arguments and is related to the Fourier kernel. There is a class of functions which is capable of doing just this and they are called space-frequency functions. There are two functions which have possibilities: one is called the Wigner distribution, which derives from quantum mechanics [50] and the other is the ambiguity function used in radar [51]. Both have been mentioned earlier as an adjunct to the Fourier transform. They will be recapitulated here in context.

For simplicity, these functions will be stated for profile data, although they have areal equivalents. The ambiguity function is given by $A(\chi, \overline{\omega})$ where

$$A(\chi,\overline{\omega})=\int_{-\infty}^{\infty} z(x-\chi/2)z^*(x+\chi/2)\exp(-\mathrm{j}\overline{\omega}x)\,\mathrm{d}x \tag{3.365}$$

and the Wigner distribution $W(x, \omega)$ by

$$W(\chi,\omega)=\int_{-\infty}^{\infty} z(x-\chi/2)z*(x+\chi/2)\exp(-\mathrm{j}\omega\chi)\mathrm{d}\chi. \tag{3.366}$$

It is clear that they are related:

$$A(\chi,\overline{\omega})=\frac{1}{2\pi}\int\int_{-\infty}^{\infty} W(x,\omega)\exp[-\mathrm{j}(\overline{\omega}\chi-\omega\chi)]\mathrm{d}x\,\mathrm{d}\omega \tag{3.367}$$

$$W(\chi,\omega)=\frac{1}{2\pi}\int\int_{-\infty}^{\infty} A(\chi,\overline{\omega})\exp[-\mathrm{j}(\omega\chi-\overline{\omega}\chi)]\mathrm{d}\overline{\omega}\,\mathrm{d}\chi. \tag{3.368}$$

Both the ambiguity function and the Wigner function are space-frequency functions occupying neither space nor frequency but both. They can be thought of as being in between the two. This makes them suitable for easy access to either domain.

The basic difference between the two, which has been demonstrated earlier, is that the ambiguity function is a correlation (because it integrates over the variable $x$ leaving shift $\chi$ constant) and the Wigner distribution is a convolution (because it retains $x$ and integrates over $\chi$ which is in effect the dummy argument in the convolution). It is this latter property—the retention of the value of $x$, the position in space or position in frequency—which makes the Wigner distribution the most useful of the two in engineering, although the ambiguity function has uses in optics.

In the same way that it makes physical sense only to consider signals that start at $t = 0$ and continue therefrom (i.e. the signals are causal), it makes more physical sense to consider signals which are analytic (i.e. exist for frequencies greater than or equal to zero). To render signals analytic, it is necessary to use the Hilbert transform thus: if

$$z_a(x) = z(x) + j\hat{z}(x) \tag{3.369}$$

and $f(x)$ is the Hilbert transform, thus

$$\hat{z}(x) = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{z(\chi)}{x - \chi} \, d\chi. \tag{3.370}$$

Rendering the signal analytic reduces the risk of aliasing because it reduces the artifact of frequencies less than zero. It is possible that both the Wigner and the ambiguity functions have many useful properties in engineering (see [52,53]). In this context it will suffice to derive the digital form and the moments of the Wigner function.

For practical applications the digital form of the Wigner distribution is used. For this the signal is $f(n)$ where $n = -2, -1, 0, 1, 2, ..$ and $W(x, \omega)$ becomes the discrete Wigner distribution (DWD) $W(n, \theta)$, where

$$W(n, \theta) = \sum_{\zeta = -\infty}^{\infty} 2z[n + \zeta] z^*[n - \zeta] \exp(-2j\theta\zeta) \tag{3.371}$$

or equivalently

$$W(n, \theta) = \frac{1}{n} \int F(\theta + v) F^*(\theta - v) \exp(2jvn) \, dv$$

where $F(v)$ is the discrete Fourier transform (DFT) of $z[n]$. The digital form of the moments become zeroth order

$$\text{zeroth order} \quad p[n] = \frac{1}{2\pi} \int_{-\infty}^{\infty} W[n, \theta] \, d\theta$$
$$= |z(n)|^2 . \tag{3.372}$$

The first-order moment $\Gamma(n)$ is

$$\Gamma[n] = \tfrac{1}{2} \arg\left( \int_{-\pi/2}^{\pi/2} \exp(2j\theta) W[n, \theta] \, d\theta \right)$$
$$= \tfrac{1}{2} \arg\{ z(n+1) z^*(n-1) \}. \tag{3.373}$$

This is the instantaneous frequency; for example, for $z(n) = v(n) \exp(j\varphi(n))$

$$\Gamma[n] = \frac{\varphi(n+1)z - (n-1)}{2} \bmod \pi = \varphi'(n) \tag{3.374}$$

and the second-order moment is

$$m(n) = \frac{p(n) - \left| (1/2\pi) \int_{-\pi 2}^{\pi 2} \exp(2j\theta) W(n, \theta) \, d\theta \right|}{p(n) + \left| (1/2\pi) \int_{-\pi 2}^{\pi 2} \exp(2j\theta) W(n, \theta) \, d\theta \right|} \tag{3.375}$$

or

$$m[n] = \frac{|f(n)|^2 - |f(n+1)f^*(n-1)|}{|f(n)|^2 + |f(n+1)f^*(n-1)|}. \tag{3.376}$$

Because the Wigner function has been applied usually to roundness the data points are restricted to $(0 \rightarrow N-1)$ where $N$ is a binary number (e.g. 256):

$$W(n, \theta) = \sum_{m=-\infty}^{\infty} 2z(n+m)z^*(n-m)\exp(2\mathrm{j}\theta m).$$

Letting $M = N/2 - 1$,

$$W(n, \theta) = \sum_{m=-\infty}^{N-1} 2z(n+m-N)z^*(n-m+M)\exp[-2\mathrm{j}(m-M)\theta].$$

Letting $\theta = k\pi/N$,

$$W\left(n, \frac{k\pi}{N}\right) = \frac{1}{\sqrt{N}} \sum_{m=-\infty}^{N-1} z(n+m-M)z^*(n-m+M)\exp\left(-\frac{2\mathrm{j}mk\pi}{N}\sqrt{\exp\left(\frac{2\mathrm{j}Mk\pi}{N}\right)}\right).$$

(3.377)

Therefore, to compute the WD it is only necessary to compute:

1. $2z[n+m-M]\,z^*[n-m+M]$ where $n, m = 0, 1, 2, \ldots, N-1$.
2. $(1/(\sqrt{N})\Sigma_{m=0}^{N-1}(2z[n+m-M]z^*[n-m+M])\exp(-\mathrm{j}2mk\pi/N)$ where $n, k = 0, 1, 2, \ldots, N-1$.
3. $W_f(n, k\pi/N) = ((1/\sqrt{N})\Sigma_{m=0}^{N-1}(2z[n+m-M]z^*[n-m+M]))\,\sqrt{N}\,\exp(\mathrm{j}2mk\pi/N)$ for $n, k = 0, 1, 2, \ldots, N-1$.

Equation (3.382) represents the WD computed for $n, k = 0, 1, 2, \ldots, N-1$. The local moments in frequency become

$$p(n) = \frac{1}{2N} \sum_{k=0}^{N-1} W\left(n, \frac{k\pi}{N}\right)$$

(3.378)

$$\Gamma(n) = \frac{1}{2} \arg\left[\sum_{k=0}^{N-1} W\left(n, \frac{k\pi}{N}\right)\exp\left(\frac{2\mathrm{j}k\pi}{N}\right)\right].$$

(3.379)

How these moments can be applied is shown in an example in , in the detection of waviness in the form of chatter.

### 3.8.5 Gabor transforms

The key to the Gabor transform is its discrete form. It is defined as

$$z(i) = \sum_{m=0}^{\infty} \sum_{n=0}^{N-1} C_{m,n} h_{m,n}(i)$$

(3.380)

where

$$C_{mn} = \sum_{i=0}^{\infty} z(i)\gamma_{m,n}^*(i)$$

(3.381)

where $z(i)$ is the discrete space or surface signal, $C_{m,n}$ is a Gabor coefficient, $h_{mn}(i)$ is the basis Gaussian function used in Gabor transformations, $\gamma_{mn},(i)$ is called the biorthogonal auxiliary function, $\gamma^*$ is the conjugate and where

$$h_{mn}(i) = h(i - m\Delta M)W_L^{n\Delta Ni} \qquad \gamma_{m,n}(i) = \gamma(i - m\Delta M)W_L^{n\Delta Ni}$$

(3.382)

$\Delta M$, $\Delta N$ are sample intervals in the space and the frequency domain.

In general a Gabor transform is not unique. Depending on the selections of $h(i)$ and $\gamma(i)$ other transform pairs are possible.

The orthogonal-like Gabor transform maps a space (surface) signal which is composed of Gaussian functions separated in space into a joint space-frequency function with complex coefficients. These coefficients represent the amplitude of a space-shifted frequency-modulated Gaussian function centred at $x_1 f_1$, $x_2 f_2$ etc for this space waveform. The Gabor coefficients can be used to reconstruct the original space series or to calculate power spectral density using the pseudo Wigner-Ville distribution. This is given by

$$\text{Wigner} - \text{Ville}(i,k) = 2\exp\left[-\left(\frac{i - m\Delta M}{\sigma^2}\right)^2 + \left(-\frac{2\pi\sigma}{L}\right)^2 (k - n\Delta N)\right]. \tag{3.383}$$

This equation can be precomputed and stored in a look-up table [54].

After the Gabor coefficients are obtained the Wigner-Ville distribution is applied to each of the time-shifted and frequency-modulated Gaussian functions.

One of the most promising applications of the Gabor transformation is to non-stationary surface signals, in much the same way as the Wigner function. However, it may be that it is better because it is virtually optimized in both space and frequency.

There is another possible benefit, which relates back to the discussions on waviness and filtering, and this concerns the Gaussian weighting function used to separate waviness. It may be that because the Gabor transform is also Gaussian, some practical benefit may be obtained.

So far the Gabor transform has not been used in surface analysis but it seems an obvious candidate for use.

The wavelet transform has been considered earlier together with the Wigner distribution and ambiguity function.

The wavelet function is arbitrary but in order to be usable as a general operator the various scales of size have to be orthogonal—much as a sine wave, and its harmonics are in the Fourier Series.

Wavelets can be very different in shape. Usually the actual shape used depends on the application. This causes two problems. One is that it is beneficial to have some knowledge of the basis of the data being operated on. Another problem is that there is usually no fall back shape to act as a contingency plan should the wavelet be unsuitable.

The general form is given by the wavelet

$$W(a,b) = \frac{1}{\sqrt{a}} \int_{-\infty}^{\infty} h^* \left(\frac{t-b}{a}\right) y(t) dt \text{ transform.} \tag{3.384}$$

where $h^*$ is the complex conjugate of the wavelet equation $h(\text{-})$.

This is in a general form for profile data. Notice that the form of $h$ is arbitrary but the arguments $a$ and $b$ are not. The wavelet transform has a positional parameter $b$ and a scale parameter $a$. It does not contain frequency as such. The argument $\frac{t-b}{a}$ or $\frac{x-b}{a}$ could be more general still if raised to index $n$, which more precisely indicates multiscale properties than just $a$, especially in terms of fractal powers. This index is used for flexibility in other distributions such as the Weibull distribution although it has to be said that the kernel is exponential. However, $h\left(\frac{t-b}{a}\right)^n$ has fractal and wavelet echoes.

The digital equivalent of equation (3.390) is

$$W(iT,a) = T_s \frac{1}{\sqrt{a}} \sum_{n}^{N} h^* \left[\frac{(n-i)T_s}{a}\right] y(nTs) \tag{3.385}$$

where $N$ is the number of samples and $Ts$ is the sample interval. This representation is the discretized continuous wavelet transform (DCGT). This can be represented as a finite impulse response filter.

**Figure 3.55** DCGT

Here the $Z$s refer to the $Z$ transform and not height.

The methods of calculation of these transforms can be graphical. (See Handbook of Surface Metrology 1st Edition [55]).

## 3.9 Surface generation

### 3.9.1 Profile generation

The generation of surfaces in a computer for use in functional experiments requires a knowledge of digital characterization. This is essential in order to match the discrete behaviour of surfaces as measured with those generated artificially. The third link is to relate the digital representation to the theory. How to do this has been described earlier in the surface characterization section (e.g. reference [99] in chapter 2). Unless the digital generation is tied up properly, a comparison between simulation experiments will be meaningless. Consequently, a small review of the issues is needed to tie up the possible approaches to the generations that have been carried out in the past. To do this, consider first a profile. In its simplest form three running ordinates would be needed (to enable peaks as well as slopes to be described).

The joint probability density of three random variables is $p(z_1, z_2, z_3)$. This can be rewritten in terms of conditional densities to give

$$p(z_{1,}z_2,z_3) = p(z_1)p\left(z_2/z_1\right)p\left(z_3/z_2,z_1\right). \tag{3.386}$$

The first approach to generation would be to select one ordinate from a random number generator. This first number could be generated from any random number generator. This usually means that there is a uniform chance of any height occurring. There are very many possibilities for this starting number and it is not within the remit of this book to consider every possibility. However, an example might be useful.

If the set of generated numbers is to have zero mean value and zero correlation between the ordinates, then methods such as congruent and mid-square techniques could be used. Both are based on elementary number theory.

As an example, if the set of numbers are denoted by $[z_n]$, $n = 0, 1, 2, 3, \ldots$, the congruent method of generation is such that adjacent ordinates are related by the recursive relationship

$$z_{n+1} = az_n + b|\text{modulo } T|$$

where $b$ and $T$ are prime numbers and $T$ is matched to the wordlength of the computer, $a$ and $b$ can be chosen to give certain statistical properties; for example by selecting $a$, $b$ and $T$ the correlation between ordinates can be made equal to $\rho_s$, where

$$\rho_s = \frac{1 - 6b_s(1 - b_s/T)}{a_s} + e \tag{3.387}$$

and where $a_s = a^s \pmod{T}$, $b_s = (1 + a + a^2 + \ldots + a^{s-1})b \pmod{T}$, and $|e| < a_s/T$. If $a - T^{1/2}$ then $\rho_s \sim T^{-1/2}$. Correlated data can also be generated for the linear correlation case by the recursive relationship

$$z_{n+1} = \rho z_{n-1} + \sqrt{12(1-\rho^2)}(R_n - 0.5) \tag{3.388}$$

where $R_n$ is obtained from a set of independent random numbers, for example obtained by squaring the two centre digits of a six-figure number—not elegant, but effective.

In the method above one ordinate is correlated to the next. This is tantamount to a first-order Markov process. Most surfaces are more complicated and are produced by multiple processes. As a result the generator has to be at least second order, which means that an ordinate is related to the next ordinate and the one next to that. Furthermore the distribution of heights is not uniform: at best it is Gaussian, at worst it is not symmetrical and may well be skewed. These issues will be considered shortly.

The most important point in surface generation is to decide upon the correlation between ordinates. This is not arbitrary but can be at least estimated from the tribology prediction routines described in reference [56] and anticipated from the discrete analysis of random surfaces. For example, if it is required that a surface model needs to have curvature and slope variances of $\sigma_c^2$ and $\sigma_m^2$— to test a contact theory—then the required correlations between ordinates can be worked out from the previously determined formulae:

$$\sigma_s^2 = (1 - \rho_2)/2h^2 \qquad \sigma_c^2 = (6 - 8\rho_1 + \rho_2)/h^4 \tag{3.389}$$

where $h$ is the spacing between ordinates taken here to be a profile; $\rho_1$ is the correlation between ordinates and $\rho_2$ is the correlation between alternate ordinates.

The described values of $\rho_1$ and $\rho_2$ from equation (3.389) are given by

$$\rho_1 = \frac{1 - h^2\sigma_m^2}{2} - \frac{\sigma_c^2 h^4}{8} \qquad \rho_2 = 1 - 2h^2\sigma_m^2. \tag{3.390}$$

Other parameters can be demonstrated in a similar way.

Given these values of $\rho_1$, $\rho_2$ and $h$, the generation can take place. Assume for example that the distribution is Gaussian. This can be arranged by shaping the uniform distribution of height frequency according to the Gaussian function. Thus the first ordinate can be chosen from the distribution

$$\rho(z_1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z_1^2}{2}\right). \tag{3.391}$$

The second $(z_2|z_1)$ represents the second ordinate given the first at height $z_1$:

$$\rho(z_2|z_1) = \frac{1}{\sqrt{2\pi(1-\rho_1^2)}} \exp\left(\frac{(-z_2 - \rho_1 z_1)^2}{2(1-\rho_1^2)}\right) \tag{3.392}$$

and the third is the probability of getting a height $z_3$, given that the second took a value $z_2$ and the first $z_1$, that is $\rho(z_3|z_2,z_1)$:

$$\rho(z_3|z_2,z_1) = \frac{\sqrt{1-\rho_1^2}}{\sqrt{2\pi(1-\rho_2)(1+\rho_2 - 2\rho_1^2)}}$$
$$\times \exp\left\{\left[\left(z_3 - \frac{\rho_1(1-\rho_2)z_2}{(1-\rho_1^2)} + \frac{(\rho_1^2 - \rho_2)}{(1-\rho_1^2)}z_1\right)\middle/ \frac{2(1-\rho_2)(1+\rho_2 - 2\rho_1^2)}{(1-\rho_1^2)}\right]^2\right\}. \tag{3.393}$$

All of these equations follow directly from the multinormal equation (2.65).

The set of equations above read as follows. The value of $z_2$ given a previous $z_1$ has a standard deviation of $\sqrt{1-\rho_1^2}$ and mean of $\rho_1 z_1$ from equation (3.392). The following $z_3$ distribution (equation (3.398)) has a mean value of

$$\frac{\rho_1(1-\rho_2)}{(1-\rho_1^2)} z_2 - \frac{\rho_1^2 - \rho_2}{(1-\rho_1^2)} z_1 \tag{3.394}$$

and a standard deviation of

$$\frac{(1-\rho_2)(1+\rho_2 - 2\rho_1^2)(1+\rho_2)}{(1-\rho_1^2)}.$$

The mean values and standard deviations are found simply by factorizing out first $z_1$, then $z_2$ and finally $z_3$ from the multinormal distribution. So, in sequence, $z_1$ is picked at random from a distribution which is shaped to be Gaussian of mean level zero and standard deviation unity (actually $R^2_q$). This value of $z_1$ determines the mean from which the distribution of $z$ is centred, $z_2$ is picked from such a Gaussian distribution at this mean having a standard deviation which is not unity but $\sqrt{1-\rho_1^2}$ ; similarly for $z_3$. Having found these three ordinates, $z_2$ replaces $z_1$, $z_3$ replaces $z_2$ and a new $z_3$ is found. This repeated exercise generates the profile.

Alternatively, but equally, a best-fit least-squares routine can be used from the linear relationship:

$$z_i + Az_{i-1} + Bz_{i-2} + ... + \varepsilon_1 \tag{3.395}$$

where $\varepsilon_1$ is the residual or error term. Minimizing $\varepsilon_1^2$ and making $\sum \varepsilon_1^2 = S$

$$\frac{\partial S}{\partial A}, \frac{\partial S}{\partial B} \cdots = 0$$

gives the covariance matrix

$$\begin{pmatrix} 1 & \cdots & \rho_{n-1} \\ \rho_1 & & \vdots \\ \rho_{n-1} & & 1 \end{pmatrix} \begin{pmatrix} A \\ \vdots \\ n \end{pmatrix} = \begin{pmatrix} \rho_1 \\ \rho_2 \\ \vdots \\ \rho_n \end{pmatrix} \tag{3.396}$$

from which it emerges that

$$A = \frac{\rho_1(1-\rho_2)}{(1-\rho_1^2)} \qquad \text{and} \qquad B = -\left(\frac{\rho_1^2 - \rho_2}{1-\rho_1^2}\right).$$

The standard deviation is obtained simply from the determinant of the matrix and is of course equal to that given in (3.394).

If a non-Gaussian surface profile is to be generated, then this is straightforward. All that has to be done is to pick the $z_3$ from a non-Gaussian distribution. Obviously if the distribution of the residuals from which $z_3$ is picked is skewed, then the final surface will be skewed; similarly for kurtosis. There are relationships between the moments of the generated profile and the moments of the distribution from which the ordinate is picked.

Take for the simplest case the simplest surface. This is one which has an exponential correlation function. It corresponds to a first-order Markov process in the sense that there are only two terms in the generator $z_{i-1}$. Thus if $\sigma_p$, $S_{kp}$, $K_p$ are the moments of the profile and $\sigma_r$, $S_{kr}$, $K$, are those of the distribution from which the surface is being generated [57]

$$\sigma_p^2 = \frac{1}{(1-\rho^2)}\sigma_r^2$$

$$S_{kp}^{1/2} = \frac{S_{kp}^{1/2}(1-\rho^2)^{3/2}}{(1-\rho^3)}$$

$$K_p = \frac{K_r(1-\rho^2)+6\rho^2}{(1+\rho^2)}.$$

(3.397)

Note that if the kurtosis of the residual distribution is 3 then $K_p$ will also be 3. It is also true to say that the generated profile will always be closer to the Gaussian distribution than the residual distribution.

The relationships for a second-order equation are much more complicated but can be evaluated from the moment equations. Therefore

$$S_{kp} = \frac{M_p^3}{(\sigma_p^2)^{3/2}} \qquad S_{kr} = \frac{M_\tau^3}{(\sigma_r^2)^{3/2}}$$

where $M_r^3$ is obtained from the relationship

$$M_p^3 \cdot \left(1 - \varphi_2^2(2\varphi_1^2 + \varphi_2) - \frac{\varphi_1^2(\varphi_1 + \varphi_2^2)[1 + \varphi_2(\varphi_1 + 2\varphi_2)]}{1 - 2\varphi_1\varphi_1 - \varphi_2^3}\right) = M_r^3$$

(3.398)

where

$$M_p^n = \frac{1}{N}\sum_{i=1}^{N}(z-\bar{z})^n p(z)$$

and so on for $n = 2$, 3 and 4, from which $S_{kp}$ can be evaluated in terms of $S_{kr}$, and

$$K_p = \frac{M_p^4}{(\varphi_p^2)^2} \qquad K_r = \frac{M_r^4}{(\varphi_r^2)^2}$$

where $M_p^4$ is obtained from the relationship

$$M_p^4 \cdot (A - BC) = M_r^4 D + \varphi_p^2\varphi_r^2(E + BF)$$

and where

$$A = (1 - \varphi_1^4 - \varphi_2^4 - 4\varphi_1^2\varphi_2^3)(1 - \varphi_2^2) - 6\varphi_1^4\varphi_2^2$$

$$B = \frac{4\varphi_1\varphi_2[\varphi_1^2(1 + 2\varphi_2^2) + \varphi_2^3(1 - \varphi_2^2)}{1 - \varphi_2^2 - \varphi_2^4 + \varphi_2^6 - 3\varphi_1^2\varphi_2(1 + \varphi_2^2)}$$

$$C = \varphi_1(\varphi_1^2 + \varphi_2^3 - \varphi_2^5 + 2\varphi_1^2\varphi_2^2)$$

$$D = 1 - \varphi_2^2$$

$$E = 6[\varphi_1^2 + \varphi_2^2 - \varphi_2^4 + 2\varphi_1^2\varphi_2(1 + \varphi_2)]$$

$$F = 3\varphi_1(1 + \varphi_2 + \varphi_2^2)$$

(3.399)

where, for both (3.398) and (3.399), the $\varphi$ can be evaluated in terms of the correlation coefficients between ordinates $\rho_1$ and $\rho_2$. Thus

$$\varphi_1 = \frac{\rho_1(1 - \rho_2)}{(1 - \rho_1^2)} \quad \text{and} \quad \left(\frac{\rho_2 - \rho_1^2}{1 - \rho_1^2}\right).$$

(3.400)

General expressions for more complicated surfaces than second order are possible but decidedly complex. However, in practice, generations to more than second order are rarely needed. Other techniques are given in reference [57].

### 3.9.2 Two-dimensional surface generation

As stated earlier, the real requirement for functional studies is in two-dimensional rather than profile generation. However, as would be expected a number of extra problems arise. The first problem is to decide on the numerical model. It could be argued that an ARMA model would be as appropriate to the areal as to the surface profile. This is difficult to justify in practice because of the more difficult numerical stability problems associated with two-dimensional analysis.

The possible spatial configurations for two-dimensional generation are considerable even in its simplest form. For the reasons given above, only second-order surfaces will be considered—and even this produces problems.

Possible models are shown in figure 3.56(*a-c*). The correlation model is shown in figure 3.56*c*.

There has to be a compromise between the amount of coverage needed to achieve a certain accuracy and parsimony of the model to minimize computing time.

Figure 3.56 shows three configurations using the common grid sampling. There are others which have advantages in some respects but they will not be considered for generation here [20]. The model most often used is that shown in figure 3.56(*c*). Notice that even for this simple arrangement eight correlation coefficients are needed to cover the cross-correlation adequately between the *x* and *y* directions (here *y is* taken as the other horizontal axis and *z* the vertical axis).

For the preferred model it is assumed that the $z_{ij}$ value is approximated by the recursive relationship

$$z_{ij} = Az_{i-1.j} + Bz_{i-2.j} + Cz_{i.j-1} + Dz_{i.j-2} + Ez_{i-1.j-1}$$
$$+ Fz_{i+2.j-2} + Gz_{i-2.j-1} + Hz_{i-1.j-2} + \varepsilon_{ij}$$

(3.401)

where $\varepsilon$ is the random distribution of residuals. Taking expected values of the difference between the true and modelled values of $z_{ij}$ and using the least-squares linear technique yields the following unwieldly covariance matrix for solution:

$$
\begin{pmatrix}
1 & \rho_1 & \rho_5 & \rho_8 & \rho_3 & \rho_8 & \rho_5 & \rho_4 \\
\rho_1 & 1 & \rho_7 & \rho_6 & \rho_5 & \rho_4 & \rho_3 & \rho_8 \\
\rho_5 & \rho_7 & 1 & \rho_3 & \rho_1 & \rho_7 & \rho_2 & \rho_5 \\
\rho_8 & \rho_6 & \rho_3 & 1 & \rho_5 & \rho_2 & \rho_7 & \rho_1 \\
\rho_3 & \rho_5 & \rho_1 & \rho_5 & 1 & \rho_5 & \rho_1 & \rho_3 \\
\rho_8 & \rho_4 & \rho_7 & \rho_2 & \rho_5 & 1 & \rho_3 & \rho_1 \\
\rho_5 & \rho_3 & \rho_2 & \rho_7 & \rho_1 & \rho_3 & 1 & \rho_5 \\
\rho_4 & \rho_8 & \rho_5 & \rho_1 & \rho_3 & \rho_1 & \rho_5 & 1
\end{pmatrix}
\times
\begin{pmatrix}
A \\ B \\ C \\ D \\ E \\ F \\ G \\ H
\end{pmatrix}
=
\begin{pmatrix}
\rho_1 \\ \rho_2 \\ \rho_3 \\ \rho_4 \\ \rho_5 \\ \rho_6 \\ \rho_7 \\ \rho_8
\end{pmatrix}.
$$

(3.402)

The significance of the correlation values are shown in figure 3.56. Obviously because of symmetry, the expected value of the product of $z_{i-1,j-1}$ is the same, for example, as $z_{i-2,j-1}$ (which in this case is $\rho_5$).

There is no restriction in this equation on the distribution of the residual (i.e. the distribution from which the current point is selected). Non-Gaussian statistics can be generated in exactly the same way as for

**Figure 3.56** Numerical methods for generation—correlation coefficients required for surface generation.

the profile, but the relationship between the skew and kurtosis of the noise distribution (residual) and the values for the generated surface are much more complicated and not presented here. By the same token, as before, the effect of non-Gaussian skew and kurtosis values of the residual distribution is considerably attenuated in the same parameter of the surface. Whatever the non-normality of the parameters the desired correlation coefficients are maintained.

Equation (3.402) has to be inverted to give the values of $A–H$ for any given set of values of $\rho_1$ to $\rho_8$. However, it has to be realized that because of the cross-terms, the values of $\rho$ have to be consistent with each other, so

$$|A| = |\text{cov}|^{-1} |\rho| \tag{3.403}$$

where $|\text{cov}|$ is the covariance matrix which by definition is symmetrical about its main diagonal.

In order that a solution is possible the covariance matrix has to be positive definite. Furthermore, the matrix values of $\rho$ have to be self-consistent as mentioned above.

Thus, if it is assumed that the surface structure comprises two independent lay patterns, for example $z(x,y)=s(x)t(y)$, then the autocorrelation-coefficient corresponding to a diagonal trace $A(\alpha,b)$ is given by

$$A(\alpha, \beta) = \frac{1}{L^2} \int_0^L \int_0^L s(x)t(y)s(x+\alpha)t(y+\beta)\mathrm{d}x \ \mathrm{d}y$$
$$= A(\alpha)A(\beta) \tag{3.404}$$

which corresponds in the figure to the conclusion that, for example, $\rho_7 = \rho_2\,\rho_1$, $\rho_6 = \rho_2\,\rho_4$ etc. If $\rho_2, \rho_1, \rho_4$, etc, are not made the matrix in equation (3.402) will be ill-conditioned.

As a further check on the conditioning, the alphabetic constants $A–H$ are related in the same way, that is $E = -AC$, $F = -BD$ and so on.

The variance of the distribution for the surface given by the model shown in figure 3.56(*c*) is

$$
\begin{aligned}
\text{var} = {} & 1 + A^2 + B^2 + C^2 + D^2 + E^2 + F^2 + G^2 + H^2 \\
& -2(A\rho_1 + B\rho_2 + C\rho_3 + D\rho_4 + E\rho_5 + F\rho_6 + G\rho_7 + H\rho_8) \\
& +2
\begin{pmatrix}
AB\rho_1 + AC\rho_5 + AD\rho_8 + AE\rho_3 + AF\rho_8 + AG\rho_5 + AH\rho_4 \\
+\ BC\rho_7 + BD\rho_6 + BE\rho_5 + BF\rho_4 + BG\rho_4 + BH\rho_8 \\
+\ CD\rho_3 + CE\rho_1 + CF\rho_7 + CG\rho_2 + CH\rho_5 \\
+ DE\rho_5 + DF\rho_2 + DG\rho_7 + DH\rho_1 \\
+\ EF\rho_5 + EG\rho_1 + EH\rho_3 \\
+ FG\rho_3 + FH\rho_1 \\
+\ GH\rho
\end{pmatrix}.
\end{aligned}
$$

$$(3.405)$$

From this it is obvious that the variance has to be less than unity and greater than zero, which provides some constraint on the chosen values of $\rho$. Unfortunately, owing to the size of the array, it is difficult to identify troublesome computations. Note that the direct terms $A\rho_1$ etc tend to increase the general correlation of the model, whereas the cross-terms $AB\rho_1$ tend to reduce it.

A further consideration is concerned with the physical sense of the correlation coefficients being modelled. The correlation coefficient values have to be realistic in the *second-order* system sense, otherwise the generation suffers. One consideration is that $1 + \rho_2 - 2\rho_1^2$ has to be positive. This has to be true for $\rho_3$ and $\rho_4$ and for any other possible linear track across the grid.

As an example of this type of generation, correlation values can be selected and the surface generated in two dimensions. Obviously in order to start the generation off two initial random tracks have to be generated in the *x* and *y* direction. These four tracks should be discarded at the end of the run. A typical 2D surface would have $100 \times 100$ tracks. One example of such a generation is shown in figure 3.57.

The overall picture of acceptability of this method is not easy, but a presentation of the fidelity of results for given values of $\rho_1$ and $\rho_2$ (and $\rho_3$ and $\rho_4$) is shown in figure 3.58.



$(\rho_1 = 0.3,\ \rho_2 = -0.5,\ \rho_3 = -0.3,$

$\rho_4 = 0.3,\ \rho_5 = -00.9,\ \rho_6 = -0.15,$

$\rho_7 = 0.15,\ \rho_8 = -00.9)$

**Figure 3.57** Typical example of areal surface generation with correlation values specified.

**Figure 3.58** Predictability of surface generation, results showing intended and computed points.

The points to note are that in quadrants $(+\rho_1, -\rho_2)$, $(-\rho_1, -\rho_2)$, $(-\rho_1, \rho_2)$ a good agreement between the arbitrary $\rho_1 \rho_2 (\rho_3 \rho_4)$ values is possible. In the quadrant $(+\rho_1, +\rho_2)$ the second-order constraint becomes obvious as it is to some extent in the $(-\rho_1, \rho_2)$ quadrant, simply because of the difficulty of making a 2D second-order fit.

Despite the problems of consistency in the modelling of the 2D case, which will depend on its severity on the nature of the degree of isotropy, this method provides very convincing surfaces for the computer simulation of tribological functions. No doubt there are other methods available, but perhaps these do not incorporate the basic surface metrology needed in such investigations.

The foregoing 2D generation has been concerned with random surfaces. In most tribological situations a geometrical form has to be superimposed. This might be the shape of a cam, a circle or an involute. On top of this there might be a sinusoidal variation due to chatter; the possibilities are endless. Nevertheless it is the surface texture which is the essential and difficult component to which this type of generation is addressed. The other are to a greater extent deterministic.

### 3.12 Summary

This chapter has concentrated on ways of processing the many surface parameters which are in use. Considerable attention has been paid to filtering because of its importance in separating different geometrical components and also because of its use in providing reference lines. Also, sections have been devoted to evaluating references using minimax methods. This is because they are now being used more often, especially as dimensional and surface metrology get closer together. Some of the ways of using new transforms have also been given. Obviously, much work has been carried out in all these fields elsewhere, but little has been targeted at surface metrology, which has its own set of completely different problems. It is hoped that this section shows how to address some of these problems computationally.

This chapter is necessarily written at different levels of difficulty because of the wide range of complexity of the metrological parameters and references which have to be evaluated. It should, however, be most useful to researchers embarking on a surface research programme and possibly people who are building or modifying instruments. It is not likely to be so useful to people operating instruments, but it may give an insight into how the parameters are evaluated. There has throughout the chapter been some emphasis on mathematical techniques. This has been intentional. It is hoped that by exploring some of these techniques,

the subject itself can be advanced. The purely digital methods described here have usually been discrete ways of obtaining parameters described mathematically in chapter 2, although in some instances more than one method has been included. Nanometrology does not need the same bandwidth as surface metrology and, from the processing point of view, is simpler. Also, nanometrology is generally planar and so does not often involve complex 3-dimensional shapes.

In the next chapter, specific ways of measuring surfaces will be considered ranging from the well-known stylus methods, both traditional and modem, to optical and other methods. Chapter 4 should therefore enable people to make a better choice of instrument to use in a particular situation.

## References

[1] *Allied and Interpolation Tables* 1960 (London: HMSO)
[2] Whitehouse D J 1978 The digital measurement of peak properties *IMechE J. Mech. Eng. Sci.* **20** 119
[3] Chetwynd D G 1978 Slope measurement in surface texture analysis *IMechE J. Mech. Eng. Sci.* **20** 115–19
[4] Whitehouse D J and Phillips M J 1982 Two-dimensional discrete properties of random surfaces *Philos. Trans. R. Soc.* A **305** 441–68
[5] Rice S O 1944 Mathematical analysis of random noise *Bell Syst Tech* J **23** 282, 1945 **24** 46, **27** 109
[6] Bendat J S 1958 *Principles and Applications of Random Noise Theory* (New York: Wiley)
[7] Dong W P and Stout K J 1995 Two dimensional FFT and power spectrum for surface roughness in two dimensions *Proc. Inst Mechanical Engineering*. **209** 81
[8] Rabiner L R and Gold B 1975 *Theory and Application of Digital Signal Processing* (New York: Prentice Hall)
[9] Tholath J and Radhakrishnan V 1997 Three dimensional filtering of engineering surfaces using a least squares non paramelvic B spline surface *Proc. Inst Mechanical Engineering*. **211** B 557
[10] Longuet-Higgins M S 1957 Statistical analysis of a random moving surface *Proc. R.* Soc. A **249** 966
[11] Nayak P R 1971 Random process models of rough surfaces *Trans. ASME J. Lubr. Technol.* **93** 398
[12] Greenwood J A 1984 A unified theory of surface roughness *Proc. R. Soc.* A **393** 133–57
[13] Sayles R S and Thomas T R 1977 Measurement of the statistical nucrogeometry of engineering surfaces *Proc. 1st Joint Poly Symp. on Manufacturing Engineering, (Leicester, UK)*
[14] Whitehouse D J and Archard J F 1970 The properties of random surfaces of significance in their contact *Proc. R. Soc. A* **316** 97–121
[15] Cheng M C 1969 The orthant probabilities of four Gaussian variates *Ann. Math. Stat.* **40** 152–61
[16] Plackett R L 1954 A reduction formula for normal multivariate integrals *Biometrika* **41** 351–60
[17] Whitehouse D J and Phillips M J 1978 Discrete properties of random surfaces *Philos. Trans. R. Soc.* A **290** 267
[18] Whitehouse D J and Phillips M J 1985 Sampling in a two-dimensional plane *J. Phys. A*: *Math. Gen.* **18** 2465
[19] Staunton R C 1987 The design of hexagonal sampling structures for image digitisation and their use with local operators *Image and Vision Computing* (Guildford: Butterworths)
[20] Li M, Phillips M J and Whitehouse D J 1989 Extension of two-dimensional sampling theory *J. Phys. A*: *Math. Gen.* **11** 5053
[21] Lukyanov V S 1980 Local irregularities of the surface and their influence on surface parameters *Ann. CIRP* **29** 423
[22] Staufert G 1979 Characterisation of random roughness profiles *Ann. CIRP* **28** 439
[23] Smith E H and Walmsley W M 1979 Walsh functions and their use in assessment of surface texture *Wear* **57** 157–66
[24] Blackman R B and Tuckey J W 1958 *The Measurement of Power Spectra* (New York: Dover)
[25] Raja J and Lui X Analysing Engineering Surface Texture Using Wavelet Filter. SPIE vol. 2825
[26] Bodschwinna H 1997 Presentation on Modified Sauni-Filter ISO/TC213 A99 San Diego CA
[27] Chen X, Raja J and Simanapalli S 1995 Multiscale of engineering surfaces *Int. Journal of Mechanical Tool and Manufacturing* **35** 2
[28] Xu J Blunt L and Stout K. Raised Wavelet. *Proc. Royal Soc.* A.
[29] Reason R E 1962 Report on reference lines for roughness and roundness *Ann. CIRP 1* 96
[30] Temoshenko 1934 *Theory of Elasticity* (New York: McGraw Hill)
[31] Whitehouse D J 1970 *PhD Thesis* Leicester University
[32] Sayles R S 1986 *PhD Thesis* Teesside Polytechnic
[33] Forbes A B 1989 Least squares best fit geometric elements *NPL Rep. /DITC* 140/89
[34] BS 7172
[35a] Decker J A and Harwit M 1969 Experimental operation of a Hadamard spectrometer *Appl. Opt.* **8** 2552–4
[35b] Scott P 1993 Height characterization using the material probabilities curve BSITC57 SC1 WG3 N20
[36] Whitehouse D J 1973 A best fit reference line for use in partial arcs *J. Phys. E*: *Sci. Instrum.* **6** 921–4
[37] Lotze W 1983 Generalised fitting algorithms in the coordinate measuring techniques in quality control *Ada Imeko* **32** 279-86
[38] Murthy T S R and Abdin S Z 1980 Minimum zone evaluation of surfaces *JMTDR* **20** 123–36
[39] Dhanish P B and Shunmugam M S 1991 Algorithm for form error evaluation—using the theory of discrete and linear Chebychev approximation *Comput. Methods Appl. Mech. Eng.* **92** 309–24

[40] Avdulov A 1967/8 *Proc. IMechE* **182** pt 3K p 425

[41] Chetwynd D G 1985 Application of linear programming to engineering metrology *Proc. IMechE B* **2199** 93–100

[42] Osborne M R and Watson G A 1968 A best fit linear Chebychev approximation *Comput. J.* **10** 172–7

[43] Chetwynd D G 1980 *PhD Thesis* Leicester University

[44] Preparata F P and Shamos M F 1985 *Computational Geometry* (New York: Springer)

[45] Whitehouse D J 1982 The parameter rash—is there a cure? *Wear* **83** 75–8

[46] Pratt W K, Kane J and Andrews H C 1969 Hadamard transform image coding *Proc. IEEE* **57** 58–69

[47] Yolles M I, Smith E H and Walmsley W M 1982 Walsh theory and spectral analysis of engineering surfaces *Wear* **83** 151

[48] Mulvaney D J 1983 *PhD Thesis* University of Leeds

[49] Daubechies I 1988 Orthonormal bases of compactly supported wavelets *Commun. Pure Appl. Math.* **41** 909–96

[50] Woodward P M 1953 *Probability and Information Theory with Applications to Radar* (Oxford: Pergamon)

[51] Wigner E 1932 On the quantum correction for thermodynamic equilibrium *Phys Rev* **40** 749

[52] Whitehouse D J and Zheng K G 1992 The use of dual space-frequency functions to machine tool monitoring *J. Phys. Meas.*

[53] Zheng K G and Whitehouse D J 1991 The application of the Wigner distribution function to machine tool monitoring *Proc. IMechE J. Mech. Eng. Set.* **206** 249–64

[54] Wexler J and Raz S 1990 Discrete Gabor expansions *Signal Process* **21** 207–21

[55] Whitehouse D J 1994 *Handbook of Surface Metrology 1st Edn* Graphical Methods p 328

[56] Whitehouse D J 1983 The generation of 2D surfaces having specified function *Ann. CIRP* **32**(1)

[57] Patir N 1978 A numerical procedure for random generation of rough surfaces *Wear* **47** 263–77

[58] Bentley F and Grant P 1995 Wavelets and the wavelets transform *Proc.* TFTS **95** 13

[59] Tholath J and Radhakrishnan V 1997 Three dimensional filtering of engineering surfaces using a least squares non parametric B spline surface *Proc. Inst. Mech. Eng*. **211** B *557*.

# Chapter 4
# Instrumentation

## 4.1   Introduction and history

This chapter is concerned with obtaining geometrical information from the surface by means of a pick-up and converting it into an electrical signal suitable for processing by means of a transducer. Included in the consideration of the pick-up is the means whereby the surface is scanned.

What is not included in this chapter is any method of processing, such as filtering, which may be dependent on the application and need not be common to a specific type of instrument. Processing details are confined to chapter 3 and then only in terms of algorithms and not hardware.

The chapter begins with a short discussion of some of the basic principles of instrument design, without which any attempt to recover the very small-scale roughness information would be doomed to failure. Included in this section are subsections on metrology loops, drives and other considerations for ensuring a stable and error-free method of scanning the surface.

The main body of the chapter considers the stylus technique (section 4.2). This is principally concerned with the conventional tactile stylus method but also extends to the new generation of scanning microscopes (section 4.2.3) in which, although styluses are used, the techniques are not easily described as contacting methods. Nevertheless the basic pick-up mechanism and drive system is so similar in concept to the tactile pick-up that it is included in this section rather than later on in section 4.8 on electron microscopy.

Following on from the stylus method are other pick-up techniques such as optical methods (section 4.3), capacitative techniques (section 4.4), inductive methods (section 4.5) and non-conventional methods (section 4.7), up to electron microscopy in section 4.8. The last part of the chapter is devoted to a consideration of transducing techniques (section 4.9).

### 4.1.1   Historical details

The development of instruments for examining surfaces began in about 1919 with a simple mechanical device by Tomlinson at the NPL. The first instrument for engineering use is ascribed to Schmalz [1]. Earlier than this, thought had already been given to the measurement of engineering surfaces using optical methods by Berndt in 1924 [2] and Andrews in 1928 [3]. That there was concern about surface measurement in countries other than Germany and the UK becomes apparent when it is realized that the first call for surface finish standards was made in the USA by Harrison in 1930 [4]. (Needless to say it took 10 years before the first B46 standard emerged.)

This may have stimulated Abbott to develop the now familiar stylus technique in 1933 [5] (although Schmalz's device was also based on a stylus). Although reported in 1933 Abbott's instrument was not in general use for some years. It was the first practical instrument for surface finish to be patented (US Patent 2,240,278, British Patent 523436). One of the very first was used in the Chrysler factory in 1936. Only a few hundred were ever made. The surface analyser made by the Brush Development Company of Cleveland in 1935 came somewhat later but was more popular because it had a chart recorder [6].

More general development of practical optical instruments by Carl Zeiss took place in Germany in 1934 [7], although earlier Linnik in the USSR in 1930 had been evolving novel interferometric methods specifically for looking at surfaces [8]. The first work on measuring surfaces using capacitance methods seems to have been due to Perthen in 1936 [9] who also became equally involved in pneumatic methods. However, it is usual to attribute the first practical pneumatic methods to General Nicolau in France in 1939 who, no doubt, did not develop them fully because of other preoccupations [10]. In 1942 Von Weingraber became more involved in pneumatic methods [11].

Development of surface finish instruments in the UK started somewhat later than in the USA. It started mainly as a result of Schlesinger's attempt in 1939 to investigate the problems of specifying surfaces adequately. Fortunately, Reason at Taylor Taylor Hobson became interested and soon became the prime mover. These and early US developments are described in Schlesinger's early articles (1942) and book (1940) [12-14].

Optical methods of measuring surfaces had not been neglected in the UK but had not developed into workshop-type instruments [3,15]. The most successful optical instrument for surface measurement was the autocollimator for use in flatness rather than roughness measurement (1942 [16]).

### 4.1.2 Some early dates of importance in the metrology and production of surfaces

1631 Introduction of Vernier system of linear measurement by Pierre Vernier
1769 Smeaton's first boring machine for cannon
1775 Watts steam engine based on Wilkinson machine
1800 High-carbon steels used for cutting tools
1865 Robert Musket of Sheffield introduced semi-high-speed steel
1867 Vernier calipers manufactured by Brown and Sharp
1886 Reynolds' paper to Royal Society on hydrodynamic theory
1890 Introduction of synthetic abrasive grinding materials
1895 Micrometer introduced
1896 Introduction of gauge blocks by Johanson—progressive tolerance principle
1898 Chip analysis mode of turning
1900 Introduction of high-speed tool steel
1904 Nicolson analysis of tool pressures in chip removal
1911 Commercial manufacture of gauge blocks
1915 Commercial introduction of centreless grinding
1916 Development of cemented carbides for cutting tools in Germany
1919 Mechanical surface profile by Tomlinson of Natural Physical Lab.
1922 Introduction of first lapping machine
1929 First tracer-type surface machine by Schmalz in Germany
1933 Abbott's profilometer conceived in USA
1934 Linnik surface microscope in USSR
1934 Gloss method, Carl Zeiss, Germany
1935 Flemming's tracer profilometer, Germany
1936 Brush surface analyser in USA
1936 Perthen capacitance surface gauge, Germany
1936 Superfinishing process introduced by Chrysler Corporation in USA
1938 Tomlinson surface finish recorder in UK
1939 Nicolau pneumatic gauge, France
1940 Talysurf developed in UK
1940 First B46 standard on surface finish, USA
1942 Machined surface standards, Norton, USA

1943    Use of diamond cutting tools reported, UK
1944    First replication methods for surfaces 'FAX FILM'

Summarizing the instrument developments, stylus methods only became practical with the advent of electronic methods of amplifying the very small displacements that were found to be involved. This seems now to be a small thing but at the time it was very important because it meant that the output signal was controlled by the stylus movement and not driven by it. For the first time the surface modulated the power needed to derive the output signal.

Thus it came about that, without regard for the relative importance of errors of form, waviness and roughness, the first consideration was given to the measurement of roughness, particularly with the use of the now familiar sharp, lightly loaded stylus. It is from this first, comparatively recent beginning that the ideas about the nature of surface topography, the need for treating it on a bandpass sampling basis and the instrumentation needed to cover texture, straightness and roundness evolved.

As mentioned above, the earliest practical stylus instrument is generally attributed to the German engineer Gustav Schmalz [1], who described an arrangement in which a stylus on the end of a pivoted arm traversed the surface and tilted a small mirror, as used in a reflecting galvanometer. The excursions of the spot, when recorded on a moving photographic chart, produced a profile graph showing a cross-section of the surface traced by the stylus. With this arrangement the maximum possible magnification of the stylus movement is limited by the very small force that can be applied to the sharp stylus, which in turn limits the size of the mirror and, hence, the numerical aperture of the optical system on which the definition of the spot depends. Schmalz's apparatus sufficed to produce some interesting profile graphs of the rougher surfaces and stimulated further work. Earlier in 1919 Tomlinson at the National Physical Laboratory in UK had developed a purely mechanical stylus instrument which used a sooted glass as a record. This was never a commercial instrument [17] but it was probably the first attempt to magnify and record a surface as a profile.

The next significant advance came when Ernest Abbott [5] in the USA developed an electrical instrument in which the output from a coil, coupled with the stylus and moving through a magnetic field, was amplified and fed into a voltmeter. Since the output of such a moving-coil transducer is responsive basically to velocity, Abbott had to interpose a circuit attenuating the output inversely as the frequency so that the meter indicated amplitude, regardless of the frequency, over a limited, yet controlled range of frequencies wide enough to represent the 'hills and valleys' of roughness. The resulting instrument appeared on the market in 1936 as the first practical workshop device, under the registered name of 'Profilometer', giving for the first time definitive meter readings of amplitude, but no profile graphs. The name of this instrument bestowed the subject with the word 'profile' which is so apt as an easy description of what was possible. The meter responded to RMS values rather than average values—a fact that tended to cause some confusion later on. It was subsequently rectified somewhat by dividing the scale by 1.11, giving the correct average value for a sine wave. Unfortunately no one in those days thought that it was necessary to cater for random signals, so this factor was applied to every surface. In the UK Schlesinger quickly came down in favour of a straightforward average value and so the CLA (centre line average) was initiated.

The mechanical datum, relative to which the excursions of the stylus were measured, was provided by the locus of a pair of skids of 6 mm radius which engaged and slid over the surface. It is interesting to note that the idea of an external datum from which to measure the roughness was either not considered necessary at the time or just not thought of!

In late 1939 in the UK Taylor Taylor Hobson (later to become Taylor Hobson Ltd.) introduced the Talysurf instruments. The name Talysurf was not derived from Taylor, as many people thought, but from the Greek *taly*—to measure. The Talysurf instruments provided the enabling feature for standardization and records, namely the profile graph; it also had calibration standards. (See the 'Biographical memoirs of R E Reason' by D J Whitehouse [17].)

These instruments used a small variable inductance as an amplitude-modulating transducer in a simple ac bridge circuit. For passing interest, this transducer form was adapted from instruments for measuring the torsion in ship propeller shafts developed by the Admiralty Research Laboratory in Teddington, London, in 1919. Unfortunately, the only AC detectors then available were too much of a laboratory nature to be used in workshops and it was not until some fifteen years later, after the convenient copper oxide rectifier from Westinghouse had become available, that the transducer reappeared in a more refined form in metrology as the Pratt and Whitney Electrolimit Gauge. This instrument, with upwards of 100 gm force on the stylus tip, gave sufficient output from a mains-operated bridge to operate a metal rectifier meter without amplification. This made the measurement of the average surface roughness (CLA) straightforward. From these rather disjointed beginnings came the modern stylus instruments.

In parallel with these earlier developments of the stylus instruments was the development of the pneumatic instrument [10] and the capacitance instrument [9]. But in both cases, because the skirt of the pneumatic gauge and the electrode of the capacitance device were large—to encompass an area of the surface rather than to generate a point-by-point profile—both methods faltered in attempts to make them universally applicable because of the considerable difficulty in measuring curved surfaces.

Optical methods in the earlier days were usually related to interferometer methods like Linnik's [8] and later Tolansky's [17]. Those which produced a profile were more or less restricted to the Schmalz 'light slit' method, described later in the text for rougher surfaces; the optical interference methods only produced fringes for the finer surfaces. Both of these methods suffered from the fact that to get a record, a photographic method was required and that, even if this were available, the estimation of an average value required the use of a planimeter. (See section 3.9.1.) Even the potential gain in reliability by having an areal assessment in the rival methods and the basically cheaper units did not enable these three techniques (pneumatic, capacitative and optical) to compete with the stylus instrument. Its versatility, both mechanical and electrical, ruled the day. Although this dominance is now being challenged, the smooth progression in the development of the stylus method has enabled the relatively respectable level of surface knowledge to be realized.

### 4.1.3   Specification

In parallel with and feeding the instrument developments were calls for surface characterization; it is difficult to control surfaces if they cannot be measured and numerically assessed.

Apart from Harrison's call for standards in 1930 [4], there was considerable discussion about what was needed, for example by Schurig [18]. 'How should engineers describe a surface?' This problem is still with us. In the UK engineers like Timms and Jost [18,19] were querying the current practice. In fact, even as early as 1945 Timms was interested in the characterization and measurement of waviness, especially as it related to gears [20].

Even measuring and understanding the surfaces is not enough. This was realized very early on when calls for standardization were made. In particular Ransome [21] and Peale [22] in 1931 in the USA and Bodart [23] in 1937 in France tried very early to get practical methods adopted. As well as the standards themselves it was also quickly recognized that suitable symbolism should be placed on drawings so that specific finishes could be recognized. Broadston in 1944 [24] was an early pioneer in this aspect of standardization.

Notwithstanding the great contributions that have been made by many engineers on both sides of the Atlantic, those who probably made the most impact were Abbott in the USA, Reason in the UK and Schmalz and Perthen in Germany. The basic reason for this is simple: they not only tried to invent and build instruments with which to measure surfaces, but they also contributed greatly to the methods of assessment and characterization. The integration of these two different aspects of metrology was and still is the most important task of the metrologist, a theme which runs throughout this book.

In the next section some simple concepts of instrument design will be given, followed by a description of the basic philosophy behind the different ways of measuring surfaces.

### 4.1.4 Design criteria for instrumentation

Some of the major points of design are given here. For a fuller exposition see Smith and Chetwynd [25] and Moore [26].

The basic questions that need to be addressed in any instrument or machine tool design are as follows:

1. Is it correct according to kinematic theory, that is, does it have the required number of constraints to ensure that the degrees of freedom conform to the movements required?
2. Where are the metrology and force loops? Do they interact? Are they shared and, if so, how?
3. Where are the heat and vibration sources and, if present, how can they influence performance?
4. Is the design symmetrical and are the alignment principles obeyed?
5. What other sources of error are present and can they be reduced by means of compensation or nulling?

These points have been put forward to enable checks to be made on a design. It is not the intention here to investigate them fully. (See sections 4.2 and 4.3.)

Some basic principles can be noted:

1. Any unconstrained body has six degrees of freedom, three translational *x, y, z* and three rotational α, β, γ.
2. The number of contact points between any two perfectly rigid bodies is equal to the number of constraints.

There are a number of additional useful rules.

3. Any rigid link between two bodies will remove a rotational degree of freedom.
4. The number of linear constraints cannot be less than the maximum number of contacts on an individual sphere. Using the rules for linking balls and counting contacts it is possible to devise a mechanism which has the requisite degrees of freedom. See figure. 4.1.

It should be noted that a movement is that which is allowed in both the positive and negative direction. Similarly, a constraint is such that movement in one direction or both is inhibited. For example, it is well known that a kinematic design can fall to pieces if held upside down.

5. A specified relative freedom between two subsystems cannot be maintained if there is underconstraint.
6. Overconstraint does not necessarily preclude a motion but invariably causes interaction and internal stresses if used and should therefore be avoided.

### 4.1.5 Kinematics

Kinematics in instrument design are an essential means to achieve accurate movements and results without recourse to high-cost manufacture. In surface metrology instruments kinematics are most often used in generating straight line movements or accurate rotations.

In both cases there have to be five constraints, the former leaving a translation degree of freedom and the latter a rotation. By far the largest class of mechanical couplings are for one degree of freedom. Couplings with more than one degree of freedom are seldom required, the obvious exception being the screw, which requires one translation and one rotation. Despite its apparent simplicity, two translations cannot be achieved kinematically. *x, y* slides are merely two systems of one degree of freedom connected together.

(a)

1) in $z$   5) $\alpha, \beta, \gamma, x, y$

Ball on flat

(b)

Points of
contact

2) in $x, z$   4) $\alpha, \beta, \gamma, y$

Ball on vee

(c)

Points of
contact

3) in $x, y, z$   3) $\alpha, \beta, \gamma,$

Ball in trihedral hole

(d)

4) in $x, y, z, \alpha$   2) $\gamma, \beta,$

Ball link in two vees

(e)

5) in $x, y, z, \beta, \gamma$   1) $\alpha$

Shaft in trihedral hole
against a vee

**Figure 4.1** Kinematic examples.

A typical translation system of one degree of freedom is shown in figure 4.2. This has five pads on the carriage resting on the slide. The pads are usually of polymer.

Another simple system has a tube, which needs to move linearly resting on four balls (figure 4.3). The tube has a slot in it into which a pin (connected to the frame of the instrument) is inserted to prevent rotation.

For one degree of rotational freedom the drive shaft has a ball or hemisphere at one end. This connects against three flats suitably arranged to give the three points of contact. The other two are provided by a vee made from two rods connected together (figure 4.4($a$)). Another version ($b$) is a shaft in two vees and resting

**Figure 4.2** Linear movement carriage showing support pads.



**Figure 4.3** Linear movement tube.



**Figure 4.4** System showing single rotational degree of freedom.

on a flat. As mentioned earlier, these would have to be encouraged to make the contacts either by gravity or by a non-intrusive spring. There are numerous possibilities for such designs, which really began to see application in the 1920s [1,2].

The other main use of kinematics in surface metrology is in relocation. In this use it is necessary to position a specimen in exactly the same position after it has been removed. Such a requirement means that the holder of the specimen (or the specimen itself) should have no degrees of freedom—it is fully constrained. The most usual arrangement to achieve this is the Kelvin clamp (figure 4.5) [27].

The six degrees of freedom are constrained by three points of contact in hole C which, ideally, is a tetrahedron. Two points of contact in a vee groove at A and one on a flat at B. Another version is that of two non-parallel vees, a flat and an edge. Another version just has three vees. The relocation requirement is usually needed when experiments are being carried out to see how the surface changes as a result of wear or some similar phenomenon.

Couplings with two degrees of freedom require four constraints. A bearing and its journal with end motion, Hooke's joint, the so-called 'fixed' knife edges and four-film suspensions are common examples of this class.

**Figure 4.5** The Kelvin clamp.

One coupling with three degrees of freedom is the ball and socket joint. The tripod rests on a surface or a plane surface resting on another plane and three filar suspensions. Alternatively, the tripod could rest on a sphere in which only rotations are possible.

Relocation configurations are sometimes a little complicated by the fact that the specimen has to relocate exactly in more than one piece of apparatus. An example could be where the surface of a specimen has to be measured on a surface roughness measuring instrument and also has to be located in a wear machine. In these cases the specimen holder has to cater for two systems. Under these circumstances the specimen itself may have to be modified in shape or added to as shown in figure 4.6. In this case the specimen is a cylindrical specimen. Attached to it is a collar which has a stud fixed to it. This collar allows a constraint to be applied to the specimen rotationally. This would not be possible with the simple specimen. A collar or fitting such as this makes multiple relocation possible.

Couplings with four degrees of freedom require two constraints. A knife edge contacting with a plane bearing is a simple example of this class. Couplings with one degree of freedom include a ball on a flat.

What has been said above about kinematics is very elementary but necessary and is the basis from which designs should originate. Often the strict kinematic discipline has to be relaxed in order to meet other constraints such as robustness or ease of manufacture, for example by using a conical hole instead of a trihedral hole for the location of spherical feet. Nevertheless the principles still hold true.

Above all, the hallmark of a kinematic design is its simplicity and ease of adjustment and, not insignificantly, the perfection of motion which produces a minimum of wear. In general, kinematic instruments are



**Figure 4.6** Specimen modified to allow relocation.

easy to take apart and reassemble. Kinematic designs, because of the formal point constraints, are only suitable for light loads—as in most instruments. For heavy loads pseudo-kinematics are used in preference because larger contacts areas are allowed between the bearing surfaces and it is assumed that the local elastic or plastic deflections will cause the load to be spread over them. The centroids of such contacts are still placed in positions following kinematic designs. This is generally called pseudo or semi-kinematic design.

### 4.1.6 Pseudo-kinematic design

This is sometimes known as elastic or plastic design.

Practical contacts have a real area of contact, not the infinitesimal area present in the ideal case in kinematics. As an example a four-legged stool often has all four legs in contact, yet according to kinematics three would be sufficient. This is because the structure flexes enough to allow the other, fourth, leg to contact. This is a case of overconstraint for a good reason, that is safety in the event of an eccentric load. It is an example of elastic design or elastic averaging.

Because, in practice, elastic or eventually plastic deformation can take place in order to support a load, semi-kinematic or pseudo-kinematic design allows the centroid of an area of contact to be considered as the position of an ideal point contact. From this the rules of kinematics can be applied.

In general, large forces are not a desirable feature of ultra-precision designs. Loads have to be supported and clamps maintained. Using kinematic design the generation of secondary forces is minimized. These can be produced by the locking of a rigid structure and when present they can produce creep and other undesirable effects. For kinematics the classic book by Pollard [27] should be consulted.

In general another rule can be recognized, which is that divergence from a pure kinematic design usually results in increased manufacturing cost.

### 4.1.7 Mobility

The mobility of a mechanism is the total number of degrees of freedom provided by its members and joints.

In general a mechanism will have $n$ members of which $(n-1)$ provide degrees of freedom and $j$ joints which provide $C_i$ or $(6-f_i)$ constraints. This leads to the Kutzback criterion:

$$
\begin{aligned}
M &= 6(n-1) - \sum_{i=1}^{j} C_i \\
&= 6(n-j-1) + \sum f
\end{aligned}
\tag{4.1}
$$

where $\sum f$ is the total freedom of the joints.

Mobility is a useful concept in that it can give a warning of bad design. It does not, however, indicate a good design if the mobility worked out simply agrees with that which is expected.

The Kutzback criterion can be modified in the case of planar motion, that is all elements and movements are in one plane only. In this case it becomes Grubler's equation [25]

$$
M = 3(n-1) - 2j_1
\tag{4.2}
$$

where $n$ is the number of links including the base or earth and $j_1$ is the number of joints having one degree of freedom. This does not mean that only one joint can be at a certain point; more than one can be on the same shaft—in fact two are allowed without the mobility sum being violated. In cases where the mobility is being worked out it should be remembered that no two dimensions can be made to the same size, that is two links can never be the same length. It may be that, if two or three links are nearly the same, then some movement is possible, but this is only the result of an overconstraint (see figure 4.7). This figure is constrained to be

considered in the plane of the paper. Here $n = 5$, $f_1 = 6$, so $M = 3(4) - 12 = 0$. It seems, therefore, that the mobility is zero—which is absolutely true because the lengths of the links cannot be equal, meaning that each upper pin describes a different arc. In practice, however, some movement would be possible. In cases where m ≤ 0 common sense should be used. There are rarely practical cases.



**Figure 4.7** Body with zero degrees of freedom.

The case for linear and arcuate movement using leaf springs and their elastic bending capability is important. In particular, they have been used for the two movements fundamental to surface metrology: translation and rotation. Both can be provided using pivoted links although they are not used in practice. Obviously the simplest is that of rotation in which a simple pivoted link is used (figure 4.8).



**Figure 4.8** Single pivoted link.

There are quite complicated mechanisms, which were devised in the late nineteenth century for describing linear motion as well as arcuate motion in a plane. One such is shown in figure 4.9 due to Peaucellier [4]. In principle, this provides a linear movement. Application of Grubler's mobility test to this give a mobility of $M = 3(7) - 2(10) = 1$, there being eight links and ten joints.

The difficulty of manufacture and tolerancing really exclude such designs from consideration in precision applications but they may soon re-emerge in importance in micromechanics and microdynamics where such mechanisms could be readily fabricated.



**Figure 4.9** Peaucellier linkage—one degree of freedom.

### 4.1.8 Linear hinge mechanisms

In instruments measuring in the nanometre regime there is one type of motion that can be used which avoids using sliding mechanisms and their associated problems of assembly and wear. This is in the field of small translations or small arcuate movements. In these the bending of elastic elements in flexure mechanisms can be exploited (figure 4.10*a*).

Much of the early work on flexures was carried out by R V Jones (see e.g. [28]). Figure 4.10(*a*) shows a simple example of such a flexible system. Although the table moves in a slight arc, by careful symmetry of the legs, the movement can be made to be parallel to the base. The curved path is reduced in (*b*) which is the compound leaf spring and further reduced in (*c*) the symmetrical double-compound rectilinear spring.



**Figure 4.10** Leaf spring elements with normal one degree of freedom: *(a)* simple leaf spring movement; (*b*) compound leaf spring; (*c*) symmetrical double-compound rectilinear leaf spring.

Manufacturing errors will reduce the performance of such schemes when compared with the theoretical expectation. Parasitic movements are also caused by component parts having different tolerances, although use of monolith devices minimizes this effect. Such a case is shown in figure 4.11. The problem with the leaf spring approach is that it has a low stiffness out of the desired plane of movement. This can be reduced somewhat by splitting the spring.

The compound spring has half the stiffness compared with the simple spring in the drive direction and therefore twice the deflection for a given maximum stress. The monolith approximation to the leaf spring shown in figure 4.12 has a number of advantages. First, it is easy to manufacture, its stiffness being determined by drilling holes, a task well suited for CNC machines. Second, because it is made from a monolith it must already be assembled, so there can be no strain due to overconstraints. The notch hinge rotation as a function of the applied bending moment is

$$\frac{\theta}{M} = \frac{24KR}{Ebt^3} \tag{4.11}$$

**Figure 4.11** Split leaf spring hinge.

where $K = 0.565t/R + 0.166$. It can be seen that the effective stiffness of the notch is very much dependent on the thickness $t$ so that great care is needed in the drilling (figure 4.13).



**Figure 4.12** Notch hinge approximation to leaf spring.



**Figure 4.13** Notch hinge calculation: $\theta/M = 24KR/Ebt^3$.

In the case of the leaf spring the movement is only a fraction of the length of the leaf. The amount is governed by the geometry and the maximum allowable tensile strength $\sigma_{max}$ of the leaf material. For the simple spring mechanism of figure 4.10 the load drive is split between two springs, so the maximum moment in each will be $Fl/4$ if the load is applied in the line of the mid-point of the springs. Given that $\lambda$, the stiffness, is given by

$$\lambda = \frac{12EI}{l^3} \qquad \text{and} \qquad I = \frac{bh^3}{12} \tag{4.4}$$

the permissible displacement $\delta$ is

$$\delta = \frac{\sigma_{max} L^3}{3Eh}.$$ (4.5)

$\sigma_{max}$, and $E$ are material properties ($\sigma_{max}$ can usually be taken as 0.3, the yield stress for the metal), $L$ and $h$ are the geometrical parameters of interest.

This movement is usually rather a small fraction of $L$ (~ 0.1) and, as can be seen from equation (4.5), scales down rather badly with size. Consequently this type of motion, although having some advantages, does result in rather large base sizes for small movements.

Two-axis linear systems are also possible using the leaf spring approach. The system in figure 4.14 does have two degrees of 'limited' freedom. However, the drive axis is not stationary and so this can cause problems. Another possibility is the parallel kinematic design mechanics (PKM)



**Figure 4.14** Two-dimensional ligament hinge movement.

Such flexures used for two axes have the following stiffnesses [25]:

$$\lambda = \frac{Et^{7/2}}{5L^2 R^{1/2}} \qquad \text{for the simple spring}$$

$$\lambda = \frac{Et^{7/2}}{10L^2 R^{1/2}} \qquad \text{for the compound spring}$$ (4.6)

Joints of two degrees of freedom based on the notch are shown in figure 4.15.

(*a*)



(*b*)



**Figure 4.15** Two-dimensional notch-type hinges.

### 4.1.9 Angular motion flexures

For angles less than about 5° elastic designs may be used. The simplest is the short cantilever. This has an angular stiffness

$$\lambda_\theta = \sqrt{PEI}\ \cot\left(L\sqrt{P/EI}\right) \qquad \text{compressive}$$
$$\lambda_\theta = \sqrt{PEI}\ \coth\left(L\sqrt{P/EI}\right) \qquad \text{tensile}$$

(4.7)

The maximum angle is

$$\theta_{max} = \frac{\sigma_{max}^2 L}{Et}.$$

(4.8)

Simple cantilevers have limited use as pivots because their centre of rotation moves (figure 4.16). This is not nearly as pronounced in the notch hinge, which is why it is used in preference.



**Figure 4.16** Angular motion flexible hinge: $\theta_{max} = \sigma_{max}2L/Et$.

A common angular hinge is the crossed strip hinge and its monolith equivalent (figure 4.17). The other sort of angle system, conceived by Jones, is the cruciform angle hinge shown in figure 4.18.

(*a*)



(*b*)



**Figure 4.17** Angular crossed hinges: *(a)* crossed ligament hinge; *(b)* fabricated rotary hinge.

Summarizing the characteristics of flexure hinges:

Advantages:

    1. They do not wear—the only possibility is fatigue.
    2. Displacements are smooth and continuous.
    3. Displacements can be predicted.
    4. No hysteresis.



**Figure 4.18** Cruciform angular hinge.

Disadvantages:

    1. Small displacements require large mechanisms.
    2. Cannot tolerate high loads in compression, otherwise buckling occurs.
    3. The stiffness in the line of action of applied force is high and the cross-axial stiffnesses out of the plane of movement are relatively low.
    4. Can have hysteresis losses.

Nevertheless they have considerable use in instrument design.

### 4.1.10 Measurement and force loops

The notion of measurement loops and force loops is fundamental in the design of instruments for surface metrology. The force loop is a direct consequence of applying Newton's third law successively to different parts of a load-carrying static structure. For equilibrium, balanced forces must follow a continuous closed path often involving the ground (mechanical earth). Also, metrology loops come as a direct consequence of linking the object to be measured via a detector to the reference with which it is being compared (figure 4.19).

Because the force loop is necessarily under load it is therefore strained, so the dimensional stability is improved if it is small. Similarly the metrology loop is better if smaller because it is less likely to have a heat source or vibration source within it. Also, the smaller the loop the more likely it is that external heat sources or vibration sources or loads act across the whole of the loop. It is only when these act on part of the loop that instability and other errors occur.

The rules concerning loops are therefore:

    1. Any changes that occur to a part or parts within a measurement loop but not all of it will result in measured results that are not distinguishable from the measurement.

**Figure 4.19** Calliper principle.

2. What is of importance is that force loops should be kept as far as possible away from measurement loops and, if possible, the force loop should not be allowed to cross into or out of a measurement loop.

So a primary consideration in any design is to determine where the metrology and force loops interact and, if possible, to separate them. Usually the best that can be achieved is that they are coincident in just one branch of the measurement loop.

This makes sense because any strain in a metrology loop will cause a distortion which is an error, even if small. Fig (4.20) shows factors affecting the typical metrology loop.



**Figure 4.20**

*4.1.10*

*(a) Introduction*

There are a number of major problems to be overcome

(1) Design problems
(2) Improving accuracy
(3) Set-up problems and usage
(4) Calibration and traceability.

The basic questions that need to be addressed in any instrument or machine tool design are as follows, with particular reference to the metrology and force loops.

(1) Is it correct according to kinematic theory, that is, does it have the required number of constraints to ensure that the degrees of freedom conform to the movements required?
(2) Where are the metrology and force loops? Do they interact? Are they shared and, if so, how?
(3) Where are the heat and vibration sources and, if present, how can they influence performance?
(4) Is the design symmetrical and are the alignment principles obeyed?
(5) What other sources of error are present and can they be reduced by means of compensation or nulling?

*(b) Metrology loop properties*

The measurement loop must be kept independent of the specimen size if possible. In figure 4.21 it is not, however, figure 4.22 shows how it can be done.

In figure 4.21*(b)* and (c) the measurement loop is determined in size by the length of the specimen. Clearly it is beneficial from the point of view of stability to have the loop small. This can sometimes



**Figure 4.21** Measurement loop—size with respect to workpiece.

curtail the ease of measurement, but this problem can be resolved in a completely different way as shown in figure 4.22 below.

*Solution*



**Figure 4.22** Inverted instrument.

In the above solution, the measurement loop is independent of the size of the specimen. The measurement loop and also the force loop which moves the carriage etc should be kept small. This is easy to see if the effective stiffness of a cantilever, representing the opened up metrology loop in considered.



**Figure 4.23** Linear movement tube.

The small loop is effectively stiffer. The same force F produces a much smaller deflection, e.g. consider a cantilever with a force on the end and the same force in the middle. The deflection itself is proportional to $l^3$ where $l$ is the loop length. See equation 4.118.

Small loops usually are beneficial (figure 4.23), but sometimes make the operation of the instrument more difficult. Thin elements in the loop should be short.

Cross-sectional shape should be hollow, if possible, like a tube.



**Figure 4.24** Small metrology loop advantages.

Another point is that, if more than one transducer is in the metrology loop or touching it, then the transducers should, if possible, be the same. Such a situation arises if auxiliary transducers are used to monitor the change of shape of the main metrology loop e.g in scanning capacitative microscopes.

*(c) Other effects*

The larger the loop is the more likely it is to be influenced by effects within the loop, e.g. differentially by heat, whereas with a small loop, all of the loop heats up so the shape does not change; similarly with vibration.



**Figure 4.25**

*Important Note:* Surface metrology has one big advantage over dimensional metrology. This is that it does not matter if the size of the metrology loop increases as long as the shape is preserved (i.e. as long as the two arms of the calliper are in the same proportion). Consider the skid stylus system (figure 4.26).



**Figure 4.26**

The angle change $\delta\theta$ is the same for both systems.

### (d) Likelihood of source in loop
The smaller the loop the less likely heat or noise sources will be within it.

### (e) Symmetry
If there is a source within the loop it is not serious, providing that it is symmetrically placed relative to the two arms of the 'calliper'. The point is that any asymmetry has to be avoided.

### (f) Noise position
If a noise source such as a gearbox vibration is located it should be removed or absorbed as close to the source as possible otherwise the noise can dissipate into the system via many routes as shown in figure 4.27.



**Figure 4.27**

*(g) Co-ordinate system in instrument.*

The obvious way to configure a co-ordinate system for an instrument working in the nanometre region is to use Cartesian co-ordinates. This is probably always true if the probe is stationary and the workpiece is moved underneath it. However, some of the best results have been achieved with the probe moving and the workpiece stationary. This can be seen in figure 4.28.

The beauty of the arcuate design used for Talystep is that at every point in the arc of movement the accuracy is determined by one point. This point is the integrated position of line 0 0′. In any translation system the error in movement at any point is with respect to a corresponding point on the reference—no integration is involved.

Figure 4.28 shows a closer look at the design. No detailed explanation is necessary.

*(h) Improvement in accuracy – capability of improvement*

It is always the case that more is demanded of equipment/people/procedures than can reasonably be asked. One question often asked is whether, say, an instrument's performance can be improved retrospectively i.e. after purchase.

The quick answer is usually no, because the instrument has been designed to meet the specification and no more. It could be argued that over-design is bad design. However, requests for improvement will continue to be made.

One general rule of metrology is that it is not usually possible to make a bad instrument into a good instrument.

Good design improves the capacity to improve performance. Figure 4.29 shows the obvious fact that all instruments will drop in performance with time.



**Figure 4.28** (*a*) and (*b*)

(*c*)



**Figure 4.28** (*c*)



**Figure 4.29**

The limiting performance is determined by the precision or repeatability of the system. The stability of the instrument determines the repeatability of measurement. The stability also determines how often the instrument needs to be calibrated (figure 4.30).

Example of design philosophy:

Figure 4.31 shows what is, at first sight, a good design for a rotary bearing. The load of the quill is taken by male and female hemisphere, A and B. The resultant reaction acts through 0.



**Figure 4.30**



**Figure 4.31** Bad design.

Figure 4.32 shows a better alternative. The position of the hemispheres has been altered, (i.e. B and A reversed). The dynamic stability of the bearing has been improved.

Having established the fact that it is only possible to improve accuracy on an already good instrument poses the question of how to determine the systematic or repeatable error of the instrument. Once found it can be removed from all subsequent measurements. The general philosophy is shown in figure 4.33. Methods of achieving the estimation of the errors will be given shortly.



**Figure 4.32** Good design.



**Figure 4.33** Replacing the accuracy value of the spindle by its precision.

### 4.1.11 Alignment errors

If a measurement of distance is required it is fundamental that the measuring system is parallel to the axis along which the displacement is required and preferably collinear.

There is one basic concept that can be applied to loops of force or measurement and this is symmetry (or balance or matching). If opposite sides of a loop can be made symmetrical in material and in size then the chances are that it will be a good design. If a heat source or vibration source has to be in a measurement loop then arranging that the paths to the transducer are symmetrical will ensure that no differences are seen. Hence, heat sources or vibration sources should be placed at an axis of symmetry if possible. This would suggest that particularly good results will be obtained if the Abbé axis lies on an axis of symmetry of the stress distribution (e.g. the balance). This is illustrated in figure 4.34. Thus what is measured is $d'$, not $d$:

$$d' - d = d'(1 - \cos \theta) \qquad (4.9)$$

This is called the cosine error and it is concerned with 'angular mismatching' of the object to the measuring system.

### 4.1.11.1 Abbé offset

This is concerned not so much with angular mismatching but with displacement of the object to the measuring system. Thus the rule is:

When measuring the displacement of a specific point it is not sufficient to have the axis of the probe parallel to the direction of motion; the axis should also be aligned with the point, that is it should pass through it.

In figure 4.35 if there is a slight bow in the measuring instrument and its axis is misaligned by $l$ then it will read an error of $2l\theta$. If $l$ is zero then the same bow in the instrument would cause at most second-order or cosine errors. This sort of error due to the offset is called Abbé offset.

In its original form, as set out by Abbé in 1890, the alignment error referred to one dimension. However, it is clear that this oversimplifies the situation. There is a need for the accuracy of alignment in two and three dimensions, particularly for specific cases such as roundness and form measurement, etc. A number of people have investigated these principles and their attempts can be summarized by Zhang [29] who, working on the modification of the principle to straightness measurement made by Bryan [30], suggested the following definition: 'The line connecting the reference point and the sensing point should be in the sensitive direction.'



**Figure 4.34** Cosine error.

**Figure 4.35** (*a*) Abbé alignment error; (*b*) Abbé error in roundness (Zhang); (*c*) Abbé error—Vernier gauge (Zhang).

This definition applies to all cases of dimensional measurement including 1D, 2D and 3D measurements, straightness and roundness and even run-out measurements. The form of the error calculation is the same, namely $\delta = \Delta l \sin \psi$ where $l$ is the distance between the reference and $\psi$ is the Abbé angle.

Although this new formulation is correct it still has its problems. These are concerned with the identification of the reference point, sensing point and sensitive direction in complicated systems. Note that the basic idea of the calliper measurement system is maintained. Some cases are shown in figures 4.35(*b*) and (*c*).

### 4.1.12 Other mechanical considerations—balance of forces

Other design principles for instruments obviously exist but can be difficult to formulate. One is the idea of matching. If the instrument (or machine tool) provides a linear movement then there are three forces involved: the drive force, the frictional force and the inertial force. Ideally, these should balance along the same axis, otherwise it is inevitable that unwanted moments or torques will result and hence create twist in the system. Hence it is essential to ensure that the centres of inertia and friction are coincident with the drive axis. The centre of inertia can be found by calculation when the drive axis is known. What is more difficult to find is the centre of frictional forces. This is determined by the position and load carried by the support pads for the specimen or the pick-up (or tool) [26]. For example, figure 4.36 shows two alternative slideway configurations, (*a*) and (*b*). It is necessary to ensure that the drive system and the slideway axis are collinear.



**Figure 4.36** Balance of inertial, drive and friction forces.

The link between the drive motion and the slideway containing the probe or specimen should be non-influencing by having free or self-aligning nuts connecting the two members; then non-collinearity will not impose unwanted forces on the traversing mechanism.

Obviously the centre of inertia will depend on the size and shape of the specimen so that some provision should be made to ensure that the centre of inertia always remains on the same axis, for example by the provision of weights to the carriage.

### 4.1.13  Systematic errors and non-linearities

These always occur in any system and they conflict with the achievement of high accuracy. However, there are strategies which can be used to overcome them. Briefly, two of the most important are nulling and compensation.

Nulling is a method well known in the weight balance, for example. A working range which would carry the measurement operation into non-linear regions of the imput-output characteristic of the system is avoided by always ensuring that only one working point is used. Essentially, the reference in the 'calliper' system is continuously altered so that it always equals the test value. This is best embodied in the closed-loop 'follower' servo system. This is often used in optical scanners where the image of the surface is always kept in focus—the distance between the objective lens and the surface is a constant. Any error of focus is used to move the objective lens back into focus. The amount of movement is taken as the output signal. The point to note here is that it does not require a calibrated relationship between the out-of-focus and the position of the lens because the out-of-focus displacement characteristic is never needed, except at only one functional position.

Very often high accuracy is not achieved because of systematic errors in movements such as slideways etc. There are techniques that can be used to evaluate such errors and cancel their effect by compensation. Invariably the errors are evaluated by making more than one measurement and in some way changing the sense of the measurement of the test piece relative to the reference from the first measurement. In its simplest form this is known as 'reversal' (see figure 4.37).

Take the case of a test cylinder having unknown straightness errors $t(x)$, being measured with reference to a reference cylinder with unknown errors $R(x)$. On the first traverse the gap between them as seen by the invariable calliper is

$$S_1(x) \qquad \text{where} \quad S_1(x) = R(x) + t(x). \tag{4.10}$$

The test bar is rotated through 180° and the second measurement taken. Thus

$$S_2(x) = R(x) - t(x) \ .$$

From these $R(x)$ and $t(x)$ are found. Subsequently the error of the reference is stored and subtracted from any further measurement of test parts. The compensation method obviously holds if the errors are stable over long periods but is not much use if they change considerably with time, in which case the assessment of the reference error will have to be taken for each measuring process. Other alternatives to the reversal method exist, for example the error indexing method used by Whitehouse [31].

The compensation of clocks by using metals with different thermal coefficients of expansion is well known.

Another way to deal with errors is averaging. As an example, if the result from an experiment is subject to random unknown errors then simply averaging many readings will reduce the uncertainty. This is taken to a limit in the design of bearings. Sometimes the journal, for example, need not be especially round if it is separated from the shaft by a full oil film. The separation between the journal and shaft is an integration of the out-of-roundness over an arc, which is always much better than point to point.

(a) Test piece

(b)

$x_0$  $S_1(x)$  $x$

$x_0$  $S_2(x)$  $x$

Reference cylinder

**Figure 4.37** Error inversion method.

### 4.1.14 Material selection

This is a very wide ranging topic and is vitally important in the design of instruments. The criteria for instrument design are not necessarily the same as for other uses. For instance, some of the most used criteria, such as strength-to-weight ratio, are not particularly important in instrument design. In most cases material cost is not critical because the actual amount of specialized material used in any instrument design is relatively low. What does tend to be important are characteristics such as stiffnesses and thermal properties.

There have been many attempts to provide selection rules for materials [32-35] which use property groupings and charts to form a selection methodology. The properties of most importance are taken to be mechanical and thermal.

In mechanical properties a number of possibilities arise depending on the application. One of special note appears to be the ratio $E/\rho$ because self-weight deflection is minimized by taking a high value of $E/\rho$. Also, the resonant frequency of a beam of fixed geometry is proportional to $(E/\rho)^{1/2}$ and it is usual to aim for high resonant frequencies in an instrument system, $\rho$ is the density, $E$ is the elastic modulus.

There are other ratios such as $y/E$ for minimum strain or $y/\rho$ for any application involving high accelerations ($y$ is the yield strength). On balance, if one 'grouping' is to be chosen then $E/\rho$ seems best.

Thermal property groups again are numerous but, for those occasions where expansion should be kept at a minimum, Chetwynd suggests that the grouping of $\alpha/K$ should be low if the element is minimally constrained or $\alpha E/K$ if rigidly clamped.

Should there be thermal shock problems it makes sense to consider the behaviour of the thermal diffusion equation

$$\frac{\partial \theta}{\partial t} = \frac{K}{C\rho} \frac{\partial^2 \theta}{\partial x^2} \qquad (4.11)$$

The temporal behaviour of the diffusion of heat in the body is obviously governed by the grouping $K/C\rho$ (the diffusivity). For rapid response a high value would be preferable.

When these three groups, specific stiffness $E/\rho$, expansion/conduction $\alpha/K$ and diffusivity $K/C\rho$, are considered, they represent some attempt to rationalize the material properties for instrument design. In such groupings it is usual to plot them on a logarithmic scale with the values or reciprocals arranged so that high (or low) values are good (bad). Three graphs have been presented from which a good idea of material selection is possible and are shown in figure 4.38(*b*) from [33, 34]. Hence, if they are low value/high merit the best materials will be near to the origins of the graphs. These three graphs show that silicon is particularly good as a material, as is silicon carbide. An unexpectedly good one is beryllium— pity that it is a hazard to machine!

More comprehensive groupings have been used in the form of a material profile in which up to 11 properties or groupings are presented in a line (figure 4.38(*a*)). The idea is that a profile which indicates good behaviour of the material would lie consistently above the reference line (usually arbitrarily picked). Relative merits of materials are easily seen at a glance from such profiles. But even with 11 groupings the list is by no means exhaustive.

Some typical properties are given in table 4.1 for some well-known materials. This covers the materials in figure 4.38.



**Figure 4.38** *(a)* Material feature profile. *(b)* Material property guide: (i) material property map of expansion/conductivity versus inverse diffusivity; (ii) material property map of inverse specific stiffness versus inverse diffusivity; (iii) material property map of expansion/conductivity versus inverse specific stiffness. Key: Al, aluminium (duralumin is similar); $Al_2O_3$, alumina; Be, beryllium; Bra, 70/30 brass; Bro, 90/10 bronze; CI, cast iron; FQ, fused quartz; Inv, Invar; MS, mild steel; Si, single-crystal silicon; SiC, silicon carbide (reaction bonded); SiN, silicon nitride (reaction bonded); StSt, 18/8 stainless steel; ULE, ultra-low expansion glass ('Zerodur'); W, tungsten.

# Table 4.1

| Material | E modulus of elasticity (GPa) | $\rho$ density (kgm$^{-3}$) | $\alpha$ expansion coefficient ($\times 10^6$) | K thermal conductivity (Wm$^{-1}$K$^{-1}$) | $\beta$ thermoelastic coefficient ($\times 10^5$) | C specific list (J kg$^{-1}$ K$^{-1}$) | $D=K/\rho C$ diffusivity (m²/5 $\times 10^6$) | $E/\rho$ (m²s$^{-2}$ $\times 10^{-6}$) | $K/\alpha$ (Wm$^{-1}$ $\times 10^{-6}$) | Comments |
|---|---|---|---|---|---|---|---|---|---|---|
| Aluminium | 71 | 2710 | 23 | 200 | | 913 | 80 | 26 | 8.7 | Cheap, easy to machine |
| Alumina | 340 | 3700 | 8 | 21 | | 1050 | 5.4 | 92 | 2.6 | Hard/brittle |
| Beryllium copper | 126 | 8250 | 17 | 120 | −40 | −350 | 41 | 15 | 7.1 | Toxic in powder |
| Brass 70/30 | 100 | 8500 | 18 | 110 | −50 | | 35 | 12 | 6.1 | Cheap, easy to machine |
| Bronze 90/10 | 105 | 8800 | 17 | 180 | −40 | | 56 | 12 | 10.6 | Cheap, easy to machine |
| Copper | 130 | 8954 | 16.6 | 386 | −50 | 383 | 112 | 14.5 | 23 | Cheap, easy to machine |
| Diamond | 1200 | 3500 | 1.2 | 1590 | | 510 | 890 | 342 | 1325 | Very hard and can be made in fibres |
| Elinver/Ni SpanC(CrFc) | 170 | 8000 | 4 | 10 | 0 | 460 | 2.7 | 21 | 2.5 | Moderately expensive |
| Fused quartz | 70 | 2200 | 0.5 | 1.5 | 10 | 840 | 1.0 | 34 | 2.7 | Hard, brittle and cheap |
| Fused silica | 70 | 2150 | 0.5 | 2.1 | 10 | | | | | Hard, brittle and cheap |
| Graphite | 4.8 | 2250 | 2 | 23.9 | | 691 | 15 | 2.1 | 12 | |
| Imar (36% Ni-Fe) | 150 | 8000 | 0.5 | 10.7 | −50 | 460 | 2.9 | 18.8 | 21 | Moderately expensive and machinable |
| Magnesium | 41.3 | 1740 | 26 | 153 | | 1020 | 86 | 24 | 5.9 | |
| Molybdenum | 325 | 10200 | 5 | 138 | | 251 | 54 | 32 | 27.6 | Difficult and toxic in powder form |
| Steel—mild | 210 | 7860 | 15 | 63 | −30 | 420 | 19 | 27 | 4.2 | |
| Spring steel | 205 | 7860 | 11.5 | 63 | −20 | 400 | 20 | 26 | 5.5 | |
| SiC | 410 | 3100 | 4.3 | 200 | | 1040 | 62 | 132 | 47 | Hard, brittle |
| SiN | 165 | 2500 | 3.2 | 150 | | 710 | 8.5 | 66 | 4.7 | Hard, brittle |
| Steatite Si ceramic | 110 | 2700 | 8.5 | 2.9 | | | | 40.7 | 0.3 | Hard, brittle |
| Syalon Si ceramic | 300 | 3950 | 3–5.6 | 20 | | 630 | 7.6 | 76 | 5 | Hard, brittle |
| Tungsten | 345 | 19300 | 4.5 | 166 | | 133 | 65 | 18 | 37 | Difficult to machine |
| Titanium | 110 | 4500 | 8.9 | 21.6 | | 528 | 9.1 | 24 | 59 | |
| ULE Titanium Silica | 70 | 2210 | 0.05 | | | | | 32 | | Hard, brittle |
| Zeradur/ Marium glass | 91 | 2530 | 0.05 | 1.64 | | 821 | 0.8 | 36 | 33 | Hard, brittle |
| Carbon fibre | 150 | 1580 | −0.7 | | | | | 95 | | |
| Silicon (pure) | 190 | 2300 | 2.3 | 157 | | 706 | 97 | 83 | 68 | Hard, brittle and very pure |
| Beryllium | 280 | 1850 | 11 | 209 | | 1675 | 67 | 151 | 19 | Toxic |
| Wood (pine) | 15 | 500 | 5 | 0.15 | | 280 | 1 | 30 | 0.03 | Very workable |

## 4.1.15  Drive systems

It is not within the scope of this book to consider the methods of driving the movements required for instruments because these depend greatly on developments in material technology and are forever changing, as one option becomes temporarily superior to another. So far the emphasis has been on describing matters of principle rather than specific detail. However, a brief mention here will be a reminder that the force drive system is an integral and important consideration in all instruments.

Piezoelectric devices are the most common method of driving over short ranges of a few micrometres. Axial strain is the usual mode, which provides very stiff drives. Bending modes can produce larger ranges for a given voltage but obviously at lower stiffnesses. One classic design is used in the scanning tunnelling microscope discussed later in which a split piezoceramic tube is used to provide all three movements $x, y$ and $z$ (figure 4.83).

Piezoelectric devices have a serious disadvantage in that they have high hysteresis and the materials are not too stable over long periods, as well as being sensitive to temperature. In general, linearities of only 1% are possible, so they are usually worked out as part of a closed-loop servo system in order to get the range/resolution to a reasonable level of better than 1 in $10^3$. Chetwynd [25] has pointed out that stiff drives can have advantages and disadvantages. One advantage is that, because of the high stiffness of the drive, displacement is more or less independent of the resistance of the rest of the system. On the other hand it transmits all vibrations directly to the platform or carriage upon which the probe is supported or the specimen sits. Thus, when flexures are used, there may be some benefits in using a more compliant drive. This has been the reason for much research but it turns out that a very acceptable drive can be provided by means of a magnetic-coil force actuator.

With modern ceramic magnets very compact actuators can now be built. It should be remembered that it does not require much force to move any structure by a few micrometres. Also, since the actuator provides force it need not be closely positioned near to the moving platform; it can be connected by a push rod. Two or more actuators can also act in series, if required, to provide a coarse and fine movement, for example.

Flexure mechanisms have been extensively used to produce movement for surface measuring instruments but they have to be driven. This is sometimes done using piezoelectric crystals. These are used because of their very high stiffness. Piezoelectric elements can move the worktable or the probe of the instrument. One big problem is the limited range and poor strain outputs. Strain value of 0.1% is typical (i.e. 100 mm crystal stack can produce only 100 $\mu$m deflection), but the force produced is high.

Unfortunately the PZT (lead zirconate titanate) crystal—typical piezo crystal—suffers from hysteresis and non-linearity. The only way to get satisfactory movement and correct position for the instrument is by using closed loop control. The scan be achieved by using optical interferometry [36].

Another approach uses capacitive position sensors which detect the crystal movement. For a good review of these methods consult [37]. For subnanometre positional accuracy the capacitive method is proving to be satisfactory.

Although piezoelectric drives will continue to be used for most applications involving scanning microscopes it may be that a magnet-driven silicon monolith offers a better system for nanotechnology. These systems will be larger and less convenient but the actual movement would be more accurate because the lattice spacing could be used as an interferometer using a system described in chapter 5.

So far if movements of greater than 1 mm are used then it seems inevitable that kinematically designed slideways will be the answer. One interesting possibility for a slide would be of zero thermal expansion glass, such as Zerodur, coated with a monolayer of polymer similar to a Langmuir-Blodgett coating for low friction. For smaller movements the piezoelectric systems could be used and in between these extremes the flexure-type movements would be acceptable.

In what has been considered so far the dynamic effects of instrument systems have not been mentioned. These will follow after a discussion of the probe system and in particular the mechanical tactile probe used so extensively in surface metrology.

## 4.2 Measurement systems

### 4.2.1 Aspect of general stylus systems

A system is shown in figure 4.39. Each of the blocks in the figure is important in its own right. However, the issue of traceability and reliability of data will be considered later in chapter 5. This section will start off by looking at various mechanical methods of measuring surfaces and establishing some limitations on their use. It could be argued that the whole of the book should have started off with instrumentation rather than parameters, using the argument that the parameters do not exist if they cannot be measured. The reason for not pursuing this path is simply that instrumental techniques do change, whereas the need to define and specify parameters does not, so here only the basic techniques will be examined without reference to individual instruments.

The simplest way to think of any surface-measuring instrument is to imagine it as a calliper of some sort as shown in figure 4.40.

In basic terms, the sensor is the device which picks up the information from the test piece (workpiece), another part of the sensor looks up a reference, the transducer establishes the difference between the two and



**Figure 4.39** Measurement system.



**Figure 4.40** Instrument calliper principle: *(a)* remote reference flat and smooth; *(b)* remote reference to test pieces; *(c)* intrinsic reference—a blunt stylus.

converts the difference information into an electrical signal by means of what is in effect an angle transducer (many gauges in surface metrology are side acting as shown to enable inaccessible parts to be measured). The calliper ABC moves relative to the reference and workpiece and the variations in angle are detected and amplified. Usually the signal is digitized and then processed.

Three features need to be examined here: the sensor picking up the information—usually a stylus or optical equivalent—the reference and the means for comparison. This last point needs amplification. If points A and B move vertically as the pick-up arrangement is moved across the surface, other vertical movements in the calliper cannot readily be tolerated—the mechanical loop ABC has to be stable. In figure 4.41, the mechanical loop coupling the reference and workpiece has to be rigid, at least within the frequency range capable of being picked up by the transducer. The reference here could be a fixed reference within the instrument, such as the precision spindle in roundness measurement. It could be an improvized reference, such as is shown in figure 4.40(c), in which the mechanical loop is artificially 'earthed' by making both parts of the sensor contact the workpiece, or it could be part of the frame of the machine tool or instrument.

In the case of figure 4.40(c) the reference is developed via a blunt sensor (more about this later).

Consider first the primary sensor, that which is picking up information from the workpiece directly. It is on this that the measurement most depends. In many instruments it is a solid stylus. Some of its properties will now be considered. The stylus requirement for scanning probe microscopes which are not measuring topography e.g. tunneling current, is somewhat more involved, as will be seen.

### 4.2.2 Stylus characteristics

#### 4.2.2.1 Tactile considerations

In tactile sensors and transducers there are four aspects to the stylus and transducer mechanism:

1. tip dimension
2. the angle
3. stylus force
4. elastic/plastic behaviour.

#### (a) Tip dimension

Consider the mechanical stulus first. It is generally a pyramid or conical diamond terminating with a flat or rounded tip. The preferred forms employ a 90° pyramid with a 2–2.5 $\mu$m flat or a 60° cone with a 12.5 $\mu$m. radius spherical tip. A tip radius of 2–2.5 $\mu$m is an option. The tip dimension will influence the value of the surface parameters obtained. The curvature of the stylus adds to that of the surface asperities and subtracts from the valleys. In effect, at any point along the traverse there is a very complicated interaction between the stylus and the geometry being measured. This effect has been investigated in detail by Agullo and Pages-Fita [38].



**Figure 4.41** Mechanical earths in calliper system.

A more straightforward approach for random surfaces is given in figure 4.42 [39]. The spherical stylus can be considered to be modelled by a three-point *possibility* of contact. For a flat stylus tip only two points need be considered. This does not mean that there can be three points of contact. Obviously there can only be one because the pick-up system has only one degree of freedom, either an arc for side-acting gauges or vertical linear for in-line gauges. It means that selecting three possible points of contact is adequate for defining the shape and size of the stylus if it is spherical. Two will suffice if the stylus has a flat tip.



**Figure 4.42** Approximate models for stylus integration: (*a*) three-point and (*b*) two-point model.

Assuming that the surface is random in the sense that it has a Gaussian height distribution and an exponential autocorrelation function, the height distribution traced by the stylus in its passage across the surface can be found. Quite obviously it will no longer be Gaussian.

From figure 4.42 for a spherical stylus tip of radius $R$ the effective tip dimension is

$$d = \frac{2sR}{\sqrt{1+s}} \tag{4.12}$$

where $s$ is the slope of the flanks. For most styluses $s$ can be taken as about unity. Therefore $d = R\sqrt{2}$.

Simple geometry shows that the circular tip shape can be specified by the three possible points separated by $R/2$ and the centre one depressed by $R(1 - \sqrt{\frac{1}{2}})$.

The probability density $p_s(z)$ becomes

$$p_s(z) = \int_{-\infty}^{z+k} \int_{z}^{z+\delta z} \int_{-\infty}^{z+k} p(z_{-1}, z_0, z_{+1}) \mathrm{d}z_{-1} \mathrm{d}z_0 \mathrm{d}z_{+1} \tag{4.13}$$

which results in the rather messy expression

$$
\begin{aligned}
p_s(z,\rho) = & \exp\left(\frac{[-(z+k)^2/2]}{2\sqrt{2\pi}}\right)\left[1 + \mathrm{erf}\left(\frac{-z}{\sqrt{2}}\sqrt{\frac{1-\rho}{1+\rho}} - \frac{k\rho}{\sqrt{2(1-\rho^2)}}\right)\right] \\
& \exp\left(\frac{[-(z+k)^2/2]}{2\sqrt{2\pi}\sqrt{1-\rho^2}}\right)\int_{-\infty}^{z} \exp\left(\frac{[z_0 - \rho(z+k)^2]}{2(1-\rho^2)}\right)\mathrm{erf}\left(\frac{(z+k-\rho z_0)}{\sqrt{2(1-\rho^2)}}\right)\mathrm{d}z_0 \\
& + \frac{1}{4\sqrt{2\pi}}\exp\left(-\frac{z^2}{2}\right)\left[1 + \mathrm{erf}\left(\frac{z}{\sqrt{2}}\sqrt{\frac{1-\rho}{1+\rho}} + \frac{k}{\sqrt{2(1-\rho^2)}}\right)\right]^2
\end{aligned} \tag{4.14}
$$

**Figure 4.43** Correlation between stylus size and independence distance of surface; stylus integration, three-element Markov chain MO simulating spherical stylus.

where $\rho$ is the correlation between the adjacent values of the three points and $k = R(1 - \sqrt{\frac{1}{2}})$, $z$ being the height on the surface at which the stylus is resting. Equation (4.14) is a true density because it integrates to unity.

Thus, given a knowledge of the autocorrelation function for the surface (obtained with a sharp stylus) the RMS and approximate $R_a$ of the same surface measured with a tip of radius $R$ can be found by taking the square root of the second central moment of equation (4.14) as shown in figure 4.43.

The results from the graph show that, for a typical ground surface, the loss in value of the $R_a$ value is likely to be about 2% using a 2 $\mu$m stylus. However, for a 10 $\mu$m stylus the loss is considerably higher.

#### (b) Stylus angle

To get an idea of the problem concerning the slope of the flank of the stylus consider the effect on the same type of surface using the formula for the average slope

$$
m = \frac{1}{h}\left(\frac{1 - \rho^2}{\pi}\right)^{1/2}
\tag{4.15}
$$

where $\rho$ is as above—the correlation between possible points of contact. For any significant integration caused by the stylus angle (assuming an infinitely sharp tip) the half angle of the stylus would have to be smaller than the surface slopes (30° or 45°). This implies from equation (4.14) that the spacings $h$ between



**Figure 4.44** Stylus pressure.

ordinates ($R/\sqrt{2}$) would have to be of the order of 0.05 $\mu$m. The spacings, however, for a typical ground surface are 20–50 times greater than this in practice, so the tip dimension is much more likely to produce the integrating or smoothing effect than the stylus flank angle.

Practical results have been obtained by Radhakrishnan [40].

The section above has considered the integrating effect of the stylus tip dimension, its shape and the slope of the flanks. Another important issue is the possibility of damage caused by the pressure at the tip. This will be considered next.

*(c) Stylus pressure*

Consider the case of a spherical stylus first. Imagine it to be contacting a workpiece as shown in figure 4.44

According to Hertz's theorem the radius of the contact zone $a$ due to elastic deformation of the two touching surfaces under a load $W$ is given by

$$a = \left[ \frac{3\pi}{4} \left( \frac{1-v_1^2}{\pi E_1} + \frac{1-v_2^2}{\pi E_2} \right) \frac{Wr_1r_2}{r_1 + r_2} \right]^{1/3} \tag{4.16}$$

where $v_1$, $v_2$ are the Poisson ratios of the materials and $E_1$ and $E_2$ their elastic moduli. Equation (4.16) and results for gramophone records were obtained earlier by Walton [41].

The compliance (the amount they press together) is given by

$$\delta = \frac{a^2}{r_1 r_2 / (r_1 + r_2)}. \tag{4.17}$$

Obviously the pressure over the area $\pi a^2$ is

$$W / \pi a^2. \tag{4.18}$$

This expression can be considered to be the "elastic" pressure. This will be considered in detail later on when the dynamic component of W is evaluated and the stylus damage index defined.

Also, the maximum pressure is 3/2 times the average pressure assuming an elliptical pressure distribution within the contact. Thus

$$P_{\max} = \frac{3}{2} \frac{W}{\pi a^2}. \tag{4.19}$$

Church [42] calculated the differences in vertical deformation between the cases when the stylus contacted the surface at a peak compared with when it contacted at a valley $\Delta\delta$ Thus

$$A\delta = \left[ \frac{3\pi}{4} \left( \frac{1-v_1^2}{\pi E_1} + \frac{1-v_2^2}{\pi E_2} \right)^2 \left( \frac{W^2}{R_T} \right)^{1/3} \left( 1 + \frac{R_T}{R_s} \right)^{1/3} - \left( 1 - \frac{R_T}{R_s} \right)^{1/3} \right] \tag{4.20}$$

where $R_s$ is the radius of curvature of the surface and $R_T$ that of the stylus.

For the case where $R_T = R_s$ the peak deformation becomes a maximum and

$$A\delta \sim \frac{3\pi}{4} \left( \frac{1-v_1^2}{\pi E_1} + \frac{1-v_2^2}{\pi E_2} \right)^2 \left( \frac{2W^2}{R_T} \right)^{1/3} \tag{4.21}$$

For $R_s \gg R_T$, the usual case, equation (4.21) is multiplied by

$$\frac{2^{2/3}}{3} \frac{R_T}{R_s}.$$

(4.22)

This deformation difference vanishes where $R_s$ (r) $\rightarrow \infty$, since then the peak and valley deform equally. Some practical values are

$$a \simeq 45(WR)^{1/3} nM.$$

So for $W = 0.7$ mN (the typical force is about 70 mg) and $R = r_1 r_2/(r_1 + r_2) = 2000$ nM, $a$ becomes equal to about 0.2, $\mu$m, giving a pressure of about 2500 N mm$^{-2}$, which is less than the yield pressure for most materials except perhaps for soft coatings.

Similar results are obtained with a stylus having a flat tip whose nominal dimension is 2 $\mu$m square. It should be noted that, in practice, the tip is usually longer in one direction (~ 10 $\mu$m) than another to allow more mechanical strength. In these circumstances the user should be aware of the fact that it is the longest dimension which is the effective integrator when an isotropic surface is being measured.

The more difficult case for deformations on hollow cylindrical workpieces has also been derived [43].

*(d) Elastic/plastic behaviour*

One important point needs to be made concerning the apparently large pressures exerted by styluses on the surface. The criterion for damage is not that this pressure exceeds the yield stress of the material being measured, it is that the pressure does not exceed the skin yield stress of the material and not the bulk yield stress. Unfortunately these are not the same. If the hardness of the material is plotted as a function of the depth of indentation as in figure 4.45, the hardness increases with decreasing depth. This increase in hardness at the skin is not due to work hardening because the same effect occurs when measuring the hardness of indium, which anneals at room temperature. This apparent increase has been attributed to the fact that at very light loads the size of the indentation is small compared with the average distance between dislocations, so the material approaches its theoretical hardness. For this reason, in many cases where the stylus should damage the surface, no damage is discernible.

Also, the presence of compliance does not in itself distort the picture of the surface. Providing that (i) it is reasonably small compared with the texture being measured and (ii) iy is more or less constant, an output of high mechanical fidelity can be obtained.

Typical tips are shown in figure 4.46. Although at first sight these two types of stylus shape appear to be very different, in practice they are not. Usually what happens is that a spherical tip develops a flat on it with use and the flat tip gets rounded edges. Hence they tend to converge on each other in shape with wear.

Most styluses are made out of diamond. To get the longest life out of the tip, if possible the diamond should have its (111) crystallographic axis parallel to the plane of the workpiece and hence movement; other planes are considerably softer. For the measurement of very fine surface roughness it is possible, with patience, to get a stylus tip down in size to smaller than 0.1 $\mu$m. However, under these circumstances the stylus force has to be considerably reduced. A typical force used in these applications is 1 milligram.

In the discussion above concerning loading it should be noted that, if the surface has been produced by polishing or similar finishing, there will always be a scale of size of asperity much smaller than the stylus (i.e. fractal) so that the curvatures will necessarily be high and the compliance will be so great that plastic flow will inevitably result. However, such marks are very small 'in energy' compared with the major finishing scratches and can be neglected. The same applies to fractal surfaces.

**Figure 4.45** Scale of size effect on hardness.



**Figure 4.46** Typical styluses.

#### 4.2.2.2  *Pick-up dynamics*

The next problem that has to be considered is whether the stylus tracks across the surface faithfully and does not lift off. This is sometimes called 'trackability'. In what follows the approach used by surface instrument-makers [44] will be followed. The bent cantilever as found in SPM instruments will not be discussed here but late in the section. Reference [45] shows a simple high speed method.

It is well known that high effective mass and damping adversely affect the frequency response of stylus-type pick-ups. However, the exact relationship between these parameters is not so well known.

This section shows how this relationship can be derived and how the pick-up performance can be optimized to match the types of surfaces that have to be measured in practice.

The dynamics of pick-up design fall into two main areas:

1. frequency response
2. mechanical resonances that give unwanted outputs.

These will be dealt with separately.

##### (a)  *Frequency response*
The frequency response of a stylus system is essentially determined by the ability of the stylus to follow a surface. That is, the output does not fall at high frequencies, although the stylus will begin to mistrack at

lower amplitudes. This concept, termed 'trackability', was coined (originally) by Shure Electronics Ltd to describe the performance of their gramophone pick-ups. A figure for trackability usually refers to the maximum velocity or amplitude for which the stylus will remain in contact with the surface at a specified frequency. A more complete specification is a graph plotting the maximum trackable velocity or amplitude against frequency. Remember that the frequency will in fact be a spatial frequency of the surface traced over the pick-up speed.

*(i) Trackability—sine wave*
The theoretical trackability of a stylus system can be studied by looking at the model of the system and its differential equation (figures 4.47–4.49). Here $M^*$ is the effective mass of the system as measured at the stylus (see later), $T$ is the damping constant (in practice, air or fluid resistance and energy losses in spring), $K$ is the elastic rate of spring, $F$ is the nominal stylus force (due to static displacement and static weight) and $R$ is the upward reaction force due to the surface. In alternative pick-up system for high speed tracking using a stylus system suited to areal mapping is discussed in section 4.2.5.6.

Trackability is usually quoted in terms of velocity, so initially the system will be looked at in terms of velocity. If the surface is of sinusoidal form, its instantaneous velocity will also be of sinusoidal form; hence



**Figure 4.47** Schematic diagram of suspension system.



**Figure 4.48** Model of system.



**Figure 4.49** Side-acting stylus pick-up.

the instantaneous velocity $v = a \sin \omega t$. Now, the equation representing the system is given in differential form by

$$M * \ddot{z} + T\dot{z} + Kz + F = -R(t) \tag{4.23}$$

or in terms of velocity

$$M * \dot{v} + Tv + K\int v\,\mathrm{d}t + F = -R(t). \tag{4.24}$$

The first analysis will be in terms of the amplitude of velocity and the condition for trackability. Hence

$$R(t) = -M * a \cos \omega t - Ta \sin \omega t + K\frac{a}{\omega}\cos \omega t - F$$
$$= a\left[\left(\frac{K}{\omega} - M * \omega\right)\cos \omega t - T \sin \omega t\right] - F \tag{4.25}$$

$$R(t) = a\left\{\left[\left(\frac{K}{\omega} - M * \omega\right)^2 + T^2\right]^{1/2} \sin(\omega t - \varphi)\right\} - F$$
$$\varphi = \tan^{-1}\left(\frac{T\omega}{K - M * \omega^2}\right) \tag{4.26}$$

but for zero trackability $R(t) = 0$ and hence

$$\mid a \mid = \frac{F}{M * \omega_0(\omega_0^2/\omega^2 + \omega^2/\omega_0^2 + 4\zeta^2 - 2)^{1/2}} \tag{4.27}$$

where $|a|$ is the amplitude of the maximum followable sinusoidal velocity of angular frequency $\omega$. (Note that this analysis only assumes sinusoidal waves on the surface.) Here $\omega_n = 2\pi\times$ undamped natural frequency of system, that is

$$\omega_n = 2\pi f_n = (K/M*)^{1/2}$$

and the damping factor (ratio) is

$$\zeta = \frac{T}{2\omega_n M *}.$$

This can be simplified further to

$$\mid a \mid = \frac{A_N K}{M*\omega_n}\left(\frac{\omega_n^2}{\omega^2} + \frac{\omega^2}{\omega_n^2} + 4\zeta^2 - 2\right)^{-1/2}$$
$$\mid a \mid = A_N \omega_n\left(\frac{\omega_n^2}{\omega^2} + \frac{\omega^2}{\omega_n^2} + 4\zeta^2 - 2\right)^{-1/2} \tag{4.28}$$
$$= A_N \omega_n C$$

where $A_N$ is the nominal stylus displacement. The factor $C$ is shown plotted for various values of $\omega/\omega_n$ in figure 4.50.

*(ii) Interpretation of figure 4.50*
Any pick-up will have a fixed value of $\omega_n$ and $\zeta$. Hence, from the family of curves, the value of $C$ can be found for any ratio of $\omega/\omega_n$ from 0.1 to 10. Thus

**Figure 4.50** Response of system.

$$\frac{\omega}{\omega_n} = \frac{f_n V}{\lambda} \tag{4.29}$$

where $V$ is the stylus tracking velocity and $\lambda$ the wavelength of the surface. Having found a value of $C$, then the peak stylus velocity for which the stylus will stay in contact with the surface is given by

$$a = 2\pi f_n A_N C.$$

This can also be written in terms of slope:

$$S = \frac{2\pi f_n A_N C}{V} \tag{4.30}$$

where $S$ is the maximum slope for which the stylus will stay in contact with the surface at tracking velocity $V$.

Expression (4.30) is the usual expression used by manufacturers of instruments.

*(iii) Trackability in terms of amplitude*

Equation (4.23) can be solved for a displacement which is of a sinusoidal form. This gives the result which can be directly connected with the velocity expression of equations (4.24) and (4.27):

$$\mid a \mid = A_N \frac{\omega_n}{\omega} \left| \frac{\omega_n^2}{\omega^2} + \frac{\omega^2}{\omega_n^2} + 4\zeta^2 - 2 \right|^{1/2}. \tag{4.31}$$

This curve is plotted for different values of $\omega_0$ ($A_N$ and $\zeta$ fixed at 0.25 mm and 0.2 respectively).

It will be noticed that the trackability at long wavelengths is only determined by the nominal stylus displacement. For example, if a stylus has a nominal displacement of 0.1 mm to obtain the correct stylus pressure, then it will be able to track sine wave amplitudes of 0.1 mm.

There is also a noticeable peak in the curve, which is due to the resonant frequency of the stylus suspension. This peak is of a magnitude

$$A_{pk} = \frac{A_N}{2\zeta}. \tag{4.32}$$

(Note that for zero damping this peak is infinite.) The position of the peak is also definable as

$$\lambda_{pk} = \frac{V}{f_n(1-\zeta^2)^{1/2}} \approx \frac{V}{f_n}.$$

(4.33)

Figure 4.51 shows the relationship between $f_n$ and $\lambda_{pk}$ for various $V$.

These simple equations allow a linear piecewise construction of the trackability curve to be easily drawn for a particular pick-up, providing that the resonant frequency, damping ratio, nominal stylus displacement and tracking velocity are known.

The resonant frequency can be calculated from the effective mass $M^*$ and spring rate $K$, that is

$$f_0 = \frac{1}{2\pi}\left(\frac{K}{M^*}\right)^{1/2}.$$

(4.34)

The damping ratio $\zeta = T/2\omega_n M^*$ (where $T$ is force/unit velocity). Unfortunately $T$ is rarely known, but since in the conventional case low damping is required, $\zeta$ can be approximated to 0.2 since in practice it is difficult to obtain a value significantly lower than this. There are occasions, however, where fixed values of $\zeta$ are required (as will be seen shortly) which are much larger than 0.2.

### (iv)  Application of instrument, trackability criterion to sinusoidal surfaces

For any surface the amplitude spectrum must lie at all points below the trackability curve. Figure 4.51 shows some typical trackability curves in relation to an area termed 'real surfaces'. The upper boundary represents a wavelength-to-amplitude ratio of 10:1, which is the most severe treatment any pick-up is likely to receive. The lower boundary represents a ratio of 100:1, which is more typical of the surfaces to be found in practice.

It is worth noting here that if a surface has a large periodic component, this can be centred on the resonant frequency to improve the trackability margin for difficult surfaces (see figures 4.51–4.54).

### (b)  Unwanted resonances in instruments for metrology

There are two principal regions where resonances detrimental to pick-up performance can occur (figure 4.55): (i) the pickup body and (ii) the stylus beam. These resonances are detrimental in that they have the



**Figure 4.51** Trackability of system.

**Figure 4.52** Tracking velocity—resonant frequency.



**Figure 4.53** Linear approximation to trackability curve.

effect of causing an extra displacement of the stylus with respect to the pick-up sensor (photocells, variable reluctance coils, etc).

*(i) Pick-up body*
The pick-up body can vibrate as a simple beam clamped at one end (figure 4.56).
    In the schematic diagram of the pick-up (figure 4.55) the resonant frequency is given by

$$f_{\mathrm{n}} = \frac{\pi}{8L^2} \sqrt{\frac{EI}{A}} \tag{4.35}$$

for the fundamental. There are only two ways in which the effect of this resonance can be reduced:

    1.  By increased damping or arranging the sampling to minimize effects to white-noise inputs.

**Figure 4.54** Wavelength for maximum trackability.



**Figure 4.55** Schematic diagram of pick-up.



**Figure 4.56** Vibration mode of pick-up body.

2. By designing the body such that the resonance frequency is well outside the frequency range of interest. This, the standard approach, will be discussed first.

This resonant frequency can be calculated for beams with one end free or both clamped, or where the beam is solid or hollow, circular or rectangular in cross-section, or whatever.

*(ii) Stylus beam*
It is evident that the stylus beam can resonate to give an unwanted output, and again this can be removed by damping or shifting the frequency out of the useful range. Since the stylus beam is of a uniform shape the resonant frequency can be reliably predicted providing the mode of vibration can be defined.

However, this is difficult to define. A beam with one end clamped and the other free is unlikely, since the stylus is presumably always in contact with the surface. On the other hand, the boundary conditions for a

clamped beam do not hold either, since the slope at the beam ends is non-zero. Hence, if $f_n$ is the resonant frequency for a free beam then $f'_n = 9f'_n$ (where $f'_n$ is the resonant frequency for a clamped beam) and the resonant frequency will probably be somewhere between $2f_n$ and $9f_n$, since there will be nodes at both ends.

*Example*

$f_n$ (one end free) was calculated at 576 Hz (measured at 560 Hz).
$f_n$ (both ends clamped) was calculated at 5.18 kHz.
Hence the true resonant frequency will be greater than 1.12kHz, and may be as high as 5kHz.

*(iii)   Stylus Damage Prevention Index [74]*
The criterion for surface damage can be given in terms of the hardness $H$ of the test material and the elastic pressure $P_e$. It is the imbalance between these two pressures that gives rise to damage.

Thus the comparison is between $H$, the Vickers or DPN (diamond pyramid number) and $P_e$. Vickers or Knoop hardness are preferred because they most nearly imitate the effect of the diamond onto the surface. Let $\psi_{st}$ be the damage index.

Then

$$P_e = \psi_{st} H \tag{4.36}$$

If $\psi_{st} > 1$ then $P_e$ is bigger than the yield stress of the material so damage occurs. For $\psi_{st}$ ( 1 only elastic deformation takes place.

So

$$\psi_{st} = \frac{P_e}{H} = \frac{W}{\pi a^2 H} \tag{4.37}$$

where $a$ is the elastic radius on the surface and $W$ the load. This has two components one static $F$ and one dynamic $WD$. Using the values of $W$, $a$ and $H$ an estimation of $\psi$ can be made [74].

In the formula (4.37) the values of $W$ and $H$ need to be understood. Consider first $W$. As indicated above

$$W = F - W_D \tag{4.38}$$

The criterion for trackability is that $W$ is never negative (i.e. $F \geq W_D$). The minimum value of $W$ therefore is $-F$ or $-W_D = F$, i.e. $W = 0$.

This occurs at a peak where the stylus is being accelerated by the surface yet being held on by the stylus spring and gravity (see figure (4.57).) If this criterion holds it automatically sets the criterion for the maximum reaction at the surface. This is when the retardation of the stylus system is a maximum which is when the stylus is in a valley. In this case $W_D$ adds to $F$ so that the maximum reaction is $W = 2F$.

The important factor here is that it is *in the valleys* where the stylus is likely to cause damage *not at the peaks*.

Damage at the peaks is due to the use of a skid which only touches the peaks.

$$\psi_{\text{stylus}} = \frac{P_{e\,st}}{H} = \frac{k}{H}\left(\frac{2FR}{E}\right)^{-2/3} \tag{4.39}$$

$$\psi_{\text{skid}} = \frac{P_{e\,sk}}{H} = \frac{k}{H}\left(\frac{W'r}{E}\right)^{-2/3} \tag{4.40}$$

**Figure 4.57** Stylus Damage Index

In equation (4.39) $R$ is the stylus radius (Figure 4.57(a)), whereas in (4.40) it is the local peak radius $r$ on the surface (figure 4.57(b)). Also the reaction due to the stylus dynamics is at most equal to the static force $F$ i.e. $W = 2F$. However, the skid does not have a big dynamic component: if it is integrating properly it should be zero. This does not make $W'$ equal to $F$ because $W'$ is the load of the whole traversing unit, which is much greater than F. Notice that for potential damage due to the skid it is the surface peak curvature that is relevant and not the skid radius (which can be assumed to be very large). $k$ is a constant dependent on Poisson's ratio.

If the damage index for the stylus is $\psi_{st}$ and that for the skid is $\psi_{sk}$ then

$$\psi_{sk} = \psi_{st}\left(\frac{W'}{2F}\right)^{\frac{1}{3}} \cdot \left(\frac{R_{st}}{r_s}\right)^{\frac{2}{3}}$$

(4.41)

where $W'$ is the load on skid, $R_{st}$ is stylus radius and $r_S$ is average asperity radius.

So knowing that $W' \gg 2F$ and $R_{st} \leq r_s$

$$\psi_{sk} > \psi_{st}$$

(4.42)

In other words, the skid is much more likely to cause damage than the stylus. Furthermore skid damage is at the peaks and any stylus damage is in valleys. Which of the two possibilities is causing the damage can therefore be identified. Also the skid is usually offset laterally from the stylus track to ensure that any skid damage does not affect the profile obtained by the stylus.

For cases in surface investigation where the dynamic reaction $W_D$ is needed, for example in friction measurement, the normal force due to the geometry is $W_D$ (see section 4.2.2.4) where

$$W_D = 3R_q M w_n^2 \left[1 + \left(4\zeta^2 - 2\right)\left(\frac{w}{w_n}\right)^2 + \left(\frac{1}{\varepsilon}\right)^2\left(\frac{w}{w_n}\right)^4\right]^{\frac{1}{2}}$$

(4.43)

Where $M$ is the effective mass of the pick-up (see next section), $w_n$ the resonant frequency of the pick-up and $\zeta$ the damping term. $\varepsilon$ represents the type of surface. For $\varepsilon \to 0$ the surface becomes more random as in grinding. $R_q$ is the RMS roughness.

Notice that surface parameters and instrument parameters both contribute to the measurement fidelity.

In equation (4.37) the $W$ term has been explained. The other parameter $H$ is more difficult because it is material processing dependent. It is not the bulk hardness but the 'skin' hardness which is different (see figure 4.45). Unfortunately the index $\psi$ is very dependent on $H$; far more than on the load as can be seen in equation 4.44.

$$\psi = \frac{k}{H} W^{1/3} \cdot \left(\frac{R}{E}\right)^{-2/3}$$

(4.44)

### 4.2.2.3 Conclusions about mechanical pick-ups of instruments using the conventional approach

An expression for the trackability of a pick-up has been derived [44,45] and from this a number of criteria can be inferred for a good pick-up:

1. The resonant frequency, comprising the stylus effect mass and the suspension spring rate, should be as high as possible.
2. The nominal stylus displacement should be as high as possible.
3. The pick-up damping ratio should be as low as possible within certain constraints, although in practice any value less than unity should be satisfactory and, as will be seen, the value of 0.6 should be aimed for.
4. To eliminate unwanted resonances the pick-up body and stylus beam should be as short and as stiff as possible.

Some of these conditions are not compatible, for example 1 and 2. Condition 1 ideally requires a high spring rate, whereas condition 2 requires a low spring rate to obtain an acceptable stylus pressure. Similarly, condition 4 is not consistent with the measurement of small bores. However, a compromise could probably be obtained by using a long stylus beam to obtain the small-bore facility and keeping the pick-up body short and stiff. The beam resonance must of course be kept out of the range of interest unless the surface is completely random.

Note 1: effective mass of a beam
In the inertial diagram of a stylus beam (figure 4.58) the effective mass is given by

$$M^* = M + \frac{mh^2}{L^2} + \frac{aL}{3} + \frac{ah}{3} \times \frac{h^2}{L^2}$$

(4.45)

where $a$ = mass/unit length. If $L/h=n$ then $L^2/h2 = n^2$, that is

$$M^* = M + \frac{m}{n^2} + \frac{aL}{3} + \frac{ah}{3n^2} = M + \frac{m}{n^2} + \frac{ah}{3}\left(\frac{n^3 + 1}{n^2}\right).$$



**Figure 4.58** Inertial diagram of stylus beam.

The effect of increasing $L$ is $(L + \Delta L)/h = n'$ and hence

$$M_2^* = M + \frac{m}{n'^2} + \frac{ah}{3}\left(\frac{n'^3 + 1}{n'^2}\right).$$

Therefore

$$M^* = M + \frac{mh^2}{L^2} + \frac{aL}{3} + \frac{ah}{3} \times \frac{h^2}{L^2} \qquad (4.46)$$

As an example, for $n = 2$ and $n' = 3$

$$\Delta M^* = \frac{ah}{3}(0.86) - m(0.14).$$

Note 2: resonant frequency of a beam (figure 4.59)



**Figure 4.59** Resonant frequency of beam.

It can be shown that the small, free vertical vibrations $u$ of a uniform cantilever beam are governed by the fourth-order equation

$$\frac{\partial^2 u}{\partial t^2} + c^2 \frac{\partial^4 u}{\partial x^4} = 0$$

where $c^2 = EI/\rho A$ and where $E$ is Young's modulus, $I$ is the moment of inertia of the cross-section with respect to the $z$ axis, $A$ is the area of the cross-section and $\rho$ is the density. Hence it can be shown that

$$F^{(4)}/F = \ddot{G}/c^2 G = \beta^4 \quad \text{where } \beta = \text{constant}$$
$$F(x) = A \cos \beta x + \beta \sin \beta x + C \cosh \beta x + D \sinh \beta x \qquad (4.47)$$
$$G(t) = a \cos c\beta^2 t + b \sin c\beta^2 t.$$

Now, if the beam is clamped at one end it can be shown that

$$\beta L = \frac{\pi}{2}(2n + 1) \quad \text{where } n = 0, 1, 2\dots. \qquad (4.48)$$

Hence, if the beam resonates in its fundamental mode

$$\beta = \frac{\pi}{2L} \quad \text{and} \quad \omega_n = \frac{c\pi^2}{4L^2}$$

and

$$f_n = \frac{1}{2\pi} \frac{\pi^2}{4L^2} \left( \frac{EI}{\rho A} \right)^{1/2} = \frac{\pi}{8L^2} \left( \frac{EI}{\rho A} \right)^{1/2}$$

### 4.2.2.4 Relationship between static and dynamic forces

The condition for stylus lift-off has been considered in equations (4.27) and (4.30). These considerations have tended to deal with the behaviour of the system to sinusoidal inputs. Only on rare occasions do such inputs happen. More often the input is random, in which case the situation is different and should be considered so. In what follows the behaviour of the pick-up system in response to both periodic and random inputs will be compared. None of this disputes the earlier analysis based on the trackability criterion $R = 0$; it simply looks at it from a different point of view, in particular the properties, and comes up with some unexpected conclusions. First recapitulate the periodic case [44].

#### (a) Reaction due to periodic surface
Using the same nomenclature as before for an input $A \sin \omega t$:

$$AM\omega_n^2 \left\{ \left[ 1 - \left( \frac{\omega}{\omega_n} \right)^2 \right]^2 + 4\zeta^2 \left( \frac{\omega}{\omega_n} \right)^2 \right\}^{1/2} \sin(\omega t + \varphi) + F = -R_v(t) \qquad (4.49)$$

where

$$\varphi = \tan^{-1} \left( \frac{2\zeta(\omega/\omega_n)}{1 - (\omega/\omega_n)^2} \right).$$

This shows that the reaction is in phase advance of the surface geometry. Hence energy dissipation can occur.

For convenience the dynamic amplitude term can be rewritten as

$$\widetilde{R}_{max}(t) = AM\omega_n^2 \left[ 1 + (4\zeta^2 - 2)\left( \frac{\omega}{\omega_n} \right)^2 + \left( \frac{\omega}{\omega_n} \right)^4 \right]^{1/2} = A'. \qquad (4.50)$$

Hence, because the input is a sine wave, comparing $A$ with $A'$ represents a transfer function.

What are the implications of this equation? One point concerns the behaviour at resonance. Putting $\omega = \omega_n$ shows that the dynamic force is $2AM\omega_n^2\zeta$, which is obviously zero when $\zeta = 0$. The dynamic component is zero because the stylus is in synchronism with the surface. The force is not zero, however, because the static force $F$ is present. This means that, even if the damping is zero and the system is tracking at the resonant frequency, it will still follow the surface. This in itself suggests that it may be possible to go against convention and to track much nearer the resonant frequency than was otherwise thought prudent. Another point concerns the situation when $\omega > \omega_n$. Equations (4.49) and (4.50) show that the forces are higher when $\omega$ gets large; the force is proportional to $\omega^2$ for $\omega > \omega_n$ and hence short-wavelength detail on the surface is much more likely to suffer damage.

Points such as these can be picked out from the equations of motion quite readily, but what happens in practice when the surface is random? Another question needs to be asked. Why insist on dealing with the system in the time domain? Ideally the input and output are in the spatial domain. It seems logical to put the system into the spatial domain also.

*(b) Reaction due to the tracking of random surfaces [47]*

For periodic waveforms other than a sinusoid the Fourier coefficients are well known and the total forces can be derived by the superposition of each component in amplitude and the phase via equation (4.49) because the system is assumed to be linear. However, for a random surface this is not so easy: the transformation depends on the sample. For this reason the use of random process analysis below is used to establish operational rules which allow the dynamic forces to be found for any surface using parameters derived from initial surface measurements. Using random process analysis also allows a direct link-up between the dynamic forces and the manufacturing process which produces the surface.

The inputs to the reaction are made up of $\ddot{z}(t)$, $\dot{z}(t)$ and $z(t)$. For a random surface these can all be assumed to be random variables. It is usual to consider these to be Gaussian but this restriction is not necessary. If these inputs are random then so is $R_v(t)$.

The mean of $R_v(t)$ is $F$ because $E[\ddot{z}(t)] = Ez[(t)] = E[z(t)] = 0$ and the maximum value of $\widetilde{R}_{v,}(t)$ can be estimated from its standard deviation or variance. Thus

$$
\begin{aligned}
\mathrm{var}(-\widetilde{R}_v(t)) &= E[(M\ddot{z}(t)+T\dot{z}(t)+kz(t))^2] \\
&= M^2\sigma_c^2 + T^2\sigma_s^2 + k^2\sigma_z^2 + 2\{MTE[\ddot{z}(t)\dot{z}(t)] \\
&\quad + MkE[\ddot{z}(t)z(t)] + TkE[\dot{z}(t)z(t)]\}
\end{aligned}
\tag{4.51}
$$

or

$$
\sigma_R^2 = M^2\sigma_c^2 + T^2\sigma_s^2 + k^2\sigma_z^2 - 2Mk\sigma_s^2.
\tag{4.52}
$$

Equation (4.51) results because it can be shown that all cross-terms are zero except $E[\ddot{z}(t)z(t)]$, which is equal to $-\sigma_s^2$, and where $\sigma_z^2$, $\sigma_s^2$, $\sigma_z^2$ are the variances of $\ddot{z}(t)$, $\dot{z}(t)$ and $z(t)$ respectively.

Letting $k/M = \omega_n^2$, $T/M = 2\zeta\omega_n$ and $\sigma_z^2 = R_q^2$ (to agree with the international standard for the standard deviation of surface heights)

$$
\begin{aligned}
\sigma_R &= R_q M\left[\left(\frac{\sigma_c^2}{\sigma_z^2}\right) + 4\zeta^2\omega_n^2\left(\frac{\sigma_s^2}{\sigma_z^2}\right) + \omega_n^4 - 2\omega_n^2\left(\frac{\sigma_s^2}{\sigma_z^2}\right)\right]^{1/2} \\
&= R_q M[v^2\omega^2 + (4\zeta^2-2)\omega^2\omega_n^2 + \omega_n^4]^{1/2}
\end{aligned}
\tag{4.53}
$$

or

$$
\sigma_R = R_q M\omega_n^2\left[1 + (4\zeta^2-2)\left(\frac{\omega}{\omega_n}\right)^2 + \frac{v^2}{\omega^2}\left(\frac{\omega}{\omega_n}\right)^4\right]^{1/2}
\tag{4.54}
$$

where

$$
\omega^2 = \frac{\sigma_s^2}{\sigma_z^2} = \left(\frac{2\pi V}{\lambda_q}\right)^2 \quad \text{and} \quad v^2 = \frac{\sigma_c^2}{\sigma_s^2} = \left(\frac{2\pi V}{\overline{\lambda}_q}\right)^2
\tag{4.55}
$$

In relationship (4.55) $\lambda_q$ is the average distance between positive zero crossings and $\overline{\lambda}_q$ that of peaks.

Taking $3\sigma_R$ as the typical $\widetilde{R}_{\max}(t)$ value and writing $v/\omega = 1/\varepsilon$ yields the dynamic force equation for a random surface:

$$
\widetilde{R}_{\max}(t) = 3R_q M\omega_n^2\left[1 + (4\zeta^2-2)\left(\frac{\omega}{\omega_n}\right)^2 + \left(\frac{1}{\varepsilon}\right)^2\left(\frac{\omega}{\omega_n}\right)^4\right]^{1/2}.
\tag{4.56}
$$

If $3R_q$ for the random wave is taken to be equivalent to $A$ for a periodic wave, then equation (4.56) corresponds exactly to the expression given in equation (4.50) except for the term $(1/\varepsilon)^2$ which is a measure of the randomness of the surface. *The term $\varepsilon$ is a surface characterization parameter.* It is also, therefore, a characterization of the type of reaction.

Equation (4.56) allows the dynamic component of reaction for any type of surface to be evaluated. For $\varepsilon = 1$ equation (4.56) reduces to equation (4.50) and is true for a periodic wave. Strictly, $\varepsilon = 1$ corresponds to the result for a sine wave or any deterministic wave, although the dynamic characteristics for such surfaces will be different. This difference is related more to the rate of change of reaction, rather than the reaction force. In the next section this will be considered. However, the primary role of the $\varepsilon$ parameter is in distinguishing random from periodic surfaces and distinguishing between random surfaces. As $\varepsilon$ tends more and more to zero the surface becomes more random until when $\varepsilon = 0$ the surface is white noise; any type of random surface can be characterized but $\varepsilon$ cannot take values larger than unity. As examples of intermediate surfaces, one which has a Gaussian correlation function has a value of $\varepsilon = 0.58$ whereas one having a Lorenzian correlation function of $1/(1 + \beta^2)$ has a value of $\varepsilon = 0.4$. Note here that, according to equation (4.56), as $\varepsilon \to 0$, $\tilde{R}_{max} \to \infty$ In practice this cannot occur because $\tilde{R}_{max}$ is curtailed at the value $F$ when stylus lift-off occurs.

### (c) Statistical properties of the reaction and their significance: autocorrelation function and power spectrum of R(t)

The autocorrelation function is given by $E[R(t)R(t + \beta)]$ where $R(t)$ is defined as earlier.

Removing the $F$ value and taking expectations gives $A_R(\beta)$ where $A_R(\beta)$ is given by

$$A_R(\beta) = M^2 \omega_n^4 \left[ A_z(\beta) - \left( \frac{1}{\omega_n} \right)^2 \frac{d^2}{d\beta^2} (A_z(\beta))(4\zeta^2 - 2) + \left( \frac{1}{\omega_n} \right)^4 \frac{d^4(A_z(\beta))}{d\beta^4} \right]. \tag{4.57}$$

It can be shown that the odd terms in the evaluation disappear because the autocorrelation function is an even function. In equation (4.57) neither $A_R(\beta)$, the autocorrelation function of the reaction, nor $A_z(\beta)$, that of the surface profile, is normalized. The normalizing factor $A(0)$, the variance of this vertical component of reaction, is given by

$$A(0) = M^2 \omega_n^2 \left( \sigma_z^2 + \frac{\sigma_s^2}{\omega_n^2} (4\zeta^2 - 2) + \frac{\sigma_c^2}{\omega_n^4} \right). \tag{4.58}$$

It is clear from equation (4.57) that there is a difference between $A_R(\beta)$ and $A_z(\beta)$. This difference can be seen more clearly by reverting to the power spectral density $P(\omega)$ derived from the autocorrelation function and vice versa. Thus

$$A(\beta) = \frac{1}{\pi} \int_0^\infty P(\omega) \cos\omega\beta \, d\omega. \tag{4.59}$$

Taking the Fourier spectrum of both sides, squaring and taking to the limit gives $P_R(\omega)$ in terms of $P_z(\omega)$. Hence

$$P_R(\omega) = P_z(\omega) \left[ 1 + \left( \frac{\omega}{\omega_n} \right)^2 (4\zeta^2 - 2) + \left( \frac{\omega}{\omega_n} \right)^4 \right]. \tag{4.60}$$

The expression in square brackets, designated $H(\omega)$, is the weighting factor on $P_z(\omega)$ to produce $P_R(\omega)$ (see figure 4.60). Later $H(\omega)$ is examined to decide how it needs to be manipulated to reduce surface damage. Equation (4.60) does not contain a term involving $\varepsilon$, despite the fact that it refers to random profiles, because only individual wavelengths in the spectrum are considered, whereas in equation (4.56) $\omega$ is related to the *average* wavelength over the whole spectrum. The latter of course does depend on $\varepsilon$.

The simple nature of H($\omega$) allows a comparison to be made of the spatial frequency content of $P_R(\omega)$ relative to $P_z(\omega)$ but it is not so clear from equation (4.57) using the autocorrelation function.

### (d) System properties and their relationship to the surface: damping and energy loss

One effect of damping is to dissipate energy. This energy will be lost in the system but the heating effects could just conceivably permeate through to the stylus thereby affecting friction and calibration in the new generation of surface instruments. The question arises as to how the type of surface affects the energy loss.

The energy loss $J$ is given by

$$J = \int R(z)\,dz. \tag{4.61}$$

For the system $z$ is a function of $x$ and thereby $t$ if a constant tracking velocity is assumed. Only the dynamic component of the reaction will be considered. Thus

$$J = \int_0^{2\pi/\omega} R(z(t))\frac{d(z(t))}{dt}\,dt = AA'\pi\sin\varphi \tag{4.62}$$

where $\varphi$ and $A'$ are given by equations (4.49) and (4.50) respectively.

It will be remembered that the reaction is of the same wavelength as the periodicity but out of phase (equation (4.49)). Letting $A^2/2 = R_q^2 = \sigma_z^2$ (the variance of the surface) the energy loss per unit cycle becomes

$$J = 4R_q^2\pi M\zeta\omega_n\omega. \tag{4.63}$$

This equation contains all the system parameters and two surface parameters, $R_q$ for amplitude and $\omega$ for spacing.

For a random wave equation (4.61) gives

$$J = \frac{1}{L}\int_0^L (M\ddot{z} + T\dot{z} + kz)\dot{z}\,dt. \tag{4.64}$$

Evaluating (4.64) and letting $T/M = 2\zeta\omega_n, k/M = \omega_n^2, \sigma_s^2 = 1/L\int\dot{z}^2 dt, \sigma_z^2 = \sigma_s^2\omega_q^2, L = \lambda_q$ and $\omega_q = 2\pi/\lambda_q$ gives the energy loss $J$ per unit *equivalent cycle* $\lambda_q$:

$$J = 4\sigma_z^2\pi M\zeta\omega_n\omega_q \equiv 4R_q^2\pi M\zeta\omega_n\omega_q. \tag{4.65}$$

Therefore $J$ is not dependent on $\varepsilon$ so that all surfaces, whether periodic or random, having the same $R_q$ and $\omega_q$ will cause the same heat losses. These losses are likely to affect the behaviour of the system rather than influence the surface. Heating effects at the surface will more probably be influenced by the frictional effects between the stylus and the surface.

### (e) Integrated damping: system optimization for random surface

Consider first the contact case. What is required is a measuring system that does not degrade the surface during measurement; if it does it should have a minimum, constant or predictable interaction. The system also, obviously, has to convey the geometry of the surface to the transducer with maximum fidelity. Both points have been discussed in isolation but not relative to the system of measurement.

There are at least two ways of attempting to optimize the system from the point of view of damage and fidelity. One is to reduce the reaction; another is to make the reaction match the surface. Of these the obvious way to reduce the magnitude of the reaction *from the point of view of damping* is to reduce the $T$ term in the

equation of motion. This has the effect of reducing the damping ratio $\zeta$, and is an obvious move. However, making $\zeta \rightarrow 0$ introduces problems at the transducer end of the system if very sharp features are to be measured. It also gives the system a long oscillatory settling time. Sensitivity to shocks and long settling times are system features which should be avoided. Increasing $\zeta$ to a large value is obviously counterproductive because of the increase in $R$ and the sluggish response. What is required is an optimum value of $\zeta$ to improve fidelity and/or reduce damage. One approach is to consider the shape of $H(\omega)$ given in equation (4.60). This is plotted for various damping ratios as a function of $\omega/\omega_n$ in figure 4.60.

So far the discussion has centred on instruments which make contact with the surface. However, there are differences in approach between contact and non-contact systems. The difference exists because of different criteria for performance. For contact instruments lack of surface damage and high fidelity are paramount but, for non-contact methods, speed and fidelity are paramount. Contact criteria for bandwidth and damping will be considered in what follows.

Because random surfaces have a spread of frequencies contained in them attempts to optimize $H(\omega)$ must encompass a band of frequencies. Also, pulses have a wide spectrum. Here the maximum bandwidth of the system is taken to be up to when $\omega = \omega_n$ for reasons given below. One criterion for optimization could be to make $H(\omega) = 1$ over the band of interest. The classical way to do this is to make $\omega_n \gg \omega_n$, where $\omega_u$ is the highest frequency on the surface, but this is not always possible. In fact with the tracking speeds required today $\omega_u$ invariably approaches $\omega_n$ in value, it being difficult to make $\omega_n$ sufficiently high by stiffening or miniaturizing the system. A further reason why this option is becoming difficult to achieve with modern contact instruments is the need to measure shape and form as well as texture with the same instrument. This requires the instrument to have a large dynamic range, which means that the spring rate $k$ is kept low in order to keep the surface forces down for large vertical deflections. Having a low $k$ value invariably requires that $\omega_n$ is small. The implication therefore is that $\omega \leqslant \omega_n$ for tactile instruments because of the potential damage caused by the forces exerted when $\omega > \omega_n$, as seen in equation (4.55). If $\omega_n$ is low, as is the case for wide-range instruments, then the only practical way to avoid damage is to reduce the tracking speed $V$, thereby reducing



**Figure 4.60** Optimum damping factor.

$\omega_u$ and consequently the forces. However, instead of insisting that $\omega_n$ is made high (the effect of which is to make $H(\omega) \simeq 1$ for each value of $\omega$) there is an alternative criterion which relaxes the need for a high value of $\omega_n$ by utilizing all the band of frequencies up to $\omega_n$ In this the damping ratio $\zeta$ *is* picked such that

$$\int_0^{\omega_n} H(\omega)\mathrm{d}\omega = 1 \tag{4.66}$$

so, although $H(\omega)$ is not unity for each individual value of $\omega$, it is unity, on average, over the whole range. *This has the effect of making $A_{yr}(\beta) = A_y(\beta)$ which is a preferred objective* set out in Chapter 5.

From figure 4.60 it can be seen that the shape of $H(\omega)$ changes dramatically with damping. The curve of $H(\omega)$ has a minimum if $\zeta < \frac{1}{2}$. This is exactly the same criterion for there to be a maximum in the system transfer function of equation (4.88). It also explains the presence of the term $(4\zeta^2-2)$ in many equations. Evaluating equation (4.66) shows that for unit area $\zeta = 0$–59. With this value of damping the deviation from unity of the weighting factor is a maximum of 39% positive at $\omega = \omega_n$ and 9% negative for a $\omega = 0.55\omega_n$. The minimax criterion gives $\zeta = 0.54$ with a deviation of 17% and the least-squares criterion gives a $\zeta$ value of 0.57. All three criteria produce damping ratios within a very small range of values! Hence if these criteria for fidelity and possible damage are to be used the *damping ratio should be anywhere between 0.59 and 0.54.* Within this range the statistical fidelity is assured.

It is interesting to note that the interaction between the system parameters and the surface can be taken a step further if the properties of the rate of change of reactional force are considered. Thus

$$R'_v(t) = \frac{\mathrm{d}}{\mathrm{d}t}[M\ddot{z}(t) + T\dot{z}(t) + kz] \tag{4.67}$$

from which, by the same method of calculation and letting $\sigma_c^2/\sigma_s^2 = v^2, \sigma_j^2/\sigma_t^2 = \eta^2$

$$\overline{R}_{\max}(t) = 3M\sigma_s\omega_n^2\left[1 + (4\zeta^2 - 2)\left(\frac{v}{\omega_n}\right)^2 + \left(\frac{1}{\gamma}\right)^2\left(\frac{v}{\omega_n}\right)^4\right]^{1/2} \tag{4.68}$$

where $\gamma = v/\eta$ and $\eta = 2\pi/\lambda_j$ and $\lambda_j$ is the average distance between positive points of inflection of the surface profile; $\eta = 2\pi/\overline{\lambda}_q$ where $\overline{\lambda}_q$ is the average distance between peaks $(\sigma_j^2 = E(\ddot{z}(t)^2))$.

As before, putting $3\delta_s = A'$ for the derivative of the periodic wave and $\gamma = 1$ yields the corresponding formula for the rate of change of reaction for a periodic wave.

In equation (4.68) $\gamma$ is another surface characterization parameter of one order higher than $\varepsilon$. It seems that this is the first time that the third derivative of the surface has been used to describe surface properties—this time with respect to the instrument system.

### 4.2.2.5 Alternative stylus systems and effect on reaction/random surface (*figure 4.61*)

For low damage risk $R$ should be small. To do this F should be made small. If it is made too small the reaction becomes negative and the stylus lifts off. The signal therefore loses all fidelity. Making the maximum dynamic component of force equal to $-F$ ensures that this never happens and that the maximum reaction possible is $2F$. For minimum damage overall, therefore, $F$ and $M\ddot{z} + T\dot{z} + kz$ should be reduced together and in parallel.

A step in this reduction is to remove $kz$, which is the largest component of force. This means that the system acts as a gravity-loaded system in which $F = mg$ where $m$ is mass and $g$ is the acceleration due to gravity. In this simplified system the range of 'following' frequencies is given by

$$f \leqslant k\frac{1}{2\pi}\sqrt{\frac{g}{a}} \tag{4.69}$$

**Figure 4.61** Alternative systems.

where $k$ is a constant which depends on the constituent parts of the pick-up and how they are connected; $a$ is the amplitude of the surface signal.

Furthermore, if $T$ is removed using air bearings or magnetic bearings, a very simple force equation becomes

$$R(t) = F + M\ddot{z}. \tag{4.70}$$

To get the height information from such a system does not mean that the signal has to integrate twice. The value of $z$ could be monitored by non-intrusive means, such as an optical method that locates on the head of the stylus.

Equation (4.70) implies a transfer function proportional to $1/M\omega^2$ which drops off with a rate of 12dB per octave.

If such systems are considered then they can be compared with the full system over a bandwidth up to $\omega_n$ as for $H(\omega)$.

### 4.2.2.6 Criteria for scanning surface instruments

There are a number of criteria upon which an instrument's performance can be judged. These are (laterally and normal to the surface) range/resolution, speed or response, and, fidelity or integrity. Cost is left out of this analysis.

The range of movement divided by the resolution is a key factor in any instrument as it really determines its practical usefulness. At the present time in surface instruments it is not unreasonable to expect ratio values of $10^5$ or thereabouts. One restriction previously encountered in the $x$ and $y$ lateral directions was due to the relatively coarse dimensions of the probe. This has been largely eliminated by the use of styluses with almost atomic dimensions. The same is true of optical resolution limitations.

The ability to manufacture these very sharp styluses has resurrected their importance in instrumentation. At one time the stylus technique was considered to be dated. Today it is the most sophisticated tool! This is partly due to a swing in emphasis from height measurement, which is most useful in tribological applications, to lateral or spatial measurements, which are more important in applications where spacing or structure is being investigated, namely in the semiconductor and microelectronics industries and in biology and chemistry.

A need for lateral information at nanometre accuracy has resulted in an enhanced requirement for the fidelity of measurement. At the same time the speed of measurement has to be high to reduce the effects of the environment. Speed has always been important in surface measurement but this is particularly so at the atomic level where noise, vibration and thermal effects can easily disrupt the measurement. Speed of measurement is straightforward to define but fidelity is not so simple. In this context fidelity is taken to mean the degree to which the instrument system reproduces the surface parameter of interest. As has been seen, failure to achieve high fidelity can be caused by the probe failing to follow the surface characteristics owing to inertial or damping effects or, alternatively, making a contact with such high pressure that the surface is damaged.

There is always a fine balance to be achieved between fidelity of measurement and speed. If the speed is too high fidelity can be lost, for example owing to surface damage in topography measurement. If it is too low then the environmental effects will destroy the fidelity.

At the atomic level distinctions between what is mechanical and what is electrical or electronic become somewhat blurred but, nevertheless, in the microscopic elements making up the measuring instrument the distinction is still clear.

In what follows some alternatives will be considered. Central to this analysis will be the basic force equation representing the measuring system. It consists essentially of a probe, a beam or cantilever, a transducer and means for moving the probe relative to the specimen.

### 4.2.2.7 Forms of the pick-up equation

#### (a) Temporal form of differential equation

The differential equation representing the force equation is linear and of second order. It is given by equations (4.71) and (4.72):

$$M\ddot{z} + T\dot{z} + kz + F = R(t) \tag{4.71}$$

or

$$z(t)(D^2 M + TD + k) + F = R(t) \tag{4.72}$$

where $M$ is the effective mass of the probe relative to the pivot, $T$ is the damping term, $k$ is the rate of the equivalent spring, $F$ is the static force and D is the differential operator on $z$, the vertical movement of the probe. This equation is present in one form or another in all surface-measuring instruments. How it is interpreted depends on whether force or topography is being measured.

Equation (4.71) can also be seen from different points of view. In force measurement $R$ is the input and $z$ the output, but for surface damage considerations $z$ is the input and $R$, the reaction at the surface, the output.

It is often convenient to decompose equation (4.71) into a number of constituent parts:

**Table 4.2** This table shows the relationships between the dynamic modes of pick-ups.

| Atomic force measurement | | Relevant equation and operations | Topographic measurement issues | | |
|---|---|---|---|---|---|
| Fidelity | Speed | | Damage | Speed | Fidelity |

Temporal form

$$M\ddot{z} + T\dot{z} + kz + F = R \quad \text{(equation (4.71))}$$

$\dot{z} = \text{constant}$ (left)          Coefficient truncation          $\dot{z} = \text{constant}$ (right)

$$V(dy/dx) = W$$
(equation (4.102))

$T, k = 0$          $V\dfrac{dz}{dx} = W$

$$R_i(t) \propto z_0(t)$$

$$M\ddot{z} = R(t) - F$$

minimum reaction $\begin{pmatrix} z \text{ input} \\ R \text{ output} \end{pmatrix}$

$$M \propto \tfrac{l}{v}, \quad T \propto \tfrac{l}{v}$$
(equation (4.97))

Statistical form

$$M^2 \omega_n^4 A_z(\beta)\left[ 1 - (4\zeta^2 - 2)D^2\left(\frac{1}{\omega_n}\right)^2 + \left(\frac{1}{\omega_n}\right)^4 D^4 \right] = A_R(\beta)$$

(equation (4.57))

$$M^2 \omega_n^4 P_z\left(\frac{\omega}{\omega_n}\right)\left[ 1 - (4\zeta^2 - 2)\left(\frac{\omega}{\omega_n}\right)^2 + \left(\frac{\omega}{\omega_n}\right) \right]$$

(equation (4.60))

Statistical characterization

$\dfrac{\bar{v}}{\omega} = 0 \to 1$

$$\left[ 1 + (4\zeta^2 - 2)\left(\frac{\bar{\omega}}{\omega_n}\right)^2 + \left(\frac{\bar{v}}{\bar{\omega}}\right)^2 \left(\frac{\bar{\omega}}{\omega_n}\right)^4 \right]$$ (equations (4.54), (4.81))

$\bar{\omega} = 0 \to 1$

Random periodic          Random-periodic          Coefficient integration

(equation (4.86))

Spatial resolution enhanced to $\lambda_n$ (left)

Spatial resolution enhanced to $\lambda_n$ (right)

$$\int \left[ 1 + (4\zeta^2 - 2)\left(\frac{\omega}{\omega_n}\right)^2 + \left(\frac{\omega}{\omega_n}\right)^4 \right] d\left(\frac{\omega}{\omega_n}\right) = 1$$

Integrated damping

Spatial form

$$MV^2(x)z_o(x)[D^2 + V(x)D(MDV(x) + T) + k] + F = R(x) \quad \text{(equation (4.83))}$$

Differential spatial damping

$$\zeta(x) = \left( \zeta + \frac{1}{2\omega_n}\frac{du}{dx} \right) \text{(equations (4.106), (4.85))}$$

Constant velocity          Variable velocity          Variable velocity          Constant velocity

$V \ll 1$
(equation (4.99))

$\zeta = 0.59$          $\zeta = 0.59$          $V \ll 1$

$$V(x)\frac{dz}{dx} = \text{constant (equation (4.102))}$$

$$M\ddot{z}+T\dot{z} \tag{4.73}$$

$$M\ddot{z}+T\dot{z}+kz \tag{4.74}$$

$$M\ddot{z}+T\dot{z}+kz+F. \tag{4.75}$$

Equation (4.73) represents the purely dynamic components involving time, whereas equation (4.74) represents the spatial components, that is those involving displacement $z$, and equation (4.75) the total force.

Consider next the way in which the force equation can be presented. In the form shown in equation (4.71) it is in its temporal form and as such is very convenient for representing forces and reactions and, by implication, damage.

### (b) Frequency form of the differential equation

Considerations of speed also involve time, so this temporal form of the equation would also be suitable. However, an alternative form, which allows both speed and fidelity to be examined is obtained if the force equation is transformed into the frequency domain.

Thus equation (4.71) becomes $F_R(\omega)$, the Fourier transform of $R(t)$, where $F_z(\omega)$ is that of $z(t)$. So

$$F_{Rz}(\omega) = F_z(\omega) M\omega_n^2 \left\{ \left[ 1 - \left( \frac{\omega}{\omega_n} \right)^2 \right]^2 + 4\zeta^2 \left( \frac{\omega}{\omega_n} \right)^2 \right\}^{1/2} \exp(-j\varphi t) \tag{4.76}$$

where

$$\varphi = \tan^{-1} \left( \frac{-2\zeta\omega/\omega_n}{1 - (\omega/\omega_n)^2} \right) \tag{4.77}$$

and $\omega_n^2 = k/M$, $2\zeta\omega_n = T/M$.

### (c) Statistical form of the differential equation

Equations (4.71) and (4.76) represent the temporal and frequency equivalents of the force equation. More realistically, to cater for the considerable noise present in atomic measurement the stochastic equivalents should be used, namely

$$\varphi = \tan^{-1} \left( \frac{-2\zeta\omega/\omega_n}{1 - (\omega/\omega_n)^2} \right) \tag{4.78}$$

and

$$A_R(\beta) = M^2\omega_n^4 \left[ A_z(\beta) - \left( \frac{1}{\omega_n} \right)^2 \frac{\mathrm{d}^2}{\mathrm{d}\beta^2} A_z(\beta)(4\zeta^2 - 2) + \left( \frac{1}{\omega_n} \right)^4 \frac{\mathrm{d}^4}{\mathrm{d}\beta^4} A_z(\beta) \right] \tag{4.79}$$

where $P_R(\omega)$ is the power spectral density of $R(t)$ and $A_R(\beta)$ the autocorrelation function.

Immediately it can be seen, from equations (4.78) and (4.79), that the properties of the output $P_R(\omega)$ are related intimately with the system and the input signal parameters $M$, $\omega_n$, $\zeta$ and $P_z(\omega)$ respectively.

In its simplest form for a random input of, say, displacement or topography to the system, consider the case when $\beta = 0$. In this case all values of $\omega$ are taken into account, not one at a time as in equation (4.77).

Thus from random process theory

$$A_R(0) = A_z(0)M^2\omega_n^4\left(1 - \frac{1}{\omega_n^2}\frac{\int\omega^2 P(\omega)}{\int P(\omega)}(4\zeta^2 - 2) + \frac{1}{\omega_n^4}\frac{\int\omega^4 P(\omega)}{\int\omega^2 P(\omega)}\right). \tag{4.80}$$

or

$$A_R(0) = A_z(0)M^2\omega_n^4\left(1 + \frac{\overline{\omega}^2}{\omega_n^2}(4\zeta^2 - 2) + \frac{\overline{\omega}^4}{\omega_n^4}\times\frac{\overline{v}^2}{\overline{\omega}^2}\right) \tag{4.81}$$

where $\overline{\omega}$ is the RMS angular frequency of $2\pi/\zeta_q$, where $\lambda_q$ is the average distance between positive zero crossings, and $\overline{v}$ is the RMS peak frequency of $2\pi/\overline{\lambda}_q$ where $\overline{\lambda}_q$ is the average distance between peaks.

Equation (4.81) shows how the variance of $R$ relates to the variance of the cantilever deflection taking into account the system parameters $\zeta$ and $\omega_n$ and the statistics of the input given by $\overline{v}/\overline{\omega}$. This formula enables all types of input shape to be taken into account. For $\overline{v}/\overline{\omega} = 0$ the signal is white noise, and for $\overline{v}/\overline{\omega} = 1$ it is a sine wave.

When topography is being considered, equation (4.80) can be written as

$$\sigma_R^2 = M^2\omega_n^4[\sigma_z^2 + \sigma_s^2/\omega_n^2(4\zeta^2 - 2) + \sigma_c^2/\omega_n^4] \tag{4.82}$$

Here $\sigma_z$, $\sigma_s$, and $\sigma_c$ are the standard deviations of the surface heights, slopes and curvatures respectively.

For investigations of fidelity involving wide bands of frequency and many varieties of profile form, equations (4.77)–(4.82) are to be used.


### (d) Spatial form of the differential equation

It could be argued that if fidelity is to be considered, properly then it should be arranged that all issues are referred to coordinates of the waveform and not of the system as in the above. Hence, for fidelity reasons, it may be beneficial to express all factors in terms of height $z$ and distance $x$ along the surface.

Hence if the velocity of scan is $V(x)$ and not constant the basic equation can be rewritten as $R(x)$, where

$$MV^2(x)\frac{d^2z(x)}{x^2} + V(x)\frac{dz(x)}{dx}\left(M\frac{dV(x)}{dx} + T\right) + kz(x) = R(x) - F \tag{4.83}$$

which, with certain constraints on $V$ to be identified later, can be simplified to

$$\frac{d^2z}{x^2} + 2\zeta(x)\omega_n(x)\frac{dz}{dx} + \omega_n^2(x)z = [R(x) - F]/MV^2 \tag{4.84}$$

where

$$\zeta(x) = \left(\zeta + \frac{1}{2\omega_n}\frac{dV}{dx}\right) \quad \text{and} \quad \omega_n(x) = \frac{\omega_n}{V}. \tag{4.85}$$

Expressions (4.83)-(4.85) will be used in the case for point-to-point fidelity. In equation (4.83) $R$ is provided by the static force $F$. Of the rest the force exerted by the restoring force due to the spring constant $k$ is the next most important. Note that in open-loop topography measurement the cantilever and spring are not essential elements; they are included to keep the stylus on the surface when tracking, unlike in force measurement where they are essential.

### 4.2.2.8 Measurement systems (surface topography/atomic force measurement—loop considerations)

The surface parameter, whether it is force or topography, can be measured in the open-loop or closed-loop mode shown in figures 4.62(*a*) and (*b*). Both figures incorporate the basic mechanical elements described by the differential equation given in equation (4.71).

(*a*)

(*b*)



**Figure 4.62** Simplified schematic diagram of (*a*) open-loop and (*b*) closed-loop systems.

#### (*a*) Topography measurement

Consider first the measurement of topography. Because it is a tactile instrument, in the first instance open-loop operation will be considered to follow conventional surface roughness practice.

Consider equation (4.78) relating the power spectral density of the reaction to that of the surface topography. Ideally from a damage point of view $R$ should be zero or constant. In practice for open loop this is impossible; the nearest option is to make the value of $P_R(\omega)$ equal to $P_z(\omega)$. This can be done over the whole frequency range $0 < \omega < \omega_n$ by making the term in the square brackets of the equation below integrate to unity:

$$\int_0^1 \left\{ 1 + \left[ \frac{\omega}{\omega_n} \right]^2 (4\zeta - 2) + \left[ \frac{\omega}{\omega_n} \right]^4 \right\} d \left[ \frac{\omega}{\omega_n} \right] = 1.$$

(4.86)

The only variable is $\zeta$. A value of $\zeta = 0.59$ satisfies this equation (shown in figure 4.60) as already seen in equation (4.60). Under these conditions, the integrand is called the 'damage fidelity weighting function' and is designated by $H(\omega)$. Using the equation for the average wavelength of any random signal given in equation (4.80) and utilizing the fact that the cross-correlation of $R$ and $z$ is $CR_z(\beta)$ given by

$$C_{Rz}(\beta) = A_z(\beta) M (D^2 - 2\omega_n \zeta D + \omega_n^2),$$

(4.87)

it is easy to show that the phase shift between a random topography and the reaction from the tip of the scanning stylus is about 40°. This indicates that spatial fidelity between tip reaction and surface topography is nominally unity and, because of the phase shift, shear damage would be likely to be small.

So far only fidelity due to damage has been considered. Also important is the fidelity between the transfer of the topographic signal $z_i$ from the tip to the transducer output $z_o$ This can be obtained from $H(\omega)$. Thus

the transfer function is given by TF where

$$\mathrm{TF} = \frac{1}{H(\omega/\omega_n)^{1/2}} = \frac{1}{[1 + (\omega/\omega_n)^2(4\zeta^2 - 2) + (\omega/\omega_n)^4]^{1/2}}$$

(4.88)

Conventional criteria for operating this transfer function require that $\omega \ll \omega_n$ so that TF $\sim 1$. This severely restricts the operating speed of the system. It is here proposed, as in damage fidelity, that an integrated approach is considered. Hence, allowing all frequencies $0 < \omega \le \omega_n$ to be present and making

$$\int_0^1 = \frac{1}{(H(\omega/\omega_n))^{1/2}} \mathrm{d}\left(\frac{\omega}{\omega_n}\right) = 1$$

(4.89)

gives the result that $\zeta = 0.59$. Indeed, if instead of displacement transfer $z$ the energy transfer is considered then

$$\int_0^1 = \frac{1}{(H(\omega/\omega_n))} \mathrm{d}\left(\frac{\omega}{\omega_n}\right) = 1$$

(4.90)

is the criterion and the best value of damping is again 0.59.

*The remarkable fact is, therefore, that in damage fidelity, displacement fidelity and energy fidelity an integrated damping coefficient of $\zeta = 0.59$ is optimum.* It seems strange that this value of damping has never before been singled out as having unique properties!

Incidentally, for all three cases cited above the values of $\zeta$ are not quite the same but they are all within 1% of each other.

Using this integrated damping coefficient for, say, topographic fidelity ensures that $z_o$ is equal to $z_i$ to within 15% at $\omega_n$ and 5% at $\omega_n/2$. There is an average phase shift of about 40°, so that distortion is not serious.

Even if $\zeta \ne 0.59$ but is constant and known, the weighting factor $H\zeta(\omega)$ can be calculated and the output signal compensated. For $\zeta = 0.59$ compensation to get all $H(\omega)$ values equal to unity and the phase modified according to equation (4.77) are probably not necessary. Thus

$$\text{for } \zeta = 0.59 \qquad P_z(\omega) \simeq P_R(\omega)$$
$$\text{For } \zeta \ne 0.59 \qquad P_R(\omega) = P_z(\omega)H_\zeta(\omega)$$

(4.91)

(*i*) *Speed*

Making $z = 0.59$ or thereabouts according to expression (4.91) allows fairly faithful reproduction of all frequencies to be obtained. This means that it is therefore possible, while maintaining credible topographic fidelity, to increase the speed of traverse such that $\omega_{max} \sim \omega_n$ and not $\omega_{max} \ll \omega_n$ as previously required. This implies that an increase of ten to one in speed over conventional systems is possible for random surfaces and probably also for spiky surfaces.

(*ii*) *General*

Design changes of the system involving truncating the spatial force terms from right to left leaving $M\ddot{z}$ are a possible way of reducing damage caused through contact.

Speed and topographic fidelity of a complete system with no truncation of force components is possible if an integrating damping coefficient of $\zeta = 0.59$ is used. What is also true is that, if this is possible, not only is topographical integrity maintained but the compliance between the stylus and the surface follows the rule that $A_R(\beta) = A_z(\beta)$ and lateral fidelity is preserved.

*(b) Force measurement*

What follows refers to the atomic force microscope type of pick-up rather than the tactile one used for conventional surface topography.

Proceeding in the same way as for topography, the same basic equation (4.71) holds (for open and closed loops as will be seen). In this case, however, there are some fundamental differences. One is concerned with the need for the cantilever. It is essential in force measurement because it is the actual sensor. When acting under free or forced oscillating mode it is integral to the measurement. It actually provides the force balance mechanism; it acts as a seismograph where $M$ is the seismic mass. It is therefore impossible to consider removing the restoring spring as in topographic measurement demonstrated earlier. Hence, rather than making $R \to 0$ for damage risk reduction, the criterion for one mode of force measurement is to make $R$ proportional to $kz$. As in purely topographic measurement the term in $F$ is meaningless from the point of view of spatial variation in the force profile.

Before considering fidelity and speed in force measurement it is timely to consider the other mode of operation of the beam when measuring force—that in which the cantilever is forced to oscillate at its natural resonant frequency $\omega_n$. Here there is little interaction with the surface and the amplitude of the envelope of the vibration is changed by the presence of a force gradient rather than the force producing a straightforward deflection of the beam.

Thus if the force gradient is $f'$ this modifies the effective stiffness of the beam to $k$–$f'$ and the spatial components in equation (4.23) become

$$M\ddot{z} + T\dot{z} + (k - f')z. \tag{4.92}$$

This effective change of stiffness changes the amplitude of the beam oscillation according to the formula derived from (4.23). Thus, if $Af$ is the amplitude of oscillation of the beam when in the presence of a force gradient $f'$

$$A_f = \frac{Q\omega/\omega_n}{\{1 + Q^2[(\omega_n/\omega)^2 - (\omega_n/\omega)^2]\}^{1/2}}. \tag{4.93}$$

$Q$ is the quality factor equal to $1/2\zeta$ and is the amplitude of $Af$ when $\omega = \omega_n$. $\omega^2$ is given by the effective resonance

$$(k - f')/M. \tag{4.94}$$

This method will not be considered in detail because of its limited lateral resolution. Also, because of having to integrate the $f'$ value obtained from equations (4.93) and (4.94) (by finding the difference $\omega - \omega_n$ from $Af$) to get $f'$, it is of limited accuracy.

*(c) Fidelity of atomic force measurement*

The criterion for high fidelity is to ensure that $R$ is proportional to $z$, the beam movement. The cross-correlation of $R$ and $z$ is the simplest and truest indication of fidelity. It is given here by $C_{Rz}(\beta)$ where

$$C_{Rz}(\beta) = E[R(t)z(t + \beta)]$$

and $A_z(\beta)$ is the autocorrelation of the beam deflection. Therefore

$$C_{Rz}(\beta) = M\frac{\mathrm{d}^2 A(\beta)}{d(\beta)^2} - T\frac{\mathrm{d}A_z(\beta)}{\mathrm{d}\beta} + kA_z(\beta) \tag{4.95}$$

For ideal force measurement

$$C_{Rz}(\beta) \propto A_z(\beta) \qquad (4.96)$$

This equation holds for random and deterministic force signals and (4.95) follows on from the treatment of equation (4.80).

As in the case of topography the simplest way to achieve (4.96) from (4.95) is to truncate the spatial force components, only this time from left to right. It is the inertial and damping terms that need to be reduced rather than the stiffness term. Fortunately, when miniaturizing the system, inertial and damping terms become smaller than the stiffness owing to the 'scale of size' effect. However, this useful tendency is limited because the shank holding the stylus has to be large enough to act as a realistic part of the transducer monitoring the beam deflection. Also there is always a requirement to increase the speed of measurement which offsets the trend to smaller inertial effects.

From equation (4.71) ignoring the possibility of making $M$ zero ($T$ can be made small relatively easily), one option would be to make the inertial and damping terms constant. Achieving this is acceptable because it is invariably the force profile as a function of $x$ (assuming one dimension) that is required. This would be possible if the conditions set out in equation (4.97) are met, namely

$$M \propto 1/\ddot{z}$$
$$T \propto 1/\dot{z}. \qquad (4.97)$$

These would have to be satisfied simultaneously. Although it seems difficult to picture such a mechanism at present, it may well be possible in the future. Another possibility is to make the inertial and damping terms cancel:

$$M\ddot{z} \propto T\dot{z}. \qquad (4.98)$$

Temporal equation (4.98) suggests that $\ddot{z}$ and $\dot{z}$ are negatively linearly correlated, which they are not, but spatially there may be ways of making the force components of $dz/dx$ and $d^2z/dx^2$ correlated.

Possibly the easiest way to reduce the effect of $M$ an3d $T$ is to reduce the velocity of tracking. Thus from equation (4.83), if $V$ is kept constant,

$$MV^2 \frac{d^2z}{dx^2} + TV\frac{dz}{dx} + kz = R. \qquad (4.99)$$

This equation, in which $F$ has been dropped, shows that making $V$ small dramatically reduces the effect of $M$ and to a lesser extent $T$. The problem is then that the measurement is much slower, which is a big disadvantage in STM and AFM operation.

A possible answer emerges if the basic equation is examined. Quite simply, if there is some way of making $\dot{z}$ constant then not only is the damping force constant and hence unimportant as far as variations are concerned but also the inertial term is zero. Hence

$$R = TW + kz \qquad (4.100)$$

where $W$ is a constant. How this could be achieved temporally is discussed in section 4.2.2.10. However, an alternative method involving the spatial coordinate form of equation (4.82) can be used. Equation (4.83) can be rewritten in the following form in which $S$ is the term $dz/dx$ which corresponds to force or topographic

gradient. Thus

$$MV\frac{\mathrm{d}}{\mathrm{d}x}(VS) + TVS + k\int S = R(x).$$
(4.101)

Both $S$ and $V$ in (4.101) are position dependent and not constant.

It is clear from equation (4.101) that making $VS = W$, a constant, automatically removes the inertial term and reduces the damping term to a constant as shown in equation (4.100).

In equation (4.100) because $TW$ is constant it is usually not of interest, so the fidelity is satisfied because changes in $R$ are measured directly by changes in $z$.

Thus, for atomic force microscopes, $V(x)$ could be variable, not constant, and subject to the constraint

$$V(x) = W/S(x)$$
(4.102)

A similar benefit is obtained for conventional surface-measuring instruments if this relationship between scanning velocity and surface slope is maintained. In this case, however, the constant term cannot be ignored. Although it produces constant compliance and therefore does not affect fidelity, it could cause damage.

Obviously equation (4.102) cannot be allowed to stand because it allows negative (and infinite) velocities. How these situations can be ameliorated will be discussed later.

### (d) Speed of measurement for force

The situation here is the same as for topography. In principle, allowing an integrated damping coefficient of $\zeta = 0.59$ will allow an order of magnitude increase in the speed at the cost of a small loss in fidelity at certain frequencies. Even if these losses are not compensated, as in equation (4.91), they probably will be more than offset by the loss of fidelity introduced by having a slower measurement cycle.

### 4.2.2.9 Spatial domain instruments

#### (a) Spatial coordinate system

It has already been indicated that the fidelity of a profile, whether it is of force, tunnelling current or topography is usefully referred to spatial coordinates rather than temporal. The frequency transformation is useful in integrated damping in the sense that it has revealed the possible benefit of an integrated damping coefficient in both fidelity and speed considerations. However, it may not deal completely with point-to-point fidelity of the profile, especially from the point of view of phase. A complementary approach, again involving damping, utilizes the spatial representation of equation (4.71) given in equation (4.84). It embodies another transform, this time being

$$V(x) = \mathrm{d}x/\mathrm{d}t$$
(4.103)

which is a hyperbolic transformation between time and space and as a result is sometimes difficult to handle. Thus the natural frequency of the system, $\omega_n$, becomes position dependent via $V(x)$ to become $\omega_n(x)$, where $\omega_n(x) = \omega_n/V(x)$ and the damping coefficient $\zeta$ becomes modified to take the spatial form

$$\zeta(x) = \left(\zeta + \frac{1}{2\omega_n}\frac{\mathrm{d}V(x)}{\mathrm{d}x}\right).$$
(4.104)

The name for $\zeta(x)$ has been taken to be the differential spatial damping coefficient (DSDC).

Reference to equation (4.104) shows that the effective damping coefficient can be changed point by point in the $x$ direction (and in two dimensions) simply by changing the velocity of scan: the velocity acts

as a coupling between the temporal instrument system and the spatial surface system. Letting

$$V = \bar{V} + u(x) \tag{4.105}$$

where $u$ is a variable and $\bar{V}$ is constant, shows that

$$\zeta(x) = \left( \zeta + \frac{1}{2\omega_n} \frac{du}{dx} \right) \tag{4.106}$$

which illustrates that the damping is not dependent on the value of the tracking velocity but just its changes. Furthermore, the effectiveness of the change in velocity on the damping is governed by the natural frequency of the system. Making $\omega_n$ low amplifies the effect of $du/dx$. Consequently, having a low $\omega_n$ has two benefits: one is that there is a much greater capability for point-to-point damping changes, and the other is that if $k$ is made small the sensitivity of the force-measuring system is increased. On the other hand, if $M$ is allowed to be larger to reduce $\omega_n$, there is more room to monitor beam displacement. In all cases there is a difference from the conventional thinking of making $\omega_n$ high. Obviously because $\omega_n$ is allowed to be lower using the integrated damping philosophy for speed enhancement, a realistic speed increase is more possible on both counts.

It is definitely advantageous to have this degree of flexibility with the damping term, especially when miniaturizing. Damping is very much the unknown quantity at the atomic level!

A point to note is the similarity between the differential spatial damping coefficient and the position-related coefficient of friction on the microscale. That is, $T\dot{z}$

$$\frac{dz}{dx} TV \left( 1 + \frac{1}{2\omega_n \zeta} \frac{du}{dx} \right) \tag{4.107}$$

and $\mu R(x)$ becomes

$$\mu R(x) = \mu R \left( 1 + \frac{1}{\mu_i} \frac{du}{dx} \right). \tag{4.108}$$

In principle, the damping coefficient can be made equal to 0.59 for any $x$ by changing $du/dx$. If this were possible then the instantaneous resolution $\lambda$ on the surface laterally would become $\lambda_n(x) = 2\pi/\omega_n V$ rather than the conventional $\lambda_{max} = 2\pi/\omega_{max} V$ where $\omega_{max} << \omega_n$.

### (b) Coupled damping

The expression for equation (4.83) is complicated in the sense that it is made up of two differential equations, both involving $x$. One is concerned with the differential relationships between $z$ and $x$ and the other between $V$ and $x$. These become highly coupled if they are related, as in the section on fidelity when $VS$ is constant. In this case control is not simple.

It has already been pointed out that this system is impractical because of negative and zero values of $S$. It can be modified, however, to allow

$$V(x) = \frac{W}{|\,dz/dx\,|} \tag{4.109}$$

Equation (4.83) then becomes

$$\frac{MW^2}{|\,dz/dx\,|} \left[ \frac{d}{dx} \left( \frac{dz/dx}{|\,dz/dx\,|} \right) \right] + TW \frac{dz/dx}{|\,dz/dx\,|} + kz. \tag{4.110}$$

The first term on the left is zero for practical purposes leaving

$$R = \text{sgn}\left(\frac{dz}{dx}\right)TW + kz. \tag{4.111}$$

Fortunately the sign in equation (4.111) is always known because the spatially dependent $dz/dx$ sign can always be found.

The only problems here concern the case when $S = 0$, in which case a velocity limiter is needed, which would be based on a value estimated from $\bar{\omega}$ and $\bar{v}$ in equation (4.81) and the fact that the speed is somewhat reduced. Control of a coupled system such as this is highly non-linear.

An alternative way of making the measuring system more linked to the actual spatial requirement is to make the profile length scan velocity constant. This does not mean the $x$ or $z$ velocity being constant but that the velocity of the locus of the scanning tip as it measures force or topography is constant. Then

$$V = W\left[1 + \left(\frac{dz}{dx}\right)^2\right]^{-1/2} \tag{4.112}$$

from which

$$W^2 \frac{d^2z}{dx^2}\left[1 + \left(\frac{dz}{dx}\right)^2\right]^{-2} + W\frac{dz}{dx}\left[1 + \left(\frac{dz}{dx}\right)^2\right]^{-1/2} 2\omega_n\zeta + kz = R. \tag{4.113}$$

From this two things emerge: the first is that there is a desired reduction in the inertial and damping forces as the surface slopes get high, which helps damage control. The second point is that the inertial term, although not zero as in equation (4.111), is considerably more reduced for a given change in velocity than it was for constant velocity $V$ in the $x$ direction given in equation (4.99).

Hence for constant horizontal scan the balance of dynamic forces is

$$V^2 \frac{d^2z}{dx^2} + V\frac{dz}{dx}2\omega_n. \tag{4.114}$$

For constant tip locus velocity the balance of dynamic forces is given by

$$V^4 \frac{d^2z}{dx^2} + V^2\frac{dz}{dx}2\zeta\omega_n. \tag{4.115}$$

It could be argued that the need to maintain constant profile path scan will become more important as atomic levels are approached. This is because, at these levels, the heights and spacings of features are at the same scale. For conventional surfaces there is a scale change of a hundred times or more between vertical and horizontal features. This difference is usually caused by the mechanism of surface manufacture. It seems likely that this reduction in change of scale of the axes has already occurred in many practical measuring examples.

The specific transformations in equations (4.109) and (4.112) result in instruments where the measuring system is more closely linked with the properties of the surface parameters than ever before. The *surface itself becomes a system parameter*. It seems that this is a way to go in the future in order to get the best results; this is the key argument in characterizing the function map in chapter 7.

### 4.2.2.10   Open-and closed-loop considerations

(*a*) *Electronic feedback*

The basic transfer function of the system can always be modified by closing the mechanical loop by rate feedback or position and rate feedback to give modified mechanical parameters such as equation (4.112):

$$TF = \frac{G_1}{pMR_1 + p(R_2T - G_2) + R_3k - G_3}$$

(4.116)

where $R_1$, $R_2$, $R_3$, $G_1$, $G_2$ and $G_3$ are electronic loop parameters. The problem with this approach is that the compensation in, say, the damping term $p(R_2T - G_2)$ is a fixed change dependent on the values $R_2$ and $G_2$ and is not flexible on a point-to-point basis such as the technique proposed in differential spatial damping. This instant flexibility is very necessary for miniature systems where mechanisms are not well understood or controlled.

In this context it is the mechanical system that has been considered and it has been shown that, by using the very simple trick of utilizing the tracking speed, great flexibility is possible to enable fidelity and speed to be improved. This does not mean to say that purely electronic means are not possible. It is possible to make $G_2$, for example, dependent on $z$ and not fixed as is conventional. This and other practical implementations of this theory are still to be resolved.

(*b*) *Mechanical 'follower' systems*

To obtain topography a servo system based on keeping force or current or focus constant is often adopted. This has the advantage that such methods do not rely on large linear ranges relating the parameter (e.g. current) to separation, which tends to be the case in an open loop. Apart from this the arguments used in this section are equally valid. The loop is invariably closed by the basic mechanical sensor transducer system having the form of equation (4.71).

### 4.2.3   Scanning probe microscopes [48—51] ( SPM )

It has been estimated [52] that there are at least twenty types of probe instruments. The general consensus is that the biggest industrial application is primarily of the AFM but also the STM. The main uses are (*a*) surface topography on atomic scale and (*b*) dimensional metrology where this is taken to mean lateral position and spacing. It is also recognized that surface feature fabrication will be a rapidly developing application with 'designer surfaces' at a nanometre scale. The industrial applications have been reduced by Vorburger *et al* to the following:

(1) Data storage
(2) Microelectronics
(3) Polymers and coatings
(4) Optics
(5) Mechanics
(6) Electrochemistry
(7) Biology and biomedicine.

Of these (2), (4) and (5) are most often cited.

In (1) the roughness of the magnetic recording head is critical, as is the disc. The critical dimension is the distance vertically between the write pole and the slider surface. This can be between 5 and 15 nm. These roughnesses and dimensions are getting out of optical range. AFM presents the natural successor.

In (2) in the semiconductor industry roughness is again very important, but line widths and spacings of circuit elements also need to be controlled as well as step heights and edge roughness. It is now recognized

that hole dimensions and the smallest line widths are the so-called 'critical dimensions' or 'CDs'—an unfortunate acronym. It is not often a roughness value such as $R_a$ with a value of 1-3nm which is measured, but usually a scan for evidence of process faults, e.g. anisotropic etching [53]. Roughness is a questionable variable.

Line width and spacings are difficult to measure because the shape of the probe is critical—and largely unknown. Some methods of tip examination will be given in chapter 5.

Needless to say, modelling of the tip surface interaction is being used to ease the measurement problem. It appears that the SEM (scanning electron microscope) is more difficult to model than the AFM interaction but it is much quicker to make the measurement, so AFM is more accurate but slower than SEM. Which is used depends on the application.

In (3) of the list above, roughness has an indirect effect on many of the chemical processes being monitored by AFM. These include polymer adhesion wettability as well as friction and wear.

One aspect of measurement which has been influenced by coating application has been the introduction of the discrete 'tapping' mode for AFM and STM instruments to minimize the damage caused by the dynamic effects.

In (4) values for optics have been well documented. One example for x-ray mirrors from Windt *et al* [54] is given here. He has shown that $R_q$ roughness of 0.1nm will be a requirement for EUV lithographic optics operating at 13.6nm for the new generation of chips with CDs of 100nm.

In (5) mechanical applications using AFM for example in the automotive industry and general manufacture have been limited except for some examples on tool tip measurement. Conventional instrumentation for measuring texture and form have developed considerably in the past few years. Commercial instruments are being made with the capability of measuring form and texture—including waviness with one measurement. Integrated measurement has many advantages including ease of spatial calibration and reductions in set-up times. Scanning instruments in the form of STM, AFM, SEM etc have multiplied dramatically, together with the number of features being measured. Engineering applications are in a minority.

In a recent CIRP survey on SPM instruments some idea of the various uses of the whole family was given [52] (table 4.3). Some of the details follow in table 4.4 showing where SPM instruments fit into the industrial scene. There is confusion over the role of these instruments. Should the results be traceable or not? Should every instrument be 'metrological'? Is there a metrology loop within the instrument that can be traced back to international standards via a laser interferometer, for example? Or, is the picture revealing the surface enough to monitor the process or should the $z$ deflection be traceable as a movement anyway? These issues have not been resolved although it seems unlikely that the extra cost and complication of a 'metrology SPM' will allow its general use. The same arguments were true for the SEM and TEM. It is noticeable that the vast majority of applications of the SEM in particular are simple pictures. Traceability has usually been achieved by means of a Talystep or similar instrument. Traceability of $x, y, z$ in position does not make the variable traceable! See chapter 5.

**Table 4.3**

| Surface | Function | Roughness $R_{q \text{ (nm)}}$ |
|---|---|---|
| Poly Si Etch damage | Gate conductor | >1 |
| Poly Si Rapid thermal chemical vapour deposition | Gate conductor | 3–15 |
| $Si\,O_2$ | Gate dielectric | 0.05–0.2 |
| TiN | Adhesion and barrier | 0.5–1.0 |
| Al-Cu (1%) | Conductor | 8–20 |
| Bare Si | Starting material <0.1 | |

**Table 4.4** Typical SPM instruments with applications.

| Class of techniques | Name | Type of probe surface interaction (if not clear from the name) | Typical applications |
|---|---|---|---|
| Force-based (i.e. causes a change in cantilever deflection or oscillation) | Atomic force (AFM) | Normal mechanical contact or long-range non-contact interaction | Semiconductors, optics, data storage |
| | | Lateral contact | |
| | Frictional force (FFM) | | Polymers, contrast enhancement on surfaces with low AFM contrast |
| | Magnetic force (MFM) | | Magnetic storage, materials science |
| | Electrostatic force (EFM) | | Locally induced charges on insulators, triboelectrification, domain structures of ferroelectrics |
| | Kelvin probe force (KPFM) | Senses work function variations | Dopant profiles, localized charge defects, band bending |
| | Eddy current | Force between induced eddy current and ferromagnetic element | Soft magnetic materials |
| | | | Surface roughness profiling |
| | Near-field acoustic (SNAM) | Detects probe-induced acoustic waves in surfaces | Subsurface defects |
| | Tunnelling acoustic (TAM) | | |
| Optical | Near-field optical (SNOM or NSOM) et al. | | Biology spectroscopy in materials science, semiconductors |
| Thermal | Thermal Profiler (SThP) *et al*. | | Thermography in microelectronics |
| Electrical | Tunnelling (STM) | Electronic tunnelling | Semiconductors, optics |
| | Capacitance (SCM) | | Lateral dopant profiling for semiconductors |
| | Electrochemical (SECM) *et al*. | Electrolytic current | Non-contact profiling, chemistry of surfaces |
| | Nanoscilloscope | RF field coupling | *In situ* testing of semiconductors |

What is questionable in all these arguments is the willingness to apply the same traceability criteria to the $z$ axis as to $x$ and $y$!

### 4.2.3.1 Scanning tunnelling microscope (STM)

Another version of the stylus topographic instrument is the scanning tunnelling microscope (STM) or the atomic force microscope (AFM). Although new generations of instruments do not appear at first sight to have much in common with the traditional stylus instrument, they do. The only difference basically is that instead of measuring geometry they measure another feature of the surface such as charge density, force, etc. The basic mechanical system is still second order and, more relevantly, the input is still determined by a stylus. Other issues such as the use of skids and references are more relevant to more conventional stylus instruments and basically peripheral to their operation. Consequently they will be considered in section 4.2.4.

Historically, scanning probe microscopy had its origins in a principle first outlined by the British scientist Edward Synge in 1928. He suggested the use of a tiny aperture at the end of a glass tip to raster an illuminated object. This concept was resurrected by J A O'Keefe of the US Army Mapping Service in 1956. To overcome the diffraction limit for microscope studies (i.e. approximately $-\lambda/2$) O'Keefe suggested using a tiny aperture in an opaque screen. By placing this aperture close to the surface being studied and scanning the aperture across the surface in a raster pattern the diffraction limit is effectively overcome. By using this near-field aperture the resolution could be determined by the size of the hole and not the wavelength of light.

Scanning probe microscopes (SPM) make use of a trick in the physics. Conventional microscopes have their performance determined by the wave properties of light. Each wavelength $\lambda$ has an amplitude, phase and polarization vector associated with it. The resolution and depth of field are determined by the way in which constituent rays making up a beam of light combine, taking into account the wave properties indicated above. These properties actually constrain the behaviour of light, for example when passing though an aperture. Diffraction effects make sure that light emerging through a very small aperture gets broadened and not diminished. It is only when light is released from the physical law straitjacket that the spot can be made much smaller than the airy disc. By scrambling wavelength and phase and dealing only with intensity, subnanometre spot size can be obtained. All that is required is a means of projecting such a spot onto the surface and picking up any remnant of the illumination from it. This requires a local transport agent such as a fibre optic with very small aperture or a pipette. It is only recently that such miniature devices have been available. This and the ability to position the probe with subnanometre accuracy constitute the basic elements of these scanning microscopes.

Investigation of other local properties of the surface can be made via a very small tip of a mechanical probe. The probe has to react to the surface property directly as in the attractive forces of the surface or indirectly by matching the probe sensor properties of the surfaces. The ability to utilize the probe (stylus) for a wide range of properties has led to a large family of scanning microscopes. Generically these are called 'scanning probe microscopes' (SPM).

Some features of a few of these microscopes will be given below, but first the tip of the stylus will be considered.

The basic system is deceptively simple, as shown in figure 4.63. the problem is that more emphasis is placed on the stylus than with the conventional stylus instruments mentioned earlier. The stylus takes on a more proactive role than it does in surface metrology.



**Figure 4.63** Cantilever and detector.

Figure 4.63 shows the basic features of the stylus assembly. The components are A, the interaction between the stylus tip and the surface; B, the size of the tip and its shape; C, the cantilever upon which the tip is fixed. It is the bending of the cantilever, due to the interaction of the tip with the surface, that provides the signal. D is the transducer shown here in one of its usual modes, which is a laser system. This responds to the angle of the cantilever at the point of incidence. Sometimes the cantilever is excited by a high frequency vibration to modulate the signal. E is a schematic representation of the modulator.

These five elements of the pickup system have been the subject of many papers because of the crucial role of the stylus in SPMs. To some extent properties of the stylus determine the nature of the signal obtained form the surface; applying a voltage between the stylus and the sample, for example, allows charge density in the region of the stylus tip within the sample to be examined. What is often missed is the correlation and calibration between the surface measured and the stylus movement.

*The tip*

The conventional stylus has a tip radius of about $2\mu$m to $10\mu$m. It can be conical or pyramidal depending on the country. Very fine tips can be made of dimension 100 nm and $20°$ angle but these are difficult to make and prone to damage and as a result are costly. [55].

Many different types of tip have been used in scanning probe microscopes. Chapter 8 mentions a few. Problems include the actual tip dimension as well as its stiffness. Figure 8.23 shows silicon tips used in AFM. Other elements have been used for example platinum/iridium and increasing use is being made of molecular derived tips such as zinc whisker and carbon nanotubes. These reduce the uncertainty of the tip dimension because of the discrete nature of atomic and molecular dimensions. Measuring the tip has always been a problem because of the properties of the tip material. The use of diamond for the tip is not straightforward because of its tendency to charge up. This has been remedied by coating with gold or a similar conductor so that SEM pictures can be obtained. Diamond pyramids as well as tetrahedral shapes (for Berkovich tips) can be treated in this way.

The requirement for stiffness of the stylus in AFM and STM is usually in the same direction as that for engineering surface metrology (i.e. at right angles to the direction of the scan).



**Figure 4.64** Stiffness requirement of probe.

The reason for the configuration (*b*) the side acting gauge in engineering metrology rather than the obvious axial gauge shown in (*a*) is to enable the probe to get into holes. This configuration is more demanding than (*a*). For completely different reasons the STM AFM configuration (*c*) has the same potential for hole measurement. It just happens to be convenient to use the deflection of the cantilever as the output. In either case, if the drive is such that discrete measurements are made in an axial 'tapping' mode the stiffness requirement is less lateral but imposes more dynamic constraints in vibration.

Such is the uncertainty of the tip that attempts have been made to estimate the tip dimension from the measurements themselves using neural methods [56]. This method is, in effect, a deconvolution of the profile. Assumptions have to be made concerning the form of the tip shape.



**Figure 4.65** (*a*) Pattern of the feature vector extraction (*b*) deconvolution system based on a trained neural network.

Figure 4.65 Shows how a measured surface can be reconstructed via a neural network.

The idea is that the actual shape of the stylus gets built up within the network during the process of measuring known surfaces; it is contained within the network (i.e. it is implicit). In the example shown, a multi-layered radius basis function (RBF) has been used. It should be emphasized that this method allows mild non-linearities. Increasing the number of hidden layers in the network enables better approximations to be obtained. As in all neural systems it is only as good as its training. Surfaces which lie well outside the 'experience' of the network suffer distortion.

Traditional methods of making stylus tips have usually been based on polishing. This technique has been used for diamond as well as sapphire and ruby tips but it can be time consuming. Very often there are two requirements, one is that the tip dimension is small and the other is that the aspect ratio (i.e. the flank slope) is small. To give an example, a stylus made of diamond can be polished to a $0.1\mu$m tip and an included angle of $20°\mu$m, but mechanical polishing is frustrating and often results in a damaged stylus with either the tip or angle incorrect.

The desired parameters have to be attained by a process of trial and error. Slightly more controlled techniques such as electro-chemical etching have had to be used [57] for scanning tunnelling microscopes. With these methods tungsten tips of 5nm diameter have been achieved. However, the presence of oxides has sometimes caused platinum to be used instead of tungsten.

It is not usual to generate artefacts (i.e. false values of topography) using a mechanical stylus in engineering applications. The usual complaint would be that the tip tends to integrate (i.e. it sometimes smooths) the very short wavelengths. The same is not true for SPM when using a conventional pyramidal probe made from silicon and working in air. Figure 4.66 shows a protein filament as measured by the silicon probe. The upper picture shows artefacts due to the side of the tip making contact rather than the apex [58]. The lower panels show an image with an attached nanotube as shown in Figure 4.67(*b*).



**Figure 4.66** Artefacts introduced by probe.

The degradation of (*a*) in figure 4.66 is due to the finite slope of the probe. When the carbon nanotube is added the aspect ratio increases dramatically and a very high fidelity signal results.

Figure 4.66 (*a, b*) demonstrates the fact that the slope of the stylus (i.e. the probe) can act as an integrator of the signal as well as the tip dimension. (*c, d*) shows that this slope integration is removed when using a tube of virtually constant diameter as the probe. What is not generally known is that, even for machined surfaces, as the surface finish gets smoother the slopes on the surface get higher, with the result that slope integration

becomes more important. Another point is that many objects being measured, such as optical discs, gratings etc. have edges that cannot be adequately measured with large angle probes. For example, a Berkovich stylus having three effective flanks and an angle of about $140°$ is prone to integration irrespective of the tip dimension.

Carbon nanotubes do not have an obvious shape for a probe. A tube attached to the end of the probe as in Figure 4.67(*b*) is, in effect, a cantilever on a cantilever. However, having an intrinsic diameter of 0.7nm and a length which can be microns makes it an attractive candidate for measuring not only fine surfaces but also quite coarse ones. Also, the tube retains its integrity despite bending. Whether this is good or bad is a moot point: distortion of any sort is suspect. The wear properties are nevertheless much better than the silicon probes. Multi-walled carbon nanotubes deflect and wear less than their single walled equivalents but have a bigger size. Single wall carbon nanotubes can consistently resolve to 2nm [59].



**Figure 4.67** Increased resolution with nanotube.

It now seems that carbon nanotubes, coupled with AFM measurement has great potential for a surface characterization tool in integrated circuit manufacture. Some applications would be in measuring gate dielectric layers of nanometres and for characterizing surface morphology and uniformity of non-conducting ultra-thin films.

How to add carbon nanotubes to a conventional AFM probe is a new art and falls outside the scope of this book. There are, however, some useful references: [60–63].

In some cases above the nanotube is attached to a cobalt coated cantilever. They are brought together and in some way the cobalt acts as a catalyst in the formation of the join.

If the local field at the tip of the tube is about $10^6$ volts/mm an arc occurs between the tube end and the cantilever, which causes disassociation of carbon at the tube end and facilitates bonding. The whole process is extremely complicated. The success rate data for attempts to attach a tube correctly are usually conspicuous by their absence.

Carbon nanotubes have a cylindrical cross section. Other tip shapes, including conical and exponential have been used and their properties investigated especially, for vibration but are not reported here except to reveal that such vibrations at the tip can degrade lateral resolution by the order of a nanometre [64].

There are other probes for use with a silicon cantilever, for example titanium and tungsten have been used as coatings on conductive cantilevers. They provide stable coatings to prevent chemical attack of the probe material.

The nanotube type of probe mentioned above has advantages due to its high aspect ratio. Needless to say, only a molecular structure such as carbon tubes could have the necessary rigidity to be able to track transversely. There are other needs for nano-structured probes. One application is in near field optical microscopy. Obviously these probes have to have an aperture. The problem is to guide light to or from the tip and to control, if required, the tip-sample distance.

One criterion for optical probes for example is a sub-wavelength aperture at the end of a tapered fibre [65]. These probes can be made by ion beam processing [66] or etching [67]. Other criteria are a monomode operation, good polarization characteristics, high brightness with high optical damage threshold and minimum light leakage.

*The cantilever*

The typical cantilever shape is shown in Figure 4.68.



**Figure 4.68** Stylus Beams.

This is basically a single wire. Several parameters have to be considered. One is the spring constant, another is the resonant frequency; the former deciding forces in the event of a contact application, the latter determines the scan speed. The tip dimension and slope of the probe are the other variables. The vee-shaped alternative (*b*) is much less prone to twist; it has a much bigger torsional spring constant compared with a single beam cantilever shown in (*a*). This reduces rotational deflection of the laser beam sensor and minimizes noise. It has been reported that the RMS noise level of a vee cantilever is about 0.25nm, whereas it is about 1nm for a single beam cantilever [68].

Thermal noise in AFM due in part to the cantilever is a problem [69]. The mean square deflection of a single beam cantilever due to thermal noise is $Z_q = \sqrt{\dfrac{kT}{K}} = \dfrac{0.64°A}{\sqrt{K}}$ where $T$ is the absolute temperature and $k$ is the Boltzmann constant, $K$ is the spring constant.

In figure 4.69 the actual deflection at the end of the beam is given by $\delta$ and is in terms of the dimensions and the material properties of the beam. Thus

$$\delta = \frac{4F_N l^3}{Ebh^3}. \tag{4.117}$$

It is not a very useful shape to adopt but it is the easiest. Picking a different type of cantilever could enable the end to remain horizontal, in which case $\delta$ would (figure 4.70) be much easier to measure but the frequency response would suffer and the fabrication would be difficult.

In this case the sensitivity is four times less

$$\delta = \frac{F_N l^3}{Ebh^3} \tag{4.118}$$



**Figure 4.69** Deflection of cantilever.

**Figure 4.70** Zero-slope cantilever.

with the same $F_N$.

Most investigators have chosen the conventional beam and found ways of measuring the deflection. One technique is to use an optical means for measuring the deflection. This is shown in figure 4.71. This method measures the angle of deflection $2\theta_N$ due to force $F_N$ where $\theta_N$ is given by

$$\theta_N = \frac{6F_N l^2}{Ebh};\qquad\qquad(4.119)$$

$b$ and $h$ are the width and height of the cantilever and $l$ is its length. $E$ is Young's modulus. The optics measure $2\theta$ because of reflection at the tip. In one embodiment of this principle [70] a quadrant detector can be used to measure the twist of the cantilever due to frictional forces as well as $\theta_N$. The twist angle or friction angle $\theta_F$, is given approximately by $\theta_F$. Notice the different scale of size sensitivities from $\theta_N$ and $\delta$.

$$\theta_F = \frac{M_t l}{\beta G h^3 b}\qquad\qquad(4.120)$$

where $M_t$ is $aF_F$, the external twisting moment due to friction. $G$ is the shear modulus of the cantilever material and $\beta$ is a constant determined by the $b/h$ ratio. For $b/h > 10$, $b \sim 0.33$; $a$ is the distance from the bottom of the tip to the centre of the cantilever. Typical values for a cantilever system would be $a \sim 2\ \mu m$, $b \sim 10 \mu m$, $l \sim 100\ \mu m$, $h \sim 0.1\ \mu m$, $G \sim 10^{10}\ pA$.

The formulae have to be modified in practice because it is angle of the beam at the tip which is usually measured, rather than the deflection of the tip 'z' itself. Thus d$z$ is changed into $\dfrac{\mathrm{d}z}{\mathrm{d}x}\bigg|_L$ by the factor $\dfrac{2L}{3}$.



**Figure 4.71** Optical readout of cantilever angle.

It is $\dfrac{\mathrm{d}z}{\mathrm{d}l}$ which is measured which can be a problem if the cantilever is in contact with a hard surface. $z$ at the

tip should be zero yet variations in $\dfrac{\mathrm{d}z}{\mathrm{d}x}$ are still possible.

There are two basic situation: (*a*) when the tip is free and (*b*) when it is supported. Case (*a*) is the non-contact situation and (*b*) the contact situation. Depending on the ratio of the spring constant to the sample stiffness, either one or the other of the two modes is pertinent. The various vibrational modes have been evaluated. Up to ten modes can be significant. It does not appear that such calculations have been carried out for the vee cantilever.

The behaviour of the cantilever has been carried out for a variety of situations from end-free oscillations to the situation where the scan element has been replaced by a tapping mechanism [71] in which, for at least part of the oscillation the probe touches the surface impulsively: the latter being used when lateral traversing can damage or move the specimen.

To avoid touching the surface in the tapping mode the cantilever is sometimes restricted to a deflection value which gives a predetermined maximum allowed force. All force-deflection curves are then measured relative to this level. Why Laplace methods are not used in such an analysis of tapping mode is not known; initial value boundary conditions are hardly sufficient to model the tapping behaviour!

Despite the geometry of the probes being simple the fact that the dimensions are very small and awkward makes stiffness and damping coefficients difficult to control.

Some ingenious methods have had to be adopted to estimate the stiffness values in AFM in particular. One method compares the cantilever on test with a special purpose calibrated flexible beam of the same nominal shape and size [72].

Not all AFM systems measure the tip movement by means of optical levers on the deflection. Although the laser deflection method embodies some useful magnification the differential of the displacement $\mathrm{d}z/\mathrm{d}x$ is noisy. It is far more straightforward to measure displacement using for example capacitance methods [73]. It has been suggested that the electrode gap can be used to control the damping of the movement. Setting the damping between 0.5 and 0.7 [74] can be achieved by utilizing the squeeze film forces between the electrodes.

There have been attempts to model the effect of stylus contact in AFM. Some analysis has been from the near field acoustic microscopists. See for example [230]. The basic limitation of acoustic microscopy is the wavelength of the imaging sound. This limit is being overcome by combining AFM with acoustic microscopy. The very limited vibrations of ultrasound in a sample which is being isonified can be measured with the tip of an AFM. There is a non-linear relationship between the vibration on the surface and that imparted to the cantilever. The measurement of the isonification is equated with the tip contacting mode of the AFM, thereby bringing to bear detailed knowledge of the material properties of the sample. Unfortunately much of the calculation is not analytical and some tenuous assumptions have to be made in order to get convincing solutions [77].

The potential applications of the SPM to different properties of the surface are very numerous. Apart from the topography, dopant concentration, resistance, magnetic properties and capacitance have been measured as well as residual stress using Raman spectroscopy.

In practical instruments there is a growing tendency to include a metrology loop to enclose the instrument and yet remain independent of it. Such instruments are often called 'metrological scanning. . . . .'. One such instrument is a capacitative metrology AFM [77].

The need for the 'metrological' instruments is the necessity of producing instruments for commercial applications which have accuracy rather than just resolution. This requirement is relatively easy to satisfy with normal engineering surfaces with the instruments having a strong traceable path and measuring just topography. At the sub-nanometre level, however, it is more difficult because the metrology loop can often be influenced by the force loop driving the scan. Also, in non-topographic applications the link between the variable absolute value and cantilever deflection is not known!

The operational loop is shown in figure 4.72.



**Figure 4.72** Operational correction of metrology loop [77].

The essential point in a metrological design measuring topography is that deflections are caused by the interaction between the tip and the sample. In an instrument using capacitative gauging, the small size and utility of capacitance displacement sensors do not allow contact deflections to be detected. In this instrument the cantilever tip displacement is measured, rather than its inclination, by means of a capacitance gauge, thereby complying with the metrological rule enforcing compatible sensors within any metrology loop (figure 4.73).



**Figure 4.73** Pick-up arrangements.

The 'molecular measuring machine' $M^3$ developed in NIST is an example where every precaution is taken in three dimensions to achieve sub-nanometre resolution throughout the workspace. This instrument is described in detail in section 5.

*Summarizing*
The basic problems with the new generation of scanning probe microscopes are:

1. Generation of artefacts caused by the probe fouling the specimen other than at the tip and the tip triggering unwanted signal components.
2. Calibration of SPMs is difficult because there are no direct natural standards of, for instance, charge density; closing the metrology loop at the stylus by making its movement traceable is indirect because there is no independent calibration of the charge density to the movement.

At the time, the required positioning technology to implement Synge and O'Keefe's technology did not exist but now, with the advent of piezoelectric micropositioning, this has become possible. The STM uses a mechanical probe rather than a hole and does not use an external source to 'illuminate' the surface. Instead it

uses electrons already present on the surface and the tip. In this case the resolution of the STM is based on a near-field 'aperture' that is a single atom at the tip of the probe, this means that, in theory at least, the STM can have a resolution of an atom. The problem is that of vibration and environmental effects. Also the tendency of the tunnelling current emanating from one atom to jump to another on the probe makes artifacts very easy to get.

The earliest in this family of instruments is the scanning tunnelling microscope. It should be obvious that because the method relies on an exchange of electrons from between the specimen and the tip the instrument can only be used for conductors or semiconductors. This originally held back the usage of the STM, especially in biology, but recently a range of instruments have emerged, which will be described briefly. In the STM the probe is actually positioned a distance (albeit small) from the surface and not actually touching, if there is such a thing as touching at this scale! In what follows some other microscopes will be outlined, after which the STM will be described in more detail.

### 4.2.3.2  Other scanning microscopes

(*a*)  *Atomic force microscope* (*AFM*)
The AFM, which will be described later, can be used to study the topography of non-conductive surfaces. Some versions employ a very sharp tip usually of fractured diamond mounted on a small cantilevered beam that acts as a spring. Piezoelectric elements move the tip up towards the sample until interatomic forces between the tip and the sample deflect the cantilever. The AFM monitors the amount of deflection (using optical interferometry or reflected beam or tunnelling methods) to sense the amount of force acting on the tip. With the STM the AFM enables most specimens to be examined.

(*b*)  *Laser force microscope* (*LFM*)
The LFM uses a tungsten probe having a very fine tip. Piezoelectric controls at the base of the probe vibrate the tip at just above its mechanical resonance, the tip is moved to within a few nanometres of the surface. Weak attractive forces from the sample lower the resonant frequency of the wire, reducing its vibration amplitude by effectively changing the spring rate, as seen in equation (4.93) earlier. The amplitude of the vibration is measured. There should be no contact at all with the surface. This instrument is primarily used for imaging microcircuits.

(*c*)  *Magnetic force microscope* (*MFM*)
Much like the LFM, the MFM uses the probe vibrating at its resonant frequency. The difference in the MFM is that the probe is magnetized. As a result this type of microscope senses the magnetic characteristics of the sample. The MFM modulates the resonant frequency, as in the case of the LFM. Its main use has been for magnetic media like recording heads etc.

(*d*)  *Electrostatic force microscope* (*EFM*)
Similar to the LFM and MFM, except that the probe has an electric charge, the EFM has been used most to detect the surfaces of electrically doped silicon for use in microcircuits.

(*e*)  *Scanning thermal microscope* (*SThM*)
The probe here is designed as a thermocouple. This is a tungsten wire with a tungsten-nickel junction. The tip voltage is proportional to the temperature. The tip is heated by an applied current and then positioned near to the sample surface. Heat lost to the sample, which varies according to the space between the probe and the surface, is a measure of the gap and hence the topography. This is used to map temperature variations.

(*f*)  *Scanning ion conductance microscope* (*SICM*)
The SICM uses a micropipette probe containing an electrode. It is intended for use in measuring the activity in living cells.

(g) *Near field scanning optical microscope* (*NSOM*)

This is the nearest device to that suggested by Synge and O'Keefe. This instrument scans a submicrometre aperture across and very close to the sample. The NSOM emits light to the sample and measures the modulation produced by the sample as the aperture is scanned across the surface.

All the microscopes in the STM family have two things in common: they have a probe, except for the NSOM which has a hole, and they have a requirement for very accurate micropositioning. This latter requirement is relatively new in surface metrology. Conventionally, the resolution in the $x$ direction has been much coarser than that for the $z$ direction, reflecting the emphasis on height information. It is only when atomic or molecular-sized objects are being examined that lateral magnifications for the probe-type instrument need to be of the same order. This aspect ratio is shown in figure 4.74.



**Figure 4.74**

The new situation is that not only should the resolution in the $x$ and $y$ direction be the same as that for the $z$ direction but also there has to be accurate positioning. Absolute positioning is needed over the whole range of the $xy$ movement so that structure can be not only seen but measured.

Hence the STM family of instruments have the unprecedented capability of potentially resolving equally in all three axes to atomic levels (in the best case). This poses a completely new problem in the fundamental mechanism of measurement.

Ideally this $xy$ positioning has to be achieved quickly as well as accurately, which poses a severe design problem. Another basic design problem is the necessity for the tip to be reliably positioned within angstroms of the surface—if there is such a thing as a surface—and then to scan over the surface at speed in a raster pattern. Most important is the question of whether the position of the active point of the stylus is ever known (figure 4.75)

As the previous text has suggested, all this family of microscopes has much in common. The major differences are in the configuration of the scanning probe, in the characteristics of the sample being measured, and in the environment in which they are designed to operate.

### 4.2.3.3 Operation and theory of the STM

Before considering some of the relevant design procedures the operation and theory of the STM will be outlined briefly.

**Figure 4.75** Demonstration of the production of artifacts.

Information from the STM is obtained in two ways. One is a null method, the other is open loop.

1. Constant-current mode: As figure 4.76 shows, the tip is moved across the surface. The distance between the tip and sample is maintained by keeping the tunnelling current constant. For any change in current (say, an increase) the STM compensates by withdrawing the tip slightly and *vice versa*. The output from the STM is the voltage applied to the tip piezoelectric control mechanism to maintain a constant current.
2. Constant $z$ (figure 4.77): in this the STM keeps the tip at a constant average $z$. The STM then measures the change in the current. This is the output. Whether this current is calibrated is open to question.

In 1 the output is a smoother curve, but because the servo system is constantly making height corrections throughout the measurement it is quite slow. However, because it is acting as a 'follower' this mode tends to be good for an irregular surface where the constant-height method would be too risky or jumping occurs (figure 4.77 and 4.75). The latter is preferred for smooth surfaces, which can then be measured quickly.



**Figure 4.76** Constant-current mode scanning.

**Figure 4.77** Constant-height mode scanning.

(*a*) *What does the STM measure?*

According to quantum mechanics, electrons on a surface behave both as particles and waves. One result of this is that electrons can be thought of as a cloud above the surface, as in figure 4.78. When two surfaces or, strictly, electrodes are brought together there is a probability that electrons will pass through the potential barrier that is trying to inhibit such movement—the work function. As the surfaces come closer together this probability gets higher until, when the gap is subnanometre, a measurable current will exist. This tunnelling $J_T$ is given by

$$J\text{T} \propto \exp(-A\psi^{1/2}z) \tag{4.121}$$

where $A = ((4\pi/h)\,2m)^{1/2} = 1.025$ Å$^{-1}$ eV$^{-1/2}$ with $m$ the free-electron mass. $\psi$ is the barrier height and $z$ is the spacing. With barrier height $\psi$ (of the order of a few electron volts) a change of $z$, the separation by a single atomic step of 2–3 Å changes the current by up to three orders of magnitude.

Using only the spacing dependence given in equation (4.121) and a spherical tip of radius $R$ the lateral spread of a surface step is about $3(2R/A\psi^{1/2})^{1/2}$, which implies that to resolve 10 nm requires a stylus with a radius of the same dimension.



**Figure 4.78** General situation showing electron clouds: (*a*) surface electronics on the sample; (*b*) tip sample interaction.

In the constant-current mode of equation (4.121) this implies that $\psi^{1/2} z$ = constant. Thus the $z$ displacement, as shown in figure 4.76, gives the surface topography only for constant $\psi$—the work function. So unfortunately if there is a change in work function on a geometrically featureless surface it will register as a topographic feature. One possible way to differentiate between changes in topography and work function is to modulate the gap $z$ while scanning at a frequency higher than the cut-off of the control system, but even then there are complexities.

The tunnelling current of electrons, to a first approximation, can therefore describe the topography of the surface—showing the actual position of the atoms. In passing it should also be pointed out that there is a finite probability that atoms tunnel as well as electrons. This is important when considering the use of the STM as a low-energy means of moving atomic material. The relationship to topography should, however, always be remembered for what it is, a measure of electron state densities. If the bias voltage is changed then the so-called topography will also change. Measurement of state density indicates the number of electrons either present or allowed at the particular energy determined by the bias voltage. It is obvious that rather than being a disadvantage this can be turned to an advantage. The current can provide more than topographic detail. By exploring bias details important in chemical bonding, chemical composition and even crystalline structure can be obtained. This will be discussed in the section on spectroscopy.

(b) *The tip*

The big uncertainty in the measurement of surfaces is the tip of the stylus. At the atomic level it is very difficult to know exactly what is happening. The geometry of the tip is often uncertain and what happens to it should it hit the surface is even more unknown. In fact there is considerable dispute as to whether the tip is not always contacting the surface and that the contact is like a liquid-liquid interaction, the liquids being partly miscible electron clouds. That useful detail emerges is indisputable but great care has to be taken in its interpretation. Artifacts can easily occur, for example when 'point jumping' occurs as in figure 4.75. The two points on the stylus cause the artifact. Tip materials vary considerably and may be tungsten, platinum, iridium, etc; which one is selected depends on the sample and the intended mode of operation.

Whereas topography has been the principal aim when using the STM it does provide more information—the most important is spectroscopy. This will be briefly described.

### 4.2.3.4 *Spectroscopy*

The previous description focused on the capabilities of the STM to scan a surface and provide data on its atomic-level topography. This capability is only part of the story. The STM also opens up new possibilities in spectroscopy.

In the operation of the STM the bias voltage and tip-to-sample distance for measurement at any point can be varied. This capability enables the STM to perform two related types of measurements: *I/V* (current versus voltage) spectroscopy and *I/z* (current versus height) spectroscopy. In classical spectroscopy, radiant energy is applied to a material and the spectral response of the material is analysed. Spectroscopic data provide information about a wide range of properties, including chemical composition, atomic and molecular energy levels, and molecular geometry.

The STM provides spectroscopic data by utilizing its sensitivity to the electron energy states of the sample. In *I/V* spectroscopy, measurements of tunnelling current made under conditions of changing bias voltage provide data on local density of states. Figure 4.79 shows how a change in bias voltage changes the computed tunnelling current for the same sample.

In *I/z* spectroscopy, measurements of tunnelling current made under conditions of changing tip height provide data on the local barrier height between the tip and the surface. Under most conditions, local barrier height data correspond to the work function, a measure of the amount of energy (expressed in electron volts, or eV) required to remove an electron from its current state. *I/z* spectroscopy makes use of the variables in this equation, which can be considered as a first approximation of tunnelling current:

**Figure 4.79** Tunnelling current versus gap for STM.

$$I \propto V\exp(-Bz) \tag{4.122}$$

where $I$ is the tunnelling current, $z$ is the tip-to-sample distance, $V$ is the bias voltage (tip-to-sample) and $B$ is the square root of the barrier height between the sample and the tip, $\psi^{1/2}$. The barrier height and work function are proportional to the slope of the graph of current versus tunnelling distance. Equation (4.122) is another way of describing equation (4.121).

The STM brings a unique advantage to spectroscopic studies. Where other spectroscopic methods provide results averaged over an area of the surface, the STM can provide spectroscopic data with atomic-level resolution. This enables close studies of surface irregularities, defects, dopants (used in semiconductors) and other local surface phenomena.

Even topography, force and spectroscopic information is not enough to satisfy the demands made on this new type of instrument. The following are just some suggestions of potential applications. They are set down by manufacturers of STMs such as Burleigh.

(*a*) *Chemistry*
It is feasible to use the STM as an initiator of some chemical reactions. This is a possibility because the energy level of the tunnelling is about the same as that encountered in a wide range of chemical reactions. The barrier heights are only of a few electron volts.

(*b*) *Biology*
In the non-vacuum mode of the STM there is real potential for measuring actual growth.

(*c*) *Biochemistry*
Biochemical operations are possible owing to the possibility of positioning a small quantity of material on a surface.

(*d*) *Lithography*
The STM has been used to machine substrates. It seems that submicrometre mechanical etching could be achieved.

(*e*) *Microfabrication*
Given the very fine positioning capability of the STM it is clear that small parts of material could be moved around using the tip as a tool.

### 4.2.3.5 Some simple scanning systems

The principles and operation of an STM are surprisingly simple. As shown schematically in figure 4.80 an extremely sharp conductive tip (ideally terminating in a single atom) traces the contours of a surface with atomic resolution. The tip can be moved in three dimensions with an $x$, $y$, $z$ piezoelectric translator as indicated in figure 4.80.

Figure 4.81 shows a drawing of an actual STM head and base showing the essential components. Also depicted are three screws used for controlling the mechanical approach of the tip to the surface. Three keys to a successful STM design are (i) a smooth mechanical approach mechanism, (ii) rigidity,



**Figure 4.80** Basic movements of STM.



**Figure 4.81** Typical STM systems.

and (iii) convenience in changing sample and tip. Figure 4.81 shows a drawing of a combination AFM-STM in which the sample moves rather than the tip. This allows the delicate force sensor to be stationary. If the force sensor is replaced with a metal tip, the instrument becomes an STM. The mechanical approach system for advancing the tip towards the sample is similar to that shown in figure 4.82.

STMs have been operated in various environments, which vary from ultra-high vacuum down to $10^{-7}$ Torr and now air. Especially important are the applications in liquids, typically electrolytes in which biological specimens can exist.



**Figure 4.82**  Systematic diagram for height control.

For the case where the stylus is being moved the piezoelectric translator is driven by a typical high-voltage circuit shown in figure 4.83.

The fine positioning mechanism is produced by the three voltages. Figure 4.84 shows the $z$ movement whereas figure 4.85 shows the ingenious way that $x$ and $y$ are achieved by means of the same crystal.

To achieve nanometre positioning and control, the whole system has to be stable. There are a number of obvious ways of achieving this. The simplest of all is to make the system small. This has a number of beneficial effects. One is that the mechanical loop between the stylus and the specimen is small, which makes it more difficult for vibrations to influence any one arm of the loop. It does not matter if the whole loop is



**Figure 4.83**

**Figure 4.84** Movement of probe, typical arrangement: $z$ positioning by means of the crystal extension.



**Figure 4.85** Movement of probe, typical arrangement: $x$–$y$ positioning by means of the crystal deflection.

affected. Also, the smaller the units making up the mechanical loop, such as screws, levers, springs, etc, the higher the resonant frequency and the less chance there will be of them being excited by motors, gearboxes, drives, etc. They are also less likely to interact with each other. One reason why piezoelectric elements are used is because of their very high stiffness.

At the same time that the mechanical loop should have a high resonant frequency by being small and stiff, any suspension system should have a low one so that it can best 'short-circuit' environmental and instrumental low-vibration frequencies.

It is also advantageous to keep the whole system simple. The simpler it is the less likely it is that interaction can take place and the more chance there is of maintaining integrity of shape.

Also, if the mechanical loop is small the chance of getting temperature differentials across it is small; the whole system moves together in temperature. This leads on to the important point about thermal coefficients. The important point is that it does not matter if the elements of the mechanical loop have high or low conductivity; providing it is the same across the loop then temperature does not cause the loop to change shape. It is the shape change that causes the signal to drift, that is if one arm becomes longer than the other in the effective calliper. For many reasons mentioned earlier the loop should be as symmetrical as possible. Skid considerations are not applied to STM or AFMs; the property being monitored is different. Also skids are used for vertical integration and are hardly concerned with lateral positional accuracy.

(a) *Control system* (*figure 4.86*)

In the case of constant current the servo system relating the current from the probe to the piezoelectric crystal (or whatever means of moving the tip) should have as high a response as possible. Most surface instruments have a frequency response of about 500 Hz. Because of the small size the STM control system should be higher; 1 kHz is often used but there is no fundamental reason why it should not be much higher.

**Figure 4.86** Typical STM system constant/variable current.

### 4.2.3.6 *The atomic force microscope [50]*

Whereas the STM actually traces profiles of constant charge density at a particular value of energy determined by the bias voltage, the atomic force microscope measures atomic force. The relationship between truly atomic force and what is measured has to be precise. In essence it comprises a stylus, a cantilever and a means of measuring the deflection of the cantilever tip. The stylus is at the end of the cantilever. Atomic forces attract the stylus when it is very close. These forces are balanced by the elastic force generated by bending the cantilever. The bend of the cantilever when stationary is therefore a direct measure of the atomic force. There is another mode in which the beam is resonated; the atomic forces then effectively modify the spring rate of the cantilever as shown in equation (4.92) in the section on force measurement. The problem arises of how to measure the deflection. One method is to use an STM system riding on the cantilever as shown in figure 4.87.

In its simplest form the AFM measures just one variable, the normal force due to atomic forces interacting between the stylus and the surface. This is not dependent on the conductivity of the surface or the tip, so the AFM does not have the disadvantage of the STM, which is restricted to conductors. It measures forces both normal and frictional. By modulating the beam with a small vibration of the order of 20 nm and measuring the



**Figure 4.87** Atomic force probe using STM measurement system.

difference in the resultant deflections of the highest and lowest values, a measure of the surface elasticity can also be found.

The STM family of instruments have had a profound effect in the sense that the emphasis on surface measurement, as far as engineering is concerned, has moved from vertical geometry to that of lateral geometry. This has brought with it great emphasis on positioning methods and scanning techniques. This change has filtered down to conventional stylus methods where the use of scanning methods has been hitherto reserved for flaw detection and not surface structure. With the realization that 'lay' properties are now being recognized as being functionally important it seems that equal importance will be placed on all directions.

The next section deals with the more mundane but nevertheless important detail that has to be considered when designing and using stylus instruments.

### 4.2.4 Aspects of stylus instruments

This section deals with some of the more practical issues encountered when measuring surface topography. Most of the points raised are relevant to the tactile type of pick-up but the principles extend through all surface measurement.

### 4.2.4.1 Metrology

This is always a comparison between a test piece and a standard. Without this comparison control of the process is impossible and, without control, quality control is impossible. All metrology is therefore an enabling discipline.

### 4.2.4.2 Generation of reference surface

### 4.2.4.3 General

Surface metrology is concerned with two aspects of geometry: one is the way in which the workpiece differs from its intended shape; the other is the effect of process marks. In the first category the obvious shapes are

(i) straight line reference and flatness reference
(ii) curved reference
(iii) circular reference
(iv) aspheric and other references.

### 4.2.4.4 Straight line reference

There are a number of basic ways of generating a reference line or plane. One is by using a solid object which is straight and smooth as the reference. Another is to use an intrinsic reference. In this the test workpiece is used as a reference, e.g. by using a stick. Another technique is to use a reference movement, e.g. a linkage mechanism. Finally, within certain restraints, the reference can be generated from the measured data by mathematical algorithm.

### 4.2.4.5 Straight line generator

The classic way to generate a straight line datum is by using Whitworth's method of three pieces. By arranging a sequence of polishing operations each involving two of the three pieces it is possible to end the process by having all three pieces flat and straight and smooth (figure 4.88)

It is based upon the fact that when two nominally flat pieces are rubbed together they will always tend to produce convex and concave spherical surfaces—if the processes are suitably random. This is clearly beneficial if lenses are being made. This fortunate outcome is produced by an infinity of lucky events.

**Figure 4.88** Whitworth's method of flatness generation.

The 'central limit theorem' says that a Gaussian distribution will always result from many random operations, so in each direction phase effects produce independent probability densities; equation (4.123).

$$K_1 \exp\left(-\frac{x^2}{2\sigma^2}\right), k_2 \exp\left(-\frac{y^2}{2\sigma^2}\right), k_3 \exp\left(-\frac{z^2}{2\sigma^2}\right) \tag{4.123}$$

the outcome is the multiplication of these additive effects because each is independent.

$$\text{Outcome} \qquad \simeq K_1 \exp\left(\frac{x^2}{2\sigma^2}\right) \exp\left(-\frac{y^2}{2\sigma^2}\right) \exp\left(-\frac{z^2}{2\sigma^2}\right) \sim \exp\left(-\frac{1}{2\sigma^2}(x^2 + y^2 + z^2)\right) \tag{4.124}$$

The resultant process output depends on $x^2$, $y^2$ and $z^2$. But $x^2 + y^2 + z^2 = R^2$ (4.125) is the equation of a sphere—hence the spherical surface generation—irrespective of the actual process mechanism.

It should be added that if one piece is removed from the exercise then 2 spherical surfaces are produced, one concave and one convex, of equal radius of curvature and opposite sense.

This method can be used to generate curved reference surfaces. Unfortunately although the parts are spherical the radius of curvature is difficult to control. This method is used to make ball bearings. A rod is cut into many equal parts. These are put in a drum with an abrasive. The drum is vibrated and rotated vigorously for some time, after which the ball bearings are taken out—near perfect spheres.

A surface generated by the method outlined above can be used to generate a reference hemisphere.

### 4.2.4.6 Curved reference

This can be achieved by using a concave cylinder (figure 4.89) or having an instrument with an arcuate trace (figure 4.90).

In figure 4.89 using a concave cylinder can provide a variable curve simply by rotating it about an axis normal to the axis of the workpiece (i.e. changing the yaw). Care has to be taken to make sure that the axis of the specimen is in the same vertical plane.

Another method, which is shown in figure 4.90, is ingenious because the reference is an axis determined by the rotation of a hinge. The curve generated by this can be backed off by tilting the stage (i.e. the pitch). A big advantage is that the actual reference itself does not move during the measurement.

Circles, aspheric and intrinsic references are described in the appropriate sections in the book.

### 4.2.4.7 Pick-up configuration

Occasionally, a very simple practical problem occurs that can cause trouble in measurement. One such factor is shown in figure 4.91, where the pick-up is at an angle relative to the surface. If this angle is large some

**Figure 4.89** Variable curvature reference.



**Figure 4.90** Curved tracking reference.

considerable errors can occur. The error is in fact $\cos\theta$, so $A\cos\theta$, not $A$, would be measured as in the correct angular configuration shown in figure 4.91. The fact that these considerations only apply to side-acting configurations should not detract from the serious practical implications they can have. Most instruments, including STMs, use the side-acting mode of measurement.

(a)                                                                 (b)

**Figure 4.91** Probe working at angle to horizontal.

Crank effects need not be serious (figure 4.92). The length of crank is offset by the effective length of the arc. Sometimes both effects are unavoidable in applications involving complicated components.



(a)                                                                 (b)

**Figure 4.92** Effect of cranked stylus.

The effect of these offsets can be important in measuring the slope of a surface (e.g. see [78]), as in figure 4.93.



**Figure 4.93** Slope measurement with stylus pivot high.

This figure shows the distortion that can result from measuring the slope of the surface using a stylus raised above the level of the surface by $h$, where $a$ is the actual slope measured with an ideal stylus when motion is strictly perpendicular to the datum line, $\delta$ is the detected slope centred on an arc K, and $\varepsilon$ is the angle between the stylus tip and rotating point $\sim h/l$.

Ideally, the slope $\tan \alpha$ is BC/AB however the stylus does not move along BC but rotates around K.

The resultant angle measured will be tan $\delta$ = BE/AB from

$$\tan\alpha = \frac{1}{1/\tan\delta + \tan\varepsilon}.$$

(4.126)

This discrepancy is generally small, of the order 0.1% for $\varepsilon \sim 7°$. In fact the parameter which suffers most from this effect is the skew of the slope distribution.

### 4.2.4.8 Generation of the skid datum

The use of a skid in surface texture measurement is as old as the subject itself. Although it is generally accepted that the first use of a stylus method for examining surfaces was due to Schmalz in 1929 it is recognized that Dr J Abbott of General Motors built the first practical instrument in 1933. In the year 1936 skids were introduced by him! So their use is not without pedigree.

The skid datum was introduced in the 'Profilometer' around 1936. This instrument had a pair of skids on each side of the stylus, as shown below in figure 4.99(c). For the Talysurf instrument, a single skid (figure 4.99(b)) was introduced around 1940. It will be convenient to consider first the behaviour of the single skid.

The principle of the behaviour is shown in figure 4.94. The pick-up body is hinged to a driving shaft and assumed to travel in a substantially straight line. At the front end of the body the skid rests on the specimen surface. In this example the stylus is mounted on parallel leaf springs for vertical movement (i.e. normal to the specimen surface). Its vertical movements are detected by some kind of transducer. For the purposes of analysis it may be assumed that the skid and stylus trace the same cross-section although, in the Talysurf, the two are deliberately adjusted to operate in planes separated by 0.1mm so that the stylus does not trace the crests traversed and possibly flattened by the skid, such flattening, if anything, then being advantageous.



**Figure 4.94** Effect of skid radius.

The quantity directly measured by the transducer is the displacement $z$ of the stylus relative to the body considered at some point vertically above it, for example at the point O.

As the pick-up traverses the surface, the stylus rises and falls over the irregularities according to their height $z$; but the body itself may also rise and fall, at the front end due to the vertical excursions $z_s$ of the skid as it slides over the crests, and at the rear end due to imperfections $z_g$ in the driving device.

If the skid and stylus are respectively at distances $L$ and $L_T$ from the hinge, the output from the transducer will be

$$z_\theta = z_s - \left[ z\frac{L_T}{L} + z_g\left(\frac{L - L_T}{L}\right) \right]$$

(4.127)

all the movements in the same direction being given the same sign. The quantity in brackets represents an error in the required measurement of $z$. Its nature and possible magnitude must therefore be examined.

The path traced by the skid is described by the locus of its centre of curvature. This centre may be transferred to the instantaneous pole, this being the intersection of the skid surface with the direction of traverse. When the centre of curvature lies above the level of the hinge point, its excursions $z_s$ will be accompanied by a horizontal component, $z_s \tan \theta$, which is usually small enough to be ignored. The instantaneous pole will displace equally.

An important point is that since the skid generally has curvatures in both the axial and the transverse plane, it can engage crests lying on either side of this plane containing the path of its lateral centre, so that whether or not the stylus follows this path, there is no single cross-section from which the true locus of the skid can be deduced.

In some instruments the transverse radius is much less than the axial radius, in which case a close approximation to the true locus could be derived from a trace covering the path of this skid, noting, however, that in these instruments the stylus is often offset from the path of the skid. In other instruments a spherical skid is used, which may respond to information coming from a track of appreciable width. The major effects of the skid, however, can be studied by assuming that the transverse radius of the skid is very small, that skid and stylus follow the same track along the surface, and that the factor $L_T/L_s$ can be neglected. It will be permissible to ignore this effect of errors $z_g$ in the motion of the drive, as their contribution to $z$ is generally extremely small and, moreover, spread gradually over the traverse.

Considering now the tracing of a profile, the locus of the skid will be affected by its own shape and by the spacing, depth and shape of the irregularities; It will suffice for the moment to examine the locus of a skid of uniform curvature tracing a sine wave specimen of constant amplitude but progressively increasing crest spacing (i.e. wavelength), as. shown in figure 4.95. The locus changes from a substantially straight reference line to a shape which almost follows that of the specimen. The general case will be examined in a subsequent section.

The combination of the foregoing skid movements with those of the stylus is shown in figure 4.95. When the skid and stylus rest simultaneously on peaks or valleys, they run up and down together, that is in phase, and their peak movements subtract so that the indicated value of the specimen is too small. When they move out of phase their peak movements add and the indicated value is too large. For each parameter there



**Figure 4.95** Effect of skid-stylus separation: (*a*) skid behaves properly; (*b*) skid sees no roughness; (*c*) skid sees twice roughness. A, no effect—(1.0); B, various effects (0 (r) 2); C, skid follows stylus (0).

will generally be an intermediate phase for which the indicated value is correct, even though the waveform itself may be distorted.

It will be clear that in the present case the skid and stylus movements will be in or out of phase whenever the separation of skid and stylus contains an even or odd number of half wavelengths, the indicated value of $z$ rising to a maximum positive value when only a single half wavelength is included. When the half wavelength exceeds this value, the output will fall progressively through its correct value towards zero. The behaviour of the stylus and skid is shown diagrammatically in figure 4.95. The transmission characteristic of the combination of skid and stylus for the given input condition is shown in figure 4.102.

In the first section of the transmission characteristic, up to $\lambda_1$, the output represents the true value of $z$, the effect of this skid being negligible. For wavelengths greater than $\lambda_2$, the skid tends to behave as a mechanical wave filter attenuating these longer wavelengths. In the intervening section there is an alternation of attenuated and amplified wavebands, separated by wavebands where the output is but little affected by the skid. In this section the skid, if still regarded outwardly as a kind of filter, must be thought of as an indeterminate filter, quite different in its behaviour from the electric wave filter.

(*a*) *Use of skids*

This has been dealt with in some detail by Reason [79] for periodic surfaces. He showed that there are two basic problems, both obvious (when pointed out). One is the relationship between the actual radius and the spacings on the surface. The other is the relationship between the skid-stylus position and the surface.

Taking the former problem first, there are two types of surface that have to be examined, periodic and random. In both cases the criterion for success is that the skid does not penetrate significantly into the surface geometry as it is moved across it; the skid is intended to bridge the gap between asperities and not to fall in (figure 4.96).

As the skid radius gets small compared with the asperity spacing it gets progressively poorer as a skid, as shown in figure 4.96.

Reason [79] has plotted the effective filtering effect produced by a circular object transversing a periodic waveform and shows quite conclusively that it is a frequency filtering which also depends on the amplitude of the wave.

It should be noticed that the behaviour of the skid when passing over a periodic surface is not straightforward. There are two definite regimes shown in figure 4.96. The case (*a*) is when the skid does not bottom and (*b*) is when it does. In case (*b*) the skid is obviously doing nothing useful at all and can give rise to very large errors.

This effect has been studied in great depth by Shunmugam and Radhakrishnan [80] for two and three dimensions by means of a computer simulation. They found that the increased integrating effect of the 3D ball over the 2D radius run across a profile was considerable. In many cases the $R_p$ increased by more than 50%



Figure 4.96 Skid behaviour—large (*a*) and small (*b*) radius skid.

from the 2D to the 3D case. Similar calculations of the filtering effect of circles and balls on surfaces have been carried out elsewhere by Hunter and Smith [81].

Attempts to quantify the relationship between the skid radius and results from a random surface have been numerous and, because of the difficulty, not entirely successful. However, two ways are possible, one using the peak height distribution, the density of peaks and the sphere radius, and the other method which is perhaps more formal using the autocorrelation function of the surface, $R_p$ and the radius [82].

In the latter the basic idea was to find the average level to which the skid (either flat or curved) drops into the surface (as seen in figure 4.94). The distance between independent events is the independence distance of the surface determined from the autocorrelation function.

The method of calculation is as follows for a flat skid for simplicity.

The probability of a contact at point 1 in between heights of $z$ and $z + \delta z$ is

$$p(z)[P(z)]^{N-1}\delta z \tag{4.128}$$

where $p(z)$ is the density function and $P(z)$ is the distribution function

$$p(z) = \int_{-\infty}^{z} p(z_1)\mathrm{d}z_1.$$

A similar situation exists for points 2 and 3 and so on, so the probability of a single contact between $z$ and $z + \delta z$ is

$$p_c(z)\delta z = N[P(z)^{N-1}]p(z)\delta z. \tag{4.129}$$

Note that only one contact need be considered because the skid is constrained to move only with one degree of freedom. When N >> 1 this can be expressed as [82]

$$p_c(z)\delta z = N\exp\{-[1 - P(z)(N-1)]\}p(z)\delta z \tag{4.130}$$

for all $z$. The contact probability density for the skid can be written in less general terms for a specific height distribution, for example a Gaussian one. Thus

$$p_c(z) = (N/2^{N-1}\sqrt{2\pi})\exp(z^2/2)[1 + \mathrm{erf}(z/\sqrt{2})]^{N-1} \tag{4.131}$$

where it is assumed that the surface has an RMS value ($R_q$) of unity and a mean value of zero.

From equation (4.131) the mean height of the skid above the mean line of the surface can be found by integration. Hence

$$\bar{z}_s = \int_{-\infty}^{\infty} z p_c(z)\,\mathrm{d}z \tag{4.132}$$

which can be evaluated approximately by observing that for this distribution the mode and mean values are very close to each other in height.

Hence, calculating the mode and taking logarithms gives

$$2\ln(\bar{z}) + \bar{z}^2 = 2\ln(N-1) - 2\ln(\sqrt{2\pi}) \tag{4.133}$$

but because $2\ln z << \bar{z}^2$, if $z > 1$ equation (4.133) becomes

$$\bar{z} = \left[2\ln\left(\frac{N-1}{2\pi}\right)\right]^{1/2} \tag{4.134}$$

where there are $N$ independence distances on the surface covering the extent of the flat skid.

Equation (4.134) represents the extent above the surface mean line to which the skid falls on average. For example, if the independence distance were 10 $\mu$m and the flat skid 10mm then the value of $N$ would be 1000 and $z = 3.5$. The value of 3.5 represents the number of $R_q$ heights above the mean line to which the skid falls. This is an acceptable level. However, if the skid were 1mm the mean height of the skid would be 2.7 $R_q$, which is lower than $3R_q$ ($3\sigma$) and so is probably unacceptable. In this it is assumed that the $3\sigma$ value is roughly where the mean peaks are.

The use of flat skids is not always to be recommended because of their tendency to dig into the surface if the odd freak peak is in the way. For this reason curved skids, as seen in figure 4.99, are more often used. This situation is more difficult to handle because each of the $N - 1$ functions of $P(z)$ in equation (4.129) has a different upper limit of integration. Each upper limit of the integral has to be chosen to fit the shape of the skid.

There is, however, an approximate way that can be used to find the probable number of events which could contact the curve (figure 4.97). Assume that half the surface (e.g. that part above the mean line) could contact the skid (this height is $R_p$). Using the spherometer approximation to the circle the extent of the curve which could be influenced by the surface is $2x$ where

$$2x = \sqrt{8RR_p}$$

where $R$ is the radius of the curve.



Figure 4.97 Skid behaviour—random waveform.

The number of events is

$$\frac{2x}{\tau_L} = \left(\frac{8RR_p}{\tau_L^2}\right)^{1/2}.$$

(4.135)

Inserting equation (4.133) into (4.132) gives (after neglecting unity in $N$—1)

$$\bar{z} = \left[\ln\left(\frac{RR_p}{\lambda^2 2\pi}\right)\right]^{1/2}.$$

(4.136)

Notice in (4.136) that the two ratios determining the value of $z$ are $R_p/\tau_L$ and $R/\tau_L$, which shows that for a random surface also the effect of the skid is amplitude as well as spacing conscious.

Note also the difference between the results for the flat skid and the rounded skid. Equation (4.134) shows that the flat skid is not amplitude dependent whereas the rounded skid exemplified in equation (4.136) is.

In conclusion, there are criteria which can be used to determine the validity of use of skids even on random surfaces. There are also alternative ways of expressing the equations given earlier which might be easiest to evaluate from a profile graph so as to establish quality criteria. For example, the correlation or independence distance $\tau_L$ is not easy to judge from a graph whereas the average distance between crossings is. Therefore using $\lambda$, the average distance between crossings, equation (4.136) becomes

$$\bar{z} \simeq \left[ \frac{\sqrt{2}}{\pi} \ln \left( \frac{RR_p}{\lambda^2} \right) \right]^{1/2} .$$

(4.137)

$R_p$ and $\lambda$ can readily be estimated from a profile graph obtained by using a true reference datum and not a skid. Figure 4.98 shows some practical results of equation (4.137).



**Figure 4.98** Skid behaviour—random surfaces.

Having estimated the magnitude of the drop of the skid into the surface the question arises as to its importance for measuring surface parameters. Taking the random surface, peak parameters will be affected by the extent that $z$ is different from the highest peak in the traverse. Providing that $R$ has been chosen properly this difference will be small because the skid, by its very nature, picks out the highest peaks from which it establishes its datum. Errors of 15% can be expected.

Average values such as $R_a$ and $R_q$ are affected differently. The $R_q$ of the skid will be as random as the $R_q$ picked up by the stylus and because of the random and assumed independence of the skid-stylus movements the variances of both can be added together. Thus

$$R_q^2(\text{measured}) = R_q^2(\text{skid}) + R_q^2(\text{stylus}).$$

(4.138)

$R_q$ can be estimated crudely from the difference between $z$ and the extreme of the range of the surface which is taken to be about $3.5Rq$. For a 10mm skid the $z$ value is $2.5R_q$ giving a difference of about $1R_q$, so $R_q$ (skid) is for this example about one-third $R_q$(stylus).

Hence $R_q(\text{measured}) \simeq 1.22 = 1.05R_q$(stylus, true value), so even for a rather poor skid the error in the $R_q$ value is likely to be small. The same goes for the $R_a$ values of random surfaces. This is not true for periodic surfaces because the movements of the skid and stylus are not necessarily independent.

So far the effect of skid integration has been examined but there is the other factor still to investigate, which is the relative position of the skid and the stylus.

A number of possibilities have been tried. These are shown in figure 4.99. In the ideal configuration the stylus projects through the skid so that they both measure the same part of the surface and are in phase with

each other. Unfortunately, although this is a good configuration from the point of view of distortion, it is poor practically. This is due to the fact that debris congregates in the passage where the stylus protrudes through the skid with the result that clogging occurs. Alternative configurations are allowed (see e.g. [83]). Sometimes the skids are situated adjacent to the stylus (figure 4.99(*c*)). They therefore eliminate the debris problem but at the same time they both suffer from the disadvantage that the skid is not covering the same ground as the stylus. One could argue that this is not too serious on a flat part which is random, but it can be troublesome if the workpiece is not flat.



**Figure 4.99** Skid configuration on random surfaces.

Either of these situations causes a constant offset of the height of the skid relative to that of the stylus. For this reason configuration (*d*) is often used in instruments. This is fine for random surfaces provided that the distance *d* between the skid and stylus is larger than the independence distance of the surface—which it has to be because of the length of the skid. However, it is not ideal by any means. Remembering that one of the principal uses of the skid arrangement is to reduce setting-up time, it does not cope with tilt as well as it might.

Figure 4.100(*a*) shows a configuration where the tilt is virtually removed completely, whereas for (*b*) it remains as a constant height difference. Similarly, a curved part configuration, figure 4.100(*d*), removes it almost completely, whereas (*b*) produces a ramp signal.

The point is that configuration (*b*) is in effect a mechanical differentiator: a curve looks like a ramp, a ramp looks like a constant, and so on. The reason why the ramp or constant value has not been a severe problem in surface texture measuring instruments is because the output from the transducer has been input through an electrical or digital filter having small wavelength transmission to the meter. For the purposes of electrical theory the ramp corresponds to a sine wave of very long wavelength and the constant value corresponds to a cosine wave of infinitely long wavelength. This is why two *CR* networks were used to isolate the roughness from the waviness and not one. If only one *CR* network had been used, it could block the DC level corresponding to a tilt problem in set-up, but it could not block a ramp. This would emerge as a constant voltage into the meter which would give a totally false $R_a$ value. Two stages are needed in cascade to preclude voltages resulting from curvature set-up errors getting to the meter. Basically, three stages of differentiation are needed to get rid of long-range curves on the specimen using configuration (*b*). One is provided by the skid position relative to the stylus, the other two by the electrical or digital filters (figure 4.101). Thus, if skids are to be used like this with conventional meters, the transmission characteristic of the filter has to be at least 12 dB/octave or 40 dB/decade.

The foregoing comments relate to long-wavelength geometry caused by the shape or tilt of the component. What about the effect on the ordinary roughness of a periodic character? This was explained originally by Reason [79]. A summary is presented here for completeness.

The situation is shown in figure 4.102. When the skid is integrating properly it is then acting as a perfectly respectable datum (region A). Once it starts to drop seriously into the gaps between asperities (as in B) then the stylus position relative to it, *d* (figure 4.102), becomes important.

**Figure 4.100** Behaviour of skid systems to slopes, (*a*) and (*b*). Behaviour of skid systems to curves, (*c*) and (*d*).



**Figure 4.101** Use of skid system and electrical filter to remove curves.



**Figure 4.102** Phase and amplitude variation with skid-stylus separation.

Two bad cases should be singled out. In one case the skid and stylus can move in complete antiphase. Under these circumstances the calliper will give a reading of twice the real value of the waveform. Alternatively, the probe and skid could be exactly in phase so that the two arms of the calliper are synchronous with the result that no signal is seen at all; these are shown in region B. For very long wavelengths, the signal is considerably attenuated, as shown in region C (figure 4.102).

The swings in amplitude from low to high values of transmission are, in principle, a function of the ratio of $d$ to $\lambda/d$ where $\lambda$ is the wavelength of the period. The fact that it only ever reaches near to its theoretical extreme values once or twice is due to the finite radius effect.

A vitally important point is that in the region between $\lambda_1$ and $\lambda_2$ the value and sign of the error will be unknown unless full details of the specimen are available. Since such details can be obtained only by sufficiently tracing the specimen relative to a true datum, which when done renders the skid trace superfluous, the region from $\lambda_1$ to $\lambda_2$ has no practical value and is liable to give grossly misleading results. A basic difficulty is that apart from not knowing the value or sign of the error, it is often difficult, even with experience, to judge where $\lambda_1$ and $\lambda_2$ begin. The spacing of periodic specimens with spacing from perhaps 1/4mm upwards can sometimes be measured well enough against a scale, but even then there may be some uncertainty about the effective separation of skid and stylus. Quite a small change in the inclination of the pick-up may roll the point of contact around the skid sufficiently to pass from an in-phase to an out-of-phase representation of a closely spaced surface.

When the profile has an irregular shape, as is typical of abraded surfaces, the separation between one crest and the next one engaged by the skid may fluctuate over a wide range. The vertical excursions of the skid will then fluctuate both in magnitude and phase, the skid error exhibiting both positive and negative values as the traverse proceeds.

Although some crests may be increased and others reduced in apparent height, the general appearance of the profile, and its $R_a$ value, may be but little affected. Thus the skid effect on random features is often of much less consequence than on periodic surfaces. There can be little doubt that it is because of this, and of the preponderance of random textures in the general run of industrial control, that the skid has become so widely accepted. This should also not be forgotten for those cases where the surface is not really random or periodic, as shown in figure 4.103. In this case a piece of debris, which the skid is much more likely to pick up than the stylus, can cause considerable distortion.

However, it should be said that sometimes the use of a blunt stylus instead of a sharp one can be an advantage because if a rare high spot on a disc or film is being sought, the blunt stylus can in effect search a lot of the surface area to find it, whereas to achieve the same with a sharp stylus would require considerable scanning.

Straight line datum             Graphical output

Skid datum             Graphical output

**Figure 4.103** Skid distortion.

### 4.2.4.9 Stylus instruments where the stylus integrates

(a) *Form measurement*

Unlike that for surface texture the stylus used in form is usually blunt. Some mention has been made of such styluses in the earlier section dealing with miniaturization. The biggest problem area is in measuring roundness, cylindricity and similar forms. The stylus is one means by which the much shorter-wavelength texture effects can be attenuated (figure 4.104).

**Figure 4.104** Form—stylus showing integration of texture.

A blunt hatchet or ball stylus is most often used to remove roughness effects around the circumference of the part. If a very sharp stylus is used (figure 4.104(*b*)) and the workpiece has been manufactured by using a single-point cutting tool, there is a considerable danger of the stylus dropping in and out of the helix produced by the tool, giving a distorted roundness graph. This effect produces a signal which is not to do with roundness but only with the process. For random surfaces like grinding it is not so serious a problem. The widely used, blunt hatchet stylus prevents this happening but causes other problems if the workpiece is tilted, as shown in figure 4.105.

The error due to the tilt is present even for a sharp stylus, as shown in the comments on the limaçon cylindroid, but it is considerably accentuated by the large radius of the hatchet. Reason [79] pointed this out carefully in his 1966 report. Unfortunately the use of integrating hatchet or toroidal styluses is becoming more restricted because of the need to measure small objects. Another factor often missed is the very penetrating effect of the stylus when measuring externally circular pieces.



**Figure 4.105** Effect of tilt of specimen on reading.

In the case of figure 4.105, if $r$ is the radius of the part, $R$ the hatchet radius and $\theta$ the tilt then the effective change in diameter for a tilt of $\theta$ is approximately $r\theta^2$.

For a hatchet stylus this is modified somewhat by the finite size to be

$$\delta_{\text{diameter}} \simeq (r + s)\theta^2 \tag{4.139}$$

for small $\theta$.

**Figure 4.106** Involute pick-up/skid showing constant contact position.

A stylus/skid can have an involute shape, similar to gears, which makes the point of contact insensitive to specimen or stylus/skid tilt as shown in figure 4.106.

Many investigators think that a flat stylus will integrate perfectly whereas in fact it is the relative curvature of the stylus and the workpiece that determines the penetration (figure 4.107).



**Figure 4.107** Illustration of penetration with flat probe.

### 4.2.4.10 Alignment of the stylus system

Viewed in cross-section, as in figure 4.108, radial displacements of the stylus should act along a line intersecting the axis of rotation (*a*). An offset *d* of the stylus (*b*) will make the real radius change of $\delta r$ look like $\delta r \sec \theta$ where

$$\theta = \sin^{-1}(d/r).$$

(4.140)

If instead of a sharp stylus one with a radius on it were used, then the apparent change in radius $\delta r'$ is given by

$$\delta r' = \delta r \sec \theta$$

(4.141)

**Figure 4.108** Effect of offset stylus.

but now $\theta$ is given by $\sin^{-1}[d/(r—R)]$ so that the effective increase in magnification caused by this effect is greater than for a sharp stylus.

Excessive values of $\theta$ will affect not only the magnification but also the magnitude of the radial component of the stylus force and, according to Reason, the friction and drag will become a complex and eventually unstable function of the offset.

In the case of very small holes it is difficult to avoid letting the diameter of the stylus become an appreciable fraction of the hole, in which case $\theta$ has to be kept to at most 10°.

### 4.2.4.11  *Limitations of 'references' used in roundness measurement*

Some of the limitations of workshop methods of measuring roundness have already been discussed in Chapter 2. Also the basic radial configurations of rotating spindle, rotating table and master disc methods have been presented briefly in the theory. Here it is necessary to get slightly more detail of the practical advantages and disadvantages of the various methods. This is an example of where problems with the 'reference' and the 'test' arms of the measurement 'calliper' become entwined.

(*a*)  *Rotating pick-up method*
The point A in figure 4.109 describes as nearly as practicable a perfect circle in space as the spindle rotates. How well it does this depends only on the bearing and not the workpiece load. Hence the measurement of roundness by this means is likely to be very accurate. But suppose that the concentricity between two different sections on the part needs to be measured. There are two ways of tackling it. One way is to measure the centre position at $S_1$, move the quill down so that the same stylus measures at plane $S_2$, and then measure the



**Figure 4.109** Rotating pick-up method for roundness measurement.

effective centre of the part from the roundness graph again. The difference between the two centre positions in principle would give the concentricity as defined in section 2.3.3. However, this presupposes that when the quill is moved within the bearing, the centre position of the axis of the spindle is not changed relative to the workpiece. Of course, there is no guarantee of this unless the quill is absolutely straight. In practice, the correlation between the centre positions is poor unless special precautions are taken.

Another possibility is to keep the quill fixed within the bearing and have the use of a stylus 'tree' (figure 4.110). The graph at $S_1$ is taken using one stylus and then another stylus already in the plane $S_2$ and on the same shank is engaged on the part at $S_2$. Now there is good fidelity because the part has not been moved relative to the axis of rotation or vice versa, but because it is angular deviations which are being measured, the sensitivity of the pick-up in position $S_2$ has changed by the factor $l_1/l_2$, so this has to be taken into account. At one time this was a considerable problem because the compensation had to be carried out manually, but with computers this can be done automatically.



**Figure 4.110** Tree stylus for measuring concentricity.

This type of instrument therefore needs to be used carefully if concentricity, or squareness, is to be measured. On the other hand, if taper is being measured, there is no problem even if the quill does move because, providing that the shank position of the stylus radially from the centre of rotation is not changed, a change in radius of the part in two or more planes will still be seen, albeit with each trace possibly at a different centre.

The exact opposite is true for the rotating-table instrument (figure 4.111).

It does not matter to a first order what happens to the stylus as it is moved along the column from position $S_1$ to $S_2$ if concentricity is being measured. Again this is because, in doing so, the relative position of the workpiece to the axis of rotation has not been moved. Therefore the rotating-table method is ideal for concentricity measurement. Notice also that the sensitivity of the probe has not changed in going from $S_1$ to $S_2$, unlike the use of the 'tree' in the rotating pick-up instrument. However, measuring taper by this method can be difficult. Care has to be taken to make sure that the column is parallel to the axis of rotation of the table.

The column can be corrected by means of the standard trick of putting a true cylinder on the table and taking a straightness graph up its axis without the table turning. This is repeated only in the second place, turning the specimen 180° and using the probe in the reverse sense, as shown in figure 4.112.

The position shown in (a) gives a graph which is tilted at $\theta_1 + \theta_2$ whereas in (b) it will be $\theta_1 - \theta_2$. Subtracting the graphs and halving gives $\theta_1$ which can then be adjusted out. It does not matter too much about $\theta_2$, the component tilt. Reasonable accuracy can be achieved in this way.

Another point to be watched is the straightness of the column. As can be seen in figure 4.113, an apparent increase in radius of the part has resulted from angular (pitch) errors in the slideway.

**Figure 4.111** Rotating-specimen method for roundness measurement.



**Figure 4.112** Error inversion method.

The effective error produced is

$$\delta r = l(\sec\theta - 1) \sim l\theta^2 / 2. \tag{4.142}$$

The distance from the column counts!

### 4.2.4.12 Other stylus methods

In the foregoing section instrumentation based upon the continuous tracking of a stylus across a surface has been considered. This does not mean to imply that this is the only way in which tactile stylus methods could be used to estimate surface roughness and form. One obvious variant is to use the stylus pick-up in a discrete rather than continuous mode. In this method the workpiece (or pick-up) is incremented by discrete equal steps along the surface. At every halt the pick-up is lowered to touch the surface and the depth of the surface measured, almost enacting the way that a height gauge for measuring depth is used.

**Figure 4.113** Effect of lack of straightness of column.

Usually a datum level is incorporated, as before, either independently or by means of feet. One such device is simple and handy. It just takes single-point measurements of the maximum peak-to-valley reading across the workpiece as shown in figure 4.114.



**Figure 4.114** Single-point measurement of surface.

Usually the system comprises three equilaterally spaced feet surrounding the probe. Obviously this technique is very crude, and many readings should be taken to ensure a reasonable measure. To get $R_a$ from these point-to-point readings a factor of 4 should be used on the highest peak-to-valley measurement found. This factor derives from the ratio of the peak-to-valley height to the average value of a triangular waveform. It is intended only for very rough checks.

Another device ingenious in conception but not very convincing in action is based upon an estimate of the coefficient of friction by a blade moving across the surface. The idea is that when the blade buckles there is a relationship between the angle of the probe and the surface finish (figure 4.115) [84].



**Figure 4.115** Rubert artificial nail method for roughness measurement.

The theory of the buckling blade is simple in principle and is made up of two facets:

1. The flexible blade behaves structurally and therefore collapses when the tip touches the slopes on the surface at 90°—all the force on the blade goes into bending it.

2. If there is a significant coefficient of friction then the angle could be less than 90°, the angle of friction offsetting the deviation from the 90° (figure 4.116).



**Figure 4.116** The principle of the buckling blade method.

Unfortunately the frictional force is likely to change with small speed variations or the presence of contaminant so the method is not absolute.

Recent versions of the point-to-point surface measurement systems have their advantages [85] (sometimes more imaginary than real). This is mostly because the dynamic problems associated with stylus lift-off are removed. This has the added effect of having more actual control over stylus forces on the surface, which could be important in certain applications. Also the actual stylus surface contact is reduced. On the other hand the measurement cycle is invariably slowed down.

As indicated before, the real issue is the balance between having a steady reliable reading in the intermittent mode, or speed using continuous measurement. The great danger is that it might be tempting to sample too sparsely using the former method. Strictly speaking, to get equivalent coverage of a surface it should be sampled twice in every stylus tip width, which is hardly practicable.

The scheme shown in figure 4.117 represents a typical system [86].



**Figure 4.117** Typical step-by-step system.

In this system a fine degree four-phase variable-reluctance step motor is employed to increment the traverse stage between measurements. The motor is controlled with a driver card. In this particular configuration there is a fast reverse facility to save time. In addition to movement in the one direction, an arrangement is made to make the table move in an orthogonal direction to allow a mapping of the surface (more about this in the next section).

Alternative drive mechanisms using piezoelectric elements etc will also be considered but only briefly.

Some stylus techniques have been devised to work on a point-by-point basis for given processes only. One such technique has been devised by Deutschke [85], an outline of which is shown in figure 4.118.



**Figure 4.118** Deutschke method for in-process measurement.

This has been developed specifically for the in-process measurement of texture during grinding. The actual sensor consists of a drum which contacts the workpiece and is forced to rotate by virtue of its contact with the workpiece. Protruding through the drum wall is a stylus attached to an electromagnetic transducer. The stylus is free to move radially. As the workpiece rotates the drum rotates and centrifugal force keeps the stylus at the extremity of its range. Once per revolution of the drum the stylus hits the workpiece and is forced inwards radially by an amount determined by the surface texture. This incursion is noted by means of circuitry which detects minima. After a few hundred revolutions of the drum a series of sample data points representing the texture is obtained, from which the amplitude probability density function $p(z)$ is obtained and $R_a$, $R_q$ and the skew and kurtosis can be evaluated. The reference from which the measurements are referred are the two rims of the drum, which contact the workpiece.

One advantage of this method is that it is robust and can, under certain conditions, deliver information whilst the part is being machined. It has problems, however, with the surface marking owing to the relatively high impact force of the stylus. Note that one good feature is that the relative circumferential velocity of the stylus and workpiece is zero at the time of measurement. Obviously using just one stylus is slow; more than one stylus around the circumference of the drum is a possibility for increasing the data capture rate. Another point to note is that only height information about the surface is obtained and nothing about the length properties. Furthermore, roundness errors can intermingle with surface texture but this might not be a bad thing from the point of view of certain functions. It is important when using a device like this to make sure that there is no possibility of an arcuate signal pattern being generated; for example, the size of the workpiece should not be an integral multiple of that of the drum, otherwise repeats are possible (although unlikely because the workpiece is continually changing size, albeit by a small amount). Also, using this in turning is

not recommended because there could be a 'beat' effect between the length of the circumference of the drum, the workpiece and the feed of the tool along the axis of work.

### 4.2.4.13 Replication

There are situations in which it is difficult to get to the surface in order to measure it. The part may be too large to get even a portable instrument to it or the part may be needed for use. In other circumstances the component may change (i.e. due to wear) and some record of the original surface is needed. In such cases a replica is used [87, 88] which gives a 'negative impression' of the surface geometry.

Surface replica materials fall into two categories. One is a chemical which cures when mixed with another, and the other is a thin acetate-type film which is used with a solvent.

The requirements for a replica material are:

(1) high fidelity between the roughness on it and on that of the surface recorded;
(2) replication does not damage the recorded surface;
(3) debris is not left on the recorded surface;
(4) the replication process is quick, permanent and safe.

Consider first the compound method. The first step is to put some kind of release agent on the surface to allow easy removal. This should be minimal so that valleys are not filled. Preferably the release agent coverage should be molecularly thin and should not react with either the replica material or the surface under test.

The next step is to build a ring barrier around the area of interest with plasticine. The chemicals, which are typically some kind of Araldite compound, are mixed with the curing agent and then poured within the barrier area and left to cure, sometimes helped by gentle heating. The length of time taken to cure depends on the chemicals and can vary from a few minutes to a few days. To ensure a good, quick release the barrier should be shaped as shown in figure 4.119.



**Figure 4.119**

The following very simple strategem suggested by George [87] saves much time and frustration.

For general use he recommends Acrulite. This has a quick curing time of 20 minutes. For cases where form as well as texture is required Araldite CY219 with a mixture of 40% aluminium powder is recommended, presumably to facilitate heat movement in the exothermic reaction. The problem here is the long curing time, which can extend into a few days.

Generally, in practice it has been found that this sort of replication method preserves the longer wavelengths on the surface despite some evidence of bowing when the curing takes place. Some materials that have been used here include Acrulite Microtech A, which is a polymethylmethacrylate resin, Araldite CY219, a casting epoxy resin, glassfibre resins, etc.

The alternative method, using various thicknesses of film, has certain advantages, one of which is that the films can be made transparent. This makes the evaluation of the roughness possible using a liquid gate and diffraction methods. Typical films are cellon or acetyl cellulose. These are used in submillimetre thicknesses

of about 0.05 mm, 0.1 mm and 0.2 mm. The solvent, usually acetone or methyl acetate, is smeared over the specimen, and then the film is lightly pressed onto the surface making sure that any trapped air is removed. Using this technique and applying cross-correlation techniques it is possible to get high fidelity—especially using Cellon. Also, medium-thickness film performs better than the very thin (0.05 mm). Some problems can occur with elongation when the film is peeled off.

   In summary, both methods are useful—the former for higher fidelity and the latter for ease of use and for getting a transparent replica. It seems that despite all the high technology of today there is still a need for such techniques (figure 4.120).



**Figure 4.120**

### 4.2.5   Area (3D) mapping of surfaces using stylus methods

### 4.2.5.1   General problem

   The problems of measuring the profiles of surfaces are negligible when compared with those encountered in areal mapping. Methods which involve purely analogue methods will not be discussed here for the simple reason that the data cannot be easily manipulated.

   The basic problems in areal mapping are common to all instruments and are as follows:

   (1)  maintaining an accurate height datum between tracks;
   (2)  maintaining an accurate spatial trigger for the digital samples on each track relative to the others;
   (3)  incorporating suitable adjustments for roll, pitch and possibly yaw;
   (4)  adopting a suitable numerical model;
   (5)  adopting a sampling pattern which enables enough area to be adequately covered in reasonable time and cost;
   (6)  using suitable algorithms to develop a contour and reveal features of functional interest;
   (7)  maintaining sufficient resolution to detect flaws and other non-standard features.

### 4.2.5.2   Mapping

This section will be concerned first with the conventional use of mapping in which a stylus method is used. Because of the nature of flaws the first serious attempt to map surfaces was made in the pioneering work of Williamson [89] (figure 4.121) and other papers ([90], for example). He pointed out correctly that although the earlier types of interferometrical methods were excellent for stepped surfaces or those with local irregularities within large flat planes, they are difficult to interpret and operate on normal rough surfaces. Progression to electron microscopy revealed for the first time the very fine detail, but there were also, in those days, two major limitations, one being the difficulty in getting quantitative data and the other the fact

**Figure 4.121** Map of surface obtained by stylus scan.

that, because of its limited field of view, electron microscopy tended to mislead scientists to concentrate on the individual behaviour of asperities rather than to consider the overall properties of the area.

The team at Burndy Corporation realized at the outset that one of the basic problems would be that of maintaining a true datum between traces. They achieved this by referring every trace to a flat polished land on the front and back ends of each track (figure 4.122).

The second innovation was to cross-correlate each track with the previous one to obtain spatial correlation of the data points.

Thus if cross-correlation of the data of profile $z_1$ with $z_2$ produced a peak of correlation shifted by $S$ from the origin, the data points were shifted by this amount to bring the two into spatial coherence. So by two simple, yet elegant, steps they were able to adjust the data to ensure that a true map had been produced. Obviously, referring to figure 4.122(b), if the specimen is wedge shaped in the $z$ direction it would show up. They then used a best-fit regression line to remove such a tilt from the data.



**Figure 4.122** (a) Areal measurement showing datum; (b) cross-correlation to keep $x$ fidelity.

These early experiments were obviously subject to certain restrictions. For example, if the surface had any pronounced lay then the cross correlation-method could not be used because the shift $S$ might well be genuine—produced by the angle of the lay within the distance $\Delta y$.

This and other earlier attempts to map surfaces were successful but very time consuming and utilized equipment which was not, quite frankly, ever built to make such measurements. Good experimental techniques overcame many of the inherent problems. Many other experimenters have since used such techniques. Notable among them are Sayles and Thomas [90]. They again used the same basic stylus instrument as Williamson but more emphasis was put on the software to produce a very convincing picture of the surface such as shown in figure 4.123 and which demonstrated for the first time profile blocking in the far scene. Sayles was also the first to use the spline function to smooth the curve, thereby considerably improving the visual impact of the picture.



**Figure 4.123** Areal picture of surface showing profile blocking (after Sayles). (*a*) Without blocking (*b*) with blocking.

### 4.2.5.3 *Criteria for areal mapping*

What are the necessary criteria that need to be followed?

(*a*) *Mechanical integrity between parallel* (*or radial tracks*) (*figure 4.124*)
The mechanical datum between $z_1$ and $z_2$ must have an uncertainty in height which is small compared with the actual value of the texture and form to be measured. A typical accepted figure is $R_q/10$. Remember that it does not matter if the mechanical datum is consistently higher (or lower) with successive tracks in the *y* direction (or with the *x* direction) because the trend can always be removed by fitting a best-fit plane to the measured data. What cannot be tolerated is uncertainty. From this criterion it is clear that the smoother the surface the more difficult it will be to measure. In the case of Williamson the surface was comparatively rough so that the datum used (the lapped flat) needed to be no better than that which is easily achievable using conventional techniques.

Unfortunately, nowadays, the requirement for measuring fine surfaces is growing and this mechanical correlation is becoming more difficult.

One such design is shown in figure 4.125. This has been designed to provide the mechanical coherence required between tracks. The most demanding mechanical element in the set-up, apart of course from the surface finish instrument, is the stage. This has been designed to provide the mechanical coherence required between tracks. Basically it comprises two precision glass blocks A and B. Both are constrained

**Figure 4.124** Critical dimensions for areal fidelity.

by five contacts allowing one degree of freedom. B is a 10:1 tapered block used for moving the block A via the ball H by small known increments. This is achieved by turning nut C. A turn through one scale unit on the periphery of C corresponds to a movement of A of 2.4 $\mu$m.

Although block A has five constraints imposed upon it by means of ball contacts, two of these are variable, via nuts E and F, to allow adjustments in the tilt of the specimen (which is placed upon A) in the roll and yaw angles. The roll adjustment includes the device of a ball-on-screw E pushing onto an offset ball in the vertical slot of D. The bottom face of A and the two edges have had to be specially manufactured, for squareness, straightness and smoothness, to ensure that no unwanted rotations are produced when the block is moved transversally relative to the direction of motion of the pick-up. Typical specifications for the block are straightness 0.1 $\mu$m per mm, squareness and smoothness 0.05 $\mu$m.

The other angular adjustment, pitch, is provided for by direct adjustment of the surface finish instrument datum, which is specifically designed to give this feature for levelling a profile graph. The accuracy of, and operation instructions for, this adjustment are well known and are included in the manufacturer's literature. Also needed is a device designed to lift the stylus just clear of the surface after every traverse. Basically it consists of a precise way of lifting the 'mechanical steady' of the datum attachment by such a small amount that the ligament hinges on the traversing unit are not strained when the gearbox is racked back. A clearance of just 500 $\mu$m is maintained.

Standard tests have to be made on such apparatus to check whether or not the minimum specification is being met. The first test is repeatability. A number of traverses in the same position should be digitized and analysed. This tests three things: the vertical mechanical coherence of the surface-measuring instrument, the traversing speed, and the repeatability of the digital triggering of the start of track. Relatively rough surfaces



**Figure 4.125** Kinematic arrangement and lateral movement: (*a*) plan view; (*b*) front elevation, side removed.

can be used in the verification because the instrumental specification, both in terms of vertical precision and digital triggering, could be relaxed, as the errors are a correspondingly smaller fraction of the dimensions involved. The sampling interval is normally fixed at 2.4 $\mu$m, this being the nominal stylus tip dimension in the direction of traverse. Successive repeat traverses can be cross-correlated to establish whether the digital initiation signal from the photodiode is satisfactory. This means that the triggering error should be much less than the spacing between samples (i.e. 2.4 $\mu$m). This can be an impossible criterion because of the lack of repeatability of the light levels needed to trigger the fibre optic relay and the physical diameter of the fibre pack itself, which was 2 mm. Tests with the use of cross-correlation showed that errors of the order of 5 $\mu$m or more could be expected in synchronization. For this reason, when building up the grid of ordinates, each track is repeated three times and averaged to reduce the effect of this spatial discrepancy.

Tests can be made on the line of traverse of the stylus by lightly carbonizing a microscope slide and tracking the stylus over the same place on the side. An overall track width of twice the diamond width is encountered after a large number of repeats and shows the lateral repeatability.

### (b) *Accuracy of sampling*

The accuracy of initiation of each digital sequence $\Delta x$ is very important in areal assessment, although not in profile digitization. Failure to get good register can produce severe distortions (figures 4.124 and 4.127 illustrate two points). Conventional digital sampling plans use a rectangular grid. If, for example, the triggering of the second profile is in error by $\Delta x$, which is of the order of $h/2$, the whole sampling plan can be changed, as will be seen shortly. This results in a different expectation of finding a summit or measuring its properties as was shown in the earlier sections on areal characterization in chapter 2 and processing in chapter 3. From figure 4.124 the maximum allowable shift is less than $h/4$. Just how serious this is in practice depends on how the spacing $h$ compares with the correlation length of the surface. The nearer $h$ is to $\tau_L$ the more serious the error. Typical values of $h$ are 1 to 2.5 $\mu$m so that the acceptable starting error is likely to be 0.2 to 0.5 $\mu$m, which is tight, especially for practical instruments.

The main general mechanical issues are given below:

1. The mechanical stage and sampling procedure must be such that the variation of spacing (or sampling) interval throughout each traverse should be at most only a fraction of the interval. For example, for a sample interval of $h$ the instantaneous error should not be larger than $h/4$.
2. Drift in the determination of the sampling interval would have to be kept low. This implies that the rate at which the data is acquired has to be considerably faster than the drift, otherwise the whole grid, of whatever pattern, will become distorted. This is often due to a lack of thermal equilibrium, which occurs when motors are running in the scanner. Half an hour at least should be allowed for the equipment to stabilize.
3. The method used to initiate the pulses at the start of each traverse has to be at least as good as in 1, otherwise successive traverses become shifted in space relative to their intended position. This requires careful design of the trigger which, in critical cases, is beyond the capability of normal interferometers.
4. Hysteresis of the stage has to be restricted to the level of 1 or the tracking has to take place always in the same direction. This is because the effective origin is shifted owing to a change in direction of the frictional force. Backlash or 'looseness' is invariably overcome by means of springs, and is not the same thing. Values of hysteresis and backlash can now be reduced to about 0.3 $\mu$m with precision drives. Repeatability of tracks after zeroing should be better than 0.2 $\mu$m.
5. The drive to the sampler should be sufficient to push the specimen without significant elastic deflection. Typical allowable stiffness values are 1 $\mu$m kN$^{-1}$.
6. The guide, to ensure one-dimensional movement in one axis, has to be free from local slope error. Typical maximum allowable values should be of the order of $\mu$ radians. Furthermore, the measuring system should be as near as possible in line to avoid Abbé errors.

7. For two dimensions the two axes have to be square to each other so that, for large $x$ and $y$, the grid does not become badly out of phase. As an example, for a 1 cm square sample the squareness should be better than 0.25 $\mu$m. (The errors tend to be small, $10^{-6}$ $\mu$rad.)
8. Also in two dimensions, the top-stage weight on the bottom stage cannot be ignored. The bottom-stage rigidity has to be greater than the upper. Furthermore, it has to be remembered that as the top one moves over the bottom a cantilever effect can be produced in certain drives, which increases with distance from the origin of the coordinate system.

There are existing devices that have been designed to give performance levels such as described [91, 92]. Commercial devices also exist which increment to 0.1 $\mu$m at various speeds.

Another example of a 3D surface-measuring system is shown in figure 4.126 after Tsukada and Sasajima [91] and others. Typical specifications of this rig are given in table 4.5.

**Table 4.5**

|  | $x$ direction (stylus) | $y$ movement (table) |
| --- | --- | --- |
| Minimum step ($\mu$m) | 0.2 | 0.25 |
| Sampling intervals ($\mu$m) | 1–99 × 0.2 | 1–99 999 × 0.25 |
| No of samples | 1–9999 | 1–999 |
| Maximum measurement length (mm) | 100 | 25 |
| Moving velocity (mm s$^{-1}$) | 0.1 | 0.025 |

### 4.2.5.4 Movement positions on surface and sampling patterns

A typical scan pattern is shown in figure 4.127. Note that all possible precautions should be taken, such as preloading the stage to eliminate backlash.

So far the only sampling pattern to be considered is the rectangular grid. The question arises as to whether there is an alternative. The factors to be decided on are coverage, cost, speed, etc. It may be possible to get considerable advantage. Some possibilities will be considered here but obviously a discrete measurement of a surface will lose information because of the numerical model, as explained earlier in Chapter 2.

The alternatives are shown in figure 4.128, which shows the different starting points and possible differences in sampling intervals (not necessarily equal) [93, 94].

Although the issues of starting point and maintenance of accurate sampling are common to all the sampling schemes suggested here, they do appear to be slightly more critical for the patterns that require staggered starting points for successive tracks, as in the case for $k = 3$. For these patterns, and also for simplicity of construction and for cost, it is preferable to have the $x$ and $y$ movements basically the same. This can pose problems because of the rather messy intervals which have to be generated by the same device. For example, in the hexagonal case if x and $y$ are to be sampled spatially, movements of $h/2$ in the $y$ direction and samples of $(\sqrt{3}/2)$ generated with zero offsets of $(\sqrt{3}/2)$ on some of the tracks are required. This is not the simplest of tasks. To cater for all these conditions requires a drive which has a very small unit increment and can be expensive.

Another factor is wear and tear; although the pattern for $k = 3$ has fewer data points it can take more tracks to cover the same area with the attendant increase in wear of the probe. This has to be taken into account over a long period.

The basic factors other than movements which have to be taken into account when considering which sampling pattern is to be used are as follows:

1. The tracking length, because this determines the stylus wear in the case of a stylus instrument. It also determines the total time to gather the data.
2. The number of ordinates per unit area, because this determines the data storage.
3. The efficiency of the pattern in terms of information content.



**Figure 4.126** Surface mapping system (after Tsukada): 1, rotary encoder (Nikon RM-1000); 2, profilometer drive unit (Kosaka Laboratories SE3C); 3, target fixed to the stylus; 4, non-contact gap sensor (Kaman KD-2300-1S); 5, stepping motor (1.8° per pulse); 6, K coupling; 7, gear ($m = 0.5$, $z = 10$); 8, gear ($m = 0.5$, $z = 100$); 9, specimen; 10, diamond-tipped stylus; 11, adjusting screws; 12, moving table; 13, spring; 14, lead screw.



**Figure 4.127** Sampling synchronism requirement.

**Figure 4.128** Sampling patterns: (*a*) three points (digonal), $k = 2$; (*b*) four points (trigonal), $k = 3$; (*c*) five points (tetragonal), $k = 4$; (*d*) seven points (hexagonal) $k = 6$.

These will be considered in turn using as a basis the tetragonal case ($k = 4$).

(*a*) *Tracking length*

It is convenient to use $h/2$ as a basic unit, where $h$ is the distance of an ordinate from the centre of symmetry as shown in figure 4.128.

Suppose the area to be tracked is square (although this is not necessary) of dimension $mh/2$. The number of $x$ axis tracks is given by the integer part of

$$\left(\frac{2m}{3} + 1\right) \quad \text{trigonal} \tag{4.143}$$

$$\left(\frac{m}{2} + 1\right) \quad \text{tetragonal} \tag{4.144}$$

$$(m + 1) \quad \text{hexagonal.} \tag{4.145}$$

The ratio of the trigonal track length to the tetragonal track length for the same area is 4/3 or 1.33 and for the hexagonal track length to the tetragonal track length is 2.0, from which it can be seen that both the $k = 3$ and $k = 6$ cases are worse than the standard technique of $k = 4$ (the tetragonal or square grid) in terms of track length per unit area. Changing the aspect ratio of the two sides making up a rectangle of the same area makes no difference to the ratios given above.

The issue of whether or not the stylus is lifted off at the end of each $x$ axis track is insignificant.

(*b*) *Density of ordinates*

This is taken to be the total horizontal (or $x$ axis) track length for a given area divided by the $x$ axis spacings between ordinates (not $h$, except for when $k = 4$). For the cases considered this can be shown from figure 4.128 to be

$$\text{trigonal} \quad \frac{(3m/2 + 1)mh/2}{\sqrt{3}h} \tag{4.146}$$

$$\text{tetragonal} \quad \frac{(m/2 + 1)mh/2}{h} \tag{4.147}$$

$$\text{hexagonal} \quad \frac{(m + 1)mh/2}{\sqrt{3}h}. \tag{4.148}$$

Comparing with the tetragonal case as for track length, then the number of ordinates (trigonal case) for a given area divided by the number of ordinates (tetragonal case) for the same area is $4/3\sqrt{3}$ or 0.76, and for the hexagonal case divided by the tetragonal case is $2/\sqrt{3} = 1.15$.

From these ratios it is clear that, from the points of view of ordinate density and tracking length, the simplest case $k = 3$ is best and the hexagonal worst. But in terms of the comprehensiveness of the numerical model the hexagonal model is obviously best, so each of these three sampling plans has advantages.

There is an added complication in the trigonal case which is simply concerned with the visual impact of any map made of the surface. Unequal track spacing does not produce such a convincing map of the surface as does equal track spacing. Whereas this is not theoretically serious it can be disconcerting. Making the trigonal model spacing equal, which can be achieved as shown in figure 4.128(*b*), has the disadvantage that the ordinate spacing is now unequal, which in turn means more difficult interpolation to produce a continuous graph. Although this change in orientation does not affect the density of ordinates and the tracking ratio it does suffer another disadvantage, which is that each track length is unequal. This makes the instrumentation slightly more difficult.

The reason for the review in this section is to try and see whether taking measurements using non-conventional sampling schemes could produce any advantages to outweigh the disadvantage of complexity. The advantages considered are less information to collect, easier analytical derivation of theoretical results and simpler numerical methods.

The sampling schemes that have been considered all have the property that the information could be collected by sampling along parallel straight lines with a fixed sampling interval. (It might be necessary, however, to have a variable starting point, though this would follow a regular pattern.) This ensured that if a measurement (ordinate) was chosen when using a particular scheme it would always have the same number of adjacent ordinates at a distance $h$ (the chosen sampling interval), provided the chosen ordinate is not on the boundary.

From the point of view of simplicity of sampling mechanism the square grid ($k = 4$) in the tetragonal case is the best. In this case the spacing between the lines is constant and equal to the sampling interval $h$ along the line. Also the starting points for the sampling all lie along a straight line. However, the other schemes do have advantages to offset their complexity.

Both the trigonal ($k = 3$) and hexagonal ($k = 6$) cases have the advantage that measurements of the slope can be taken in three directions, as opposed to two for the tetragonal ($k = 4$) case. Although the theoretical results have been restricted to consideration of isotropic surfaces, it may still be of practical value to be able to check the assumption of isotropicity in more than two directions.

The trigonal ($k = 3$) case can be obtained by an alternative sampling method but this involves alternating the sampling interval from $h$ to $2h$. This alternative method is equivalent to rotating the grid through $\pi/6$. This rotation in the case of the hexagonal ($k = 6$) case just reverses the sampling interval along the line and the distance between the parallel lines.

From the point of view of collecting digital information the trigonal ($k = 3$) case is preferable as 'less' information is collected. The density of ordinates is $(4/3\sqrt{3})/h^2$ ($= 0.77/h^2$) compared with $1/h^2$ for the square grid ($k = 4$), so in the same area 23% fewer ordinates would be needed. The advantage of this would need to be weighed against the disadvantages.

Another advantage of the trigonal ($k = 3$) case is that fewer ordinates are needed when defining the properties of the extremities. To check the definition of a four-point summit only three conditions have to be obeyed as opposed to four conditions for a five-point summit. It should also be noted that some properties of the discretely defined random variables, such as the limiting value of ($k = 1$)-point summit (or peak) height as the sampling interval tends to infinity, are simply a function of the numerical definition and are independent of the surface being measured.

In summary, any discrete measurement of a surface must lose information compared with a complete 'map' of the surface. This is inevitable! However, any discrete measurement should produce results which converge to the results for the continuous surface as the sampling interval $h$ tends to zero.

For sampling along a straight line ($k = 2$) it is seen that the discrete results do converge to those for the continuous profile. They do not, however, converge to the results of the two-dimensional continuous surface. For example, $D^2_{\text{peak}} = 0.83 S_{\text{sum}}$, so that assuming independent measurements at right angles would produce a limit which was 17% too small.

For three-dimensional measurements when sampling with $k = 3$ or 4, the limiting results for expected summit density and expected summit height do not converge to the continuous surface. In the case of expected summit density the limit is 73% too large for $k = 3$ and 31% too large for $k = 4$. Again for expected summit height the case $k = 3$ is worse than for $k = 4$ but the differences are not so large. This suggests that some surface parameters may be estimated by discrete methods fairly well but others may not. For the case of average profile slope all three sampling schemes agree (for $k = 2$, 3 and 4) but this is, of course, an essentially one-dimensional parameter.

In order to consider the merits of sampling schemes it is necessary to study their theoretical properties. By doing this it is possible to obtain new insights into the general problem. This is possible only by using mod-

els which lead to tractable mathematics. The three simpler sampling schemes $k = 2$, $k = 3$ and $k = 4$ considered theoretically earlier have been chosen because they have a common property which enables them to be investigated using analytical results previously obtained in theoretical statistics. Using the trigonal ($k = 3$) symmetry case leads to a simpler mathematical model than for the tetragonal ($k = 4$) symmetry case, as this reduces the dimension by one. However, taken as a whole it may be that the hexagonal sampling plan where $k = 6$ offers the maximum benefit in terms of the three criteria mentioned above. In order to verify this, the mathematical expressions for the surface parameters will have to be derived. One message which has emerged from this exercise is that the *conventional grid pattern method of sampling is not necessarily the best.*

(*c*) *Static sampling*

The other way to ensure that the *z* scale is at the same datum for each data point is to use very flat, smooth and accurate *x* and *y* slides. This method has been used by many researchers. The movements are usually stepper motors or servo drives. In the event that *x* and *y* tables are used often the stylus is not allowed to traverse. All that happens is that the specimen is moved underneath it and, at every grid point, the stylus is made to contact the work. In this way a complete grid can be built up with any desired pattern.

Obviously a number of precautions and aids have to be incorporated in any good system. For example, the table supporting the *xy* movements has to have tilt adjustments in two directions. Also the device should be incremented in one direction only, to reduce the effect of backlash. As a general rule of thumb the variability of each sample interval should not be greater than one-tenth of the interval itself, so for an interval of 1 $\mu$m the standard deviation should be less than 0.1 $\mu$m. Also cumulative error over, say, 1000 sampling intervals should be no bigger than half an interval. This sort of accuracy has been achieved for example by Baker and Singh (see [95]). Typical sizes of data stored using these methods are ~$10^6$ in 100 traces of 10 000 data points or thereabouts. Actually there is rarely any need to take as many data points as this. Freak behaviour can be spotted over data sets 10 times smaller than this.

There have been variations on this basic set-up, for example using an oscillating-stylus method [86].

### 4.2.5.5  *Contour and other maps of surfaces*

Having decided on the sampling pattern, taking account of the various considerations of coverage etc, there remains a number of problems. One is that of processing, that is what features of the surface are to be measured digitally based on those reviewed in the processing section? Another question which is not to be ignored is that of display. There are two possibilities: one is a simple contour map originally used by Williamson [89] and the other the isometric view shown originally by Sayles and Thomas [90].

Taking the contour method first, the simplest way of generating the contour line is to establish a given height plane and then find the intersections by linear interpolation between the points making up the grid. The coordinates thereby obtained can be either connected by straight lines to get the necessary contour or jointed by more elaborate curve-fitting routines such as the cubic spline function. There are, however, numerous alternative methods of smoothing the derivatives at the joins between meshes.

Obviously certain criteria have to be followed such as that the contour lines never intersect each other, never form branched lines and never break—all obvious but necessary to state. More detail on this is obtainable from pattern recognition texts.

Following Sayles and Thomas' philosophy of presentation it is well to realize that most contouring programmes are not written to display predominantly random data, with the result that discontinuities are sometimes created at data point boundaries and the contours generated can be misleading. Another point often mentioned is that contour programs usually generate contours in a piecewise fashion. This has the effect of neglecting the specific identity of each contour and, by implication, any and each asperity.

The data cannot be assumed to follow any preconceived trends and generally there is no justification for reducing the number of data points required in an area, or assuming anything about the spatial distribution. In many ways the construction of such diagrams for engineering surfaces is more difficult than the equivalent

larger-scale problem of producing the Ordnance Survey. Here the surveyor may choose the mesh and separation of data regions to suit the severity of local conditions, but on the scale of sizes found on machined surfaces this is rarely possible.

Another possibility is the isometric view.

### (a) Isometric representation

Although random data is best represented in contour form, there are many cases where an isometric diagram can be useful to identify salient features. An engineering surface often contains machining or forming marks which, although similar in size, are usually distributed randomly in height on the surface. Clearly a contour representation is not ideally suited to show such trends, but in isometric form the presence of similarity between features becomes clearer.

An important asset of such a plot is its advantages over SEM images. It can supply the engineer with a magnified surface picture of chosen scale, and, even more important in many cases, of chosen bandwidth. It is well known that the small asperities, which generally possess the higher surface slopes, tend to dominate SEM scans, and longer spatial variations, which in many applications are the most important, are not well represented. An isometric plot, drawn with the emphasis on the significant bandwidth, is the only technique at present available which will portray the surface as relevant to its application.

Realism can be increased by the suppression of parts of profiles where they are screened by the features of previous profiles. This is obtained by maintaining and updating an array representing the maximum height, on an imaginary vertical grid, which any previously drawn profile has reached.

By comparison of each profile prior to plotting with this array the decision to lift or lower the pen carriage (i.e. blank off the write mechanism) can be made.

### 4.2.5.6 High speed area tracking stylus [96]

Instead of a conventional stylus pick-up to measure the surface it is possible to simplify the mechanics to increase the speed at which the stylus can track without jumping off the surface. Such a system is shown in figure 4.129, which should be compared with figure 4.48. In practice, figure 4.129 is just a simplified version of the true side-acting gauge.

The idea is to reduce the inertia term dramatically by removing the pivot and simply attaching the probe to the hinge holding the bearing. This simple trick reduces the referred mass term in equation 4.45. Also the electromagnetic transducer is replaced by an optical system.



**Figure 4.129** Schematic diagram of prototype high-speed stylus gauge.

**Figure 4.130** Maximum trackable amplitude as a function of frequency.

The tracking relationship i.e. the maximum velocity $V_T$ which can be achieved before the stylus jumps, here for a sine wave of amplitude $A$ and wavelength $\lambda$ is

$$\left(\frac{2\pi}{\lambda} V_T\right)^2 = \frac{F}{A}$$

(4.149)

where F is the force on the surface exerted by the hair spring via a shank of length $L$.

It can be seen using this simple approximate formula that increasing $F$ only has a square root increase in $V_T$.

Figure 4.130 shows that the tracking speed for this areal tracking stylus is almost an order of magnitude higher than for a conventional side-acting gauge (i.e. it can track 0.1 $\mu m$ amplitude wave of 1 $\mu m$ wavelength at 5 mm/sec rather than the typical 1 mm/sec used by other stylus methods and also optical followers).

The biggest problem with this type of system is damping. It cannot be controlled to the same extent as a conventional stylus, so for example maintaining a damping factor of 0.59 indicated earlier for wide band-width signals is unrealistic. The uncertain damping causes a spread in repeatability, which is why it has not been taken up extensively [96].

## 4.3   Optical techniques for the measurement of surfaces

### 4.3.1   General

Stylus methods, because of their great versatility, can encompass an enormous range of sizes of roughness and processes within their measurement capability, but there are drawbacks in the conventional measuring instruments used for roughness and form.

Some of these are:

(1) the technique is relatively slow;
(2) the stylus force can in some instances damage the surface;
(3) limitation on measuring areas.

The obvious advantages are:

(1) versatility to accommodate a wide diversity of shapes;

(2) high range to resolution in the vertical direction;

(3) high spatial bandwidth.

To a large extent the deficiencies of stylus methods might be made up by employing other methods. The obvious possibilities which will be discussed are:

(1) optical methods;

(2) capacitative methods;

(3) other techniques such as eddy current, pneumatic, etc.

There are many possibilities available and nowhere can it be said that any method is no good. It is a fact that for a given process, material and range of size, most techniques can be developed to give some sort of useful output. However, all techniques have a comfortable working range. To go outside this range requires some effort.

Having said that, the most obvious complement to the tactile instrument would appear to be optical instruments [97, 98].

*Comparison between stylus and optical methods*

A source of great irritation in engineering metrology is the comparison of surface finish values obtained using a stylus method with that obtained using an optical technique. Because the stylus technique has only one basic mode it will be compared with the optical technique devised to imitate it, which is the optical follower in which a focused spot on the surface is the optical equivalent of the stylus. Followers will be discussed later.

One of the reasons for discrepancies between stylus and optical methods is because they work on different physical principles. Figure 4.131 shows an example in which two films of the same thickness $d$ but different refractive index $n$ are placed together on a substrate.

Figure 4.131(*a*) shows what a stylus measures, whereas figure 4.131(*b*) shows the optical follower equivalent.

The two outputs show that although this dimensional thickness of the two films are the same (i.e. $d$) they have different optical path thickness of $nd$. Consequently, the optical method shows a step at the film joint whereas the stylus method does not! The two methods should give different results! In the case shown, there



**Figure 4.131** Optical and stylus measurements.

is no change in the stylus method in moving from one film to the other but in the optical follower method the difference is

$$2(n_2 - n_1)d \tag{4.150}$$

Also, as shown in figure 4.132 any optical method suffers from optical diffraction effects at a sharp edge on the surface. This is good for defining the boundary edge but poor for measuring the height of the step.

For a general surface both the stylus method and the optical method impinge on the surface. In the case of the stylus there is a mechanical deflection dependent on the elastic modulus E. Light also penetrates the surface. The degree of penetration is determined by the conductivity of the material, so the penetration is different for the stylus and the optical methods.

Table 4.6 shows a simple comparison. Properties useful for surface finish instruments are listed. Where there is an advantage a tick is place in the appropriate stylus or optical box. It is clear from the table that there is about the same number of ticks in each column. There is no outright winner or loser: both have their strengths.

A general metrology rule is that the method of measurement should wherever possible mimic the application, so in applications involving contact, for example gears or bearings, then the stylus method is preferred. However, in non-contact applications such as mirrors the optical method should be used. When lateral structure is of interest, such as in lithography, optical and similar methods can be used to advantage because of the edge enhancement.

More of the optical properties will be discussed in the following section.

The reason is obvious. Most optical methods derived from real life are concerned with areal visualization in terms of edge recognition, colour and intensity variations. The eye has the past experience to judge



**Figure 4.132** Edge enhancement.

**Table 4.6**

| Stylus | Optical |
|---|---|
| Possible damage | No damage √ |
| Measures geometry √ | Measures optical path |
| Tip dimension and angle independent √ | Tip resolution and angle dependent |
| Stylus can break | Probe cannot be broken √ |
| Insensitive to tilt of workpiece √ | Limited tilt only allowed |
| Relatively slow speed | Can be very fast scan √ |
| Removes unwanted debris and coolant √ | Measures everything good and bad |
| Can be used to measure physical parameters as well as geometry for example hardness and friction √ | Only optical path |
| Roughness calibration accepted at all scales √ | Difficult to calibrate by standards |
| Temporal and spatial influence/dynamic | Spatial influence/geometric effects |

quality for ordinary vision but not for surfaces. However, it seems natural to regard the areal discrimination of conventional optics as the complement to the axial resolution of the stylus. The question arises for both: can they take over each other's role as well as maintaining their own? The previous section dealt, to some extent, with using stylus methods for areal assessment. Now it is the turn of optics. Can it provide the vertical capability of the stylus method and capitalize on its non-contact and fast response potential?

As a basis for comparison, consider figure 4.132. The graph has two abscissae, one related to incoherent optical methods and the other to coherent. The ordinate is a number of different things. In some cases it is better to consider it as geometrical fidelity or the $z$ movement capability of the instrument; in others it is best to consider it as the area covered by the basic instrument.



**Figure 4.133** Optical systems.

What are the sort of factors that have to be considered when light hits a surface? Because of the variety of different optical techniques that are possible it will be necessary to introduce some of the relevant theory before each section, so for optical probes comments about the geometrical optics will be introduced. For interferometry spatial and temporal coherence will be discussed and for scatter, diffraction and speckle some wave theory.

At this stage it must be emphasized that the relevance of the particular discipline to texture and form only will be discussed and not its general applicability.

In its simplest form consider figure 4.134. On a macroscopic level one might expect a surface to behave as a simple mirror which reflects about the normal to the general plane but it is rarely as simple as this. Local slope $\gamma$ of the surface causes the reflected ray B in figure 4.134 to emerge at an unexpected angle $(\alpha - 2\gamma)$.

If the slope at O is $f'(x)$ OB will be reflected to

$$\alpha - 2f'(x).$$

(4.151)

In other words the angle of basic reflection is determined by local slope.

In the case of figure 4.134 where curvature is taken into account, the beam O'D is reflected at an angle to the normal of $\alpha + \beta - 2f'(x - s)$.

Using Taylor's theorem the slope at $x$ - $a$ from $f'(x)$ is given by

$$2sf''(x) - \beta$$

(4.152)

so, simply, the scatter around the basic reflection is related to local curvature $f''(x)$ for $f'(x)$ small, $s$, the spot size, $\beta$, the degree of collimation of the incident beam and also the characteristics of the source and the depth of focus of the system (related to $\beta$).

When the statistics of the surface are taken into account the situation becomes complicated. Reverting back to figure 4.134, the abscissa does not just represent spot size, it more nearly represents the number of independent local variations of the surface within the area of illumination of the spot. Near to the origin in the figure the area of illumination is small compared with $\tau_L$, the independence length of the surface, say $10~\mu m$ for grinding. For a periodic surface when the period is large compared with $\lambda$ equations (4.151) and (4.152) are relevant.



**Figure 4.134** General reflection.

### 4.3.2 Properties of the focused spot

The fundamental problem associated with all techniques attempting to simulate the mechanical stylus is the fact that the focused spot size and the depth of resolution are not independent.

The simple Rayleigh criterion for resolution insists that the diffraction-limited spot size—corresponding to the stylus tip—has a value

$$s_d = \frac{1.22\lambda}{\mu \sin \alpha}$$

(4.153)

where $s_d$ is the spot size, $\lambda$ is the wavelength of light, $\mu$ is the refractive index between the medium and the object and $\sin \alpha$ is the effective numerical aperture of the lens—which corresponds to the stylus slope.

The depth of focus $D_f$ defined as the change in focal position for an increase in beam diameter of, say, $\sqrt{2}$ is given by

$$D_f = \frac{1.22\lambda}{\mu \sin \alpha} \bigg/ \tan \alpha.$$

(4.154)

Thus, one inherent difference between stylus and optical methods is that the range is obviously independent of tip size and slope for a stylus, whereas for an optical probe the range to resolution ratio depends on the numerical aperture.

It could be argued that the depth of focus is, in effect, no different to the elastic deformation of the stylus, which is larger the bigger the stylus tip radius, and that the transducer system itself provides the range limitations. This is true to a large extent but the fact remains that the mechanical force to close the loop between the stylus and workpiece is much more positive than the optical one. It is no doubt for this reason that optical systems such as Nomarsky and differential techniques used in video systems are often used. Another is that optical 'tips' cannot be damaged!

Typical values of $\lambda = 0.85$, $\mu = 1$, $\sin \alpha = 0.28$ give the spot size as $3.7$ $\mu$m and the depth of focus as about 13 $\mu$m, which is not large.

For laser light the result is slightly different because coherent light does not behave in the same way as incoherent light (figure 4.135).

The effective formula relating focus and depth of focus in this case is determined from

$$\omega_z = \omega_0 \left[ 1 + \left( \frac{z}{\omega_0^2 \pi} \right)^2 \right]^{1/2}$$

(4.155)

from which for the same $\sqrt{2}$ criterion the depth of focus works out to be

$$D_f = 2\pi \frac{\omega_0^2}{\lambda}.$$

(4.156)



**Figure 4.135** Laser collimation and depth of focus.

The effective far-field half angle for divergence is given by

$$\theta = \frac{\lambda}{\pi \omega_0}.$$

(4.157)

From this the spot size determines the effective numerical aperture (or vice versa, because a small spot can only be obtained if the beam is expanded first) and also the depth of focus $D_f$.

Notice that the beam is not uniform but Gaussian and that o»o has to be measured accordingly.

In figure 4.136 the intensity pattern is

$$I(x, y) = I_0 \exp[-2(x^2 + y^2)/\omega^2].$$

(4.158)



**Figure 4.136** Laser beam intensity distribution: $I = I_0 \exp[-2(x^2 + y^2)/\omega^2]$

The question has to be posed as to whether or not one gets a longer depth of focus for a laser for the same spot size. Taking the case above, if $2\omega_0 = 3.7$ then $\omega_0 = 1.85$, which gives for the same $\lambda$ a depth of focus of 25 $\mu$m, or twice as good? Or is it just the woolly definition of spot size?

Some devices are being introduced which incorporate two beams, one of which has a high NA to get the resolution and the other a low NA and a big area of illumination to provide a measure for increasing the range, by using in effect an optical skid (see section 4.4.8).

### 4.3.3 Optical followers

If the intention is to mimic the stylus instrument then it is necessary for some feature of the optical reflection to convey the same information to the detector as the stylus. Then, as the surface moves horizontally changes in the pattern on the detector should initiate a movement of the instrument vertically to maintain the original condition. These instruments are based on the principle that, as the surface moves laterally relative to the probe, the probe will be forced to 'follow' the surface geometry. As a result these instruments are sometimes called 'followers'. Obviously they always include a scanning mechanism.

A large number of optical instruments fall into this category. They have an advantage in that, providing the sensor responds to a change in geometry of the profile, the nature of the change does not have to be fully understood because the system works on a closed-loop null principle. In many other optical and related methods extracting surface geometry information from the detected signal requires the making of assumptions about the statistics of the surface. They are basically indirect methods. The more direct methods—the 'followers'—will be considered first.

Many reviews have been made of the various optical techniques [42, 43].

Most optical probe techniques are based upon some examination of the wavefront at the focus: sometimes there is one point of focus, sometimes two axial focus points, produced by artificial means. In the former case there is some object situated near the back focus position which modulates the focus transversely as the surface moves up or down (e.g. on a hill or in a valley). In other cases the transverse movement is

caused by other means such as having the optical system non-normal to the surface, and yet others involving Nomarsky are normal incidence and measure slopes.

Consider focus methods first. An ingenious method has been devised by Dupuy [99] based on the principle of the Foucault knife-edge test. It is shown schematically in figure 4.137.

In this method a very small spot of diameter of approximately 1 $\mu$m is focused from a bright (not necessarily laser) source onto the surface by a projective lens system. A numerical aperture of between 0.4 and 0.8 is incorporated into the projection optics design so that there is sufficient light on reflection to fill the pupil in the image space—probably of a reduced numerical aperture.

A polarizing beam splitter and retardation plate close to the objective transmits this light to form the image at the knife edge, the spot O′ on the surface effectively acting as a secondary source of O.

The in-focus condition is detected with high resolution by differential photodiodes located behind the knife edge. A small lens in the projector system is axially oscillated causing the projected spot to scan vertically through the surface with useful longitudinal magnification. The output may be obtained by controlling the lens vibrator in a servo system mode or alternatively by deriving an analogue output from the time coding of the oscillating waveform. In another mode the knife edge can be moved to maintain the null, but a positional axial shift of $\Delta x$ of the surface will, due to Newton's axial magnification formula, result in $(\Delta x)^2$ in the movement of the effective focus of the spot from the edge. So, although the frequency response of the system can be improved, the amplitude of the movement is larger. For example, if a normal microscope objective is used with a numerical aperture of 0.95 and matched to its design conjugates the magnification will be of the order of 30$\times$. Consequently the motion required of the knife edge would be 900$\times$ that of the movement of the surface. If a range of 25 mm is expected for the instrument the knife edge would have to move by nearly 25 mm, which is hardly practical.

Another technique is to attempt to interpret the signal in the detector plane more closely still using a knife-edge obstruction. This involves the use of storage and calculation. There are obviously numerous ways in which this could be achieved.

Keeping the surface stationary has been the norm in the traditional way of testing lenses. Interpretation of the resulting image has been subjective but very sensitive. The figuring of high-quality telescope mirrors has been achieved to a high degree of perfection. One method [100] uses Zernicke polynomials to develop coefficients from the knife-edge position. They are then converted to wavefronts by using matrix operations. The use of Zernicke polynomials is advantageous in the sense that the different coefficients can be related easily to first-order and Siedel aberrations. Principally the method makes use of the Lagrangian differential.

If the transverse ray aberrations are $X(l, m)$ and $Y(l, m)$ they must be determined accurately if the total integrated wavefront is to be found. The problem is that the shadowing of the source (pupil) caused by the knife edge does not occur as a sharp edge as the knife edge cuts through the imaged point. This is shown simply in schematic form in figure 4.138.



**Figure 4.137** Dupuy system optical follower.

**Figure 4.138** Detector signals for Dupuy system (*a*) before focus, (*b*) in focus, (*c*) after focus.

As pointed out by Granger [100], the entire pupil shadows uniformly (figure 4.138), but does not go dark abruptly; instead the shadowing is gradual due to diffraction effects at the edge. The diffraction effect is worse for small spot sources used by optical probes. Another way of finding the edge is to plot the intensity at a point $(l, m)$ in the pupil and record the intensity as a function of the knife position in $x$ and $y$. The 50% point of the cumulative intensity curve is used as the assumed location of the ray aberration $X(l, m)$ and $Y(l, m)$.

Thus the process could be to start with scanning the image at the knife edge. The intensity of light at a grid of sample points is recorded for each position of the knife edge (equivalent to each position of the surface). This is done for a number of steps of the surface, so the intensity pattern at the knife edge modulated by the position of the surface can be stored. This large array of intensity data is usually processed to the smaller set of transverse aberrations if mirrors are being tested. However, for a simple point portion of the surface it need not be necessary. In effect the position of the surface detail axially is being coded as an intensity pattern via the knife edge step by step. By interrogating the intensity pattern, therefore, the height detail can be obtained. Alternative methods can be used in this way. It is not necessary to have a knife edge in the exit pupil focus; in principle any obstruction could be used, such as a Hadamard mask [101]. Here the calculated waveform variations are complicated but the coding is more convincing.

In its simplest form the system is prone to veiling glare, and, because of the spacing of the detectors and the finite size of the spot on the surface, local tilt of the surface is important. As usual there always has to be a compromise in the basic simple system. If the two detectors are brought closer together the sensitivity to tilt reduces but the longitudinal sensitivity also reduces.

Another variant uses an optical wedge to form a lateral optical intensity pattern. This is shown in figure 4.139. A similar optical configuration is used except that the reflected light rays pass back through the beam splitter and are brought to focus by lens A. A wedge bisects the pencil of rays forming two separate images of the spot on the surface, each image falling on the very small gap separating a pair of diodes. Displacement of the surface being examined causes the images on the photodiodes to be displaced longitudinally along their optical axes, but the asymmetry introduced by the wedge causes a lateral change in light distribution across the photodiodes. Both magnitude and direction of the surface displacement are available by means of analysing the photodiode signals. The resultant error signal is fed to the objective lens displacer, therefore bringing the spot back to focus on the surface. In this application the photodiodes could produce a position-sensitive signal in either analogue or digital form.

Two longitudinal focus investigations are possible, for example that due to Simon [102]. In this method the basic system is the same (figure 4.140) but, whereas transverse optical obstructions have been used in the other methods, this uses longitudinal obstructions. What happens is that the light rays reflected from the surface are brought to focus by lens C. The pencil of rays from C is divided into two by means of a half-silvered mirror. The Nightingale technique concerns the analysis of the energy density on either side of the nominal

**Figure 4.139** Optical follower using optical wedge.

forms of C. Accordingly two stops, A and B, are fitted, one before and one after the nominal focal plane of C, on their respective axes.

Displacement of the surface in the normal direction causes the images of C to move axially. A change in light intensity detected by the photodiode behind stop A will be accompanied by an opposite change in that behind B. The error signal contains amplitude and direction information. This method is less sensitive to glare than other methods. In fact, the sensitivity has been shown to be given by $KI_0 M(\Omega)$, where $I_0$ is the ini-



**Figure 4.140** Simon method of optical follower.

tial intensity, $M$ the magnification, $\Omega$ the NA of the objective and $K$ a constant of proportionality—a simple system sensitivity.

Alternative methods exist which, although they pick up spot detail, are not strictly probes. They involve some degree of analysis of local fringe patterns yet still rely on some axial movement to provide the necessary phase information. One obvious one is based on Nomarsky [103] (figure 4.141).

This is the principle of differential interference microscopy. See later for detailed discussion of interference. Light from S is reflected from beam splitter B and is incident on a Wollaston prism W. This produces two angularly sheared beams of opposite polarization. These are focused by O into the surface. The separation on the surface of these two beams is about the size of a stylus which can be taken to be about 1-2 $\mu m$. These beams are reflected back to recombine at the Wollaston prism and then the orthogonally polarized beams are made to interfere by means of the linear polarizer P. The interference pattern is viewed either by eye or by a camera. The element C is a polarization compensator used to adjust the phase delay between the two beams to alter the contrast of the interference pattern.

Because there are two spots effectively produced on the surface and the interference is a measure of their path difference, this device is, in effect, a slope-measuring system rather than one of amplitude. When properly adjusted, positive slopes appear bright, negative slopes appear dark and zero slope uniformly grey. Note that slopes perpendicular to the direction of shear of the Wollaston prism are not seen.

The slope measurement can be made absolute using the conventional pentaprism-laser configuration shown with the Nomarsky probe in figure 4.142.

As before, if the system is adjusted properly the intensity at the detectors is linearly dependent on slope—provided that the height variations are relatively small when compared with the wavelength of light [104, 105].

One useful point to note is that if white light is used rather than a laser, the contrast of the image is seen as a colour variation which can be exploited if the eye is the detector.

Other examples of cases where colour can be advantageous to get an indication of the surface is in the use of FECO (fringes of equal chromatic order). Note in figure 4.143 that, as the height differential on the



**Figure 4.141** Optical follower using Wollaston prism.

**Figure 4.142** Pentaprism with Nomarsky system.

surface increases, one detector increases in signal whereas the other decreases. Summing the two detector signals allows compensation for the possible light intensity variations of the laser.

Taking the difference of the two signals gives $2\Delta I$ which is linearly dependent on $\Delta h$, the surface slope. Hence $2\Delta I/2I$ is also linearly dependent on $\Delta h$ and independent of intensity variations. Similar tricks have to be employed in all optical methods.

A slope profile can be obtained by translating the whole head in the direction of the laser beam, which acts as a vertical reference in much the same way as for a stylus instrument. Intrinsically, this technique is very sensitive to seeing small detail because it is in effect an optical differentiator. To get a surface profile the signal has to be integrated.

Typical height sensitivity is about 1 nm with a vertical resolution of about the same as a stylus. Unfortunately the method has a limited vertical range much less than for a stylus method.

Other methods exist, which rely on polarization (e.g. figure 4.144). The system shown in figure 4.144 is a variant of the previous system but it introduces another factor. In this the two polarized beams produced by the calcite/silica doublet are focused at different axial positions normal to the surface rather than alongside as



**Figure 4.143** Signals on detectors.

**Figure 4.144** Optoelectronic interference contrast optical follower.

in the Nomarsky method. This produces a completely different type of stylus in which one polarization acts as the sharp stylus spot, whereas the other, the defocused one, acts in a way as a skid from which to measure the other—again a very clever mimic of the stylus technique.

In this system the polarizer P is rotated until there are beams of equal intensity in the two orthogonal polarization directions P and S. The half-wave plate is adjusted to give a maximum intensity of the light transmitted by the polarizer.

The fact that the doublet is calcite/silica means that it has different powers for P and S paths. Thus, employing a conventional objective lens O, one beam is focused to a spot of nominally of 1 $\mu$m size while the other covers approximately 100 $\mu$m, depending obviously on the particular doublet used. The light reflected from both beams from the surface is combined by the analyser which gives the result that phase shifts between the two beams are converted to intensity and subsequently to electrical signals.

As with other methods involving polarization this technique suffers to some extent from the fact that the chemical properties of the surface skin rather than its geometry can influence the magnitude of the reflected polarized ray.

The interesting feature of the method outlined here is the use of an optical skid. Some other attempts at this have been made with scatter methods [97], which will be considered later.

The reason why there is such an interest in measuring fine surfaces is principally due to the pressures of the semiconductor and defence industries. Whilst scanning electron microscopes and stylus methods are capable of high resolution, they both suffer from operational problems. In the latter there is a contact problem whereas in the former the measurement takes place in a vacuum.

The obvious method is the use of optical microscopes but they are inadequate for linewidths of below 1.5 $\mu$m, particularly as the width, registration and distortion tolerances are shrinking faster than the linewidths themselves.

### 4.3.4 *Hybrid microscopes*

One possible way forward is the scanning confocal microscope (SCM) [106]. In this there is increased resolution which results from the confocal properties of the system and from the use of rather smaller wavelengths.

This has the advantage of being non-contacting. The apparent enhancement of edge resolutions is 3:1 or thereabouts [106–108]. The reason why is explained shortly, but first we must ask the question, why scan?

The reason for a scanning system in microscopes might not be obvious. The whole method of obtaining resolution by lenses involves so much loss of contrast, lack of magnification and so many other difficulties that it is often hard to get a good display. Some of these difficulties can be avoided by using a different means of getting resolution. The basic idea of resolution is to separate in some way the light passing through, or coming from, very close regions of an object. The conventional microscope does this by using refraction by lenses to separate the light from neighbouring regions. An alternative method is to use the lens system the other way round, that is to produce a spot of light. Discrimination between adjacent points is then produced by projecting the light at them at different times by making the spot scan [109]. After reflection (or transmission) the spot is made to fall onto a photocell with subsequent amplification and display. Basically the scan method (sometimes called the flying-spot microscope) converts the intensity distribution of light from microscopic objects into a voltage time series.

The idea behind the flying-spot microscope has been used many times, the version mentioned here being just one. In another method the spot source can be that of a spot on a cathode-ray tube (CRT). As this spot is imaged onto the surface it is automatically scanned using the coils in the CRT over quite a large area. The whole area is then imaged onto a large photodetector insensitive to position. The voltage from this can be automatically synchronized to the voltage deflecting the beam in the CRT.

Whereas usual projection microscopes have a maximum magnification of about 2000 x, the scanner can give 500 000 x and it is not limited to low- or medium-power objectives. Most use has been with transparent parts. Quantitative information has not yet been obtained for engineering specimens as with all laser scanners mentioned earlier.

Both the scan mode and the fine spot are integral issues in the design of a confocal system, which in principle is a refinement of the flying-spot method developed by J Z Young 30 years earlier [109].

The method relies on the depth discrimination properties of the so-called confocal microscope. This simply means that the imaging and receiving optics are exactly the same, which gives very definite advantages in optical properties. As in the flying-spot method the optics is used to project a point onto the surface. The reflected light is picked up by a point detector. The resulting signal is used to modulate the brightness of the spot on a television screen which is scanned in synchronization with the object. The essential issue is that the first lens focuses a diffraction-limited spot onto the object. The same lens collects light only from the same portion of the object that the first lens illuminates. It is a consequence of using the single-point detector that both the image-forming and -receiving systems contribute to the signal at the detector (figure 4.145).

The depth discrimination as explained in references [106, 108] is best understood by means of considering a transmission object (figure 4.146).

Consider the light from outside the focal plane of the collector lens. This forms a defocused spot at the detector plane. As a central point detector is used this results in a much weaker signal and so provides discrimination against detail outside the focal plane. The integrated intensity, which is a measure of the total power in the image of figure 4.146 is found to fall off monotonically as the object is defocused [108]. This



**Figure 4.145** Confocal imaging system.

**Figure 4.146** Developed confocal system.

fall-off is dramatic. It has been reported as reaching the half-power point at $0.7\lambda$ from the focal plane (NA = 1). Therefore, when the object passes through focus the image intensity reaches a sharp maximum. Moving the object axially until a maximum is reached can be used to indicate surface height—an optical probe effect in fact, only mixed in with a scanner. In this way this type of device can be used as an alternative to a microscope and has resolutions of 1 $\mu$m in the $x$ and $y$ directions and 0.1 $\mu$m in depth—not quite electron microscope standards but more practical in many respects.

The key to the working is the refocusing of the spot on the surface onto the detector. In principle the spot can be made smaller than the projected spot by means of a pinhole at the detector. No diffuse light is collected at all, only the immediate specular component. This means that the confocal scanning microscope becomes close to being the first instrument having resolutions that are good in both the axial $z$ direction and the $x$, $y$ directions. The position of the maximum intensity can be found very effectively by differentiating the intensity signal. Obviously the output, that is the roughness, is simply the amount by which the object has to be moved axially to obtain maximum intensity on the detector.

Because of the restricted focus of the confocal system it is sometimes beneficial to increment in height as well as scanning horizontally, as shown in Figure 4.147.

Take, for example, scan 3 in the figure. Only the points A and B would register because it is only these points which would be in focus. In scan 1 points C, D, E, F and G would be picked up. If one level is chosen and the scan taken over an area then a contour is generated at that level [110]. The whole process is called 'laser scanning confocal microscopy' (LSCM).

The probability density curve of the profile and/or the areal view of the surface can be obtained. The reflectance image of the surface can be obtained using LSCM.

The stack of contours can be used to determine conventional surface parameters such as $R_a$ and $R_q$ and also skew and kurtosis. These can be compared with the results obtained using a stylus instrument. Typical results indicate that the stylus readings are higher than LSCM by between 5% and 10%.

The idea of 'stacking' contours used in this laser scan confocal scheme is very similar to the method used to develop rapid prototypes in design. Material ratio curves can easily be craters and the amount of material moved in plastic deformation.

Whilst on the subject of fine surface measurement, one technique which has only been hinted at so far is a kind of microscopy and a kind of interferometry [111, 112], and this is the method called FECO (fringes of equal chromatic order). The basic interference is illustrated in figure 4.148 and relies upon multiple beam interferometry [113]. The plates are coated with reflecting films and mounted parallel and close together. Instead of the monochromatic light normally used, white light is used.



**Figure 4.147** Incrementing in height.

**Figure 4.148** Basic multiple beam interference.

The basic interference equation for the fringes is

$$(m+1)\lambda = 2t + (\rho_1 + \rho'_1)\lambda/2\pi \tag{4.159}$$

where $m$ is an integer and A is the wavelength at the minimum of the reflected intensity distribution. For silver, $\rho_1$ and $\rho'_1 0 1.2\pi$. Thus the lowest order which is obscurable is for $m = 1$ (the dispersion of wave change can be neglected).

Results using this method indicate that measuring $R_t$ involves determining the peak-to-peak roughness from the extreme width of one fringe and the order number of the fringe [114]. This may be only a rough approximation and, considering the surface as a $2\sigma$ spread, it could be better.

Surfaces having roughnesses of about ± 1nm and slopes of ± $2 \times 10^{-4}$ rad have been evaluated but the technique does require interpretation and automation. Whether this is now necessary remains to be seen.

The typical apparatus is shown in figure 4.149.



**Figure 4.149** Simple schematic diagram of FECO fringes.

### 4.3.5 Oblique angle methods

The methods outlined above have first and foremost used normal incidence spots. However, many well-known methods have used oblique angles to get advantages in one form. Perhaps the most famous method is that due to Schmalz [115] who illuminated the surface past a knife edge (figure 4.150).

The only way to get a magnified view of the surface was to project a slit or edge onto the surface at an angle $\alpha$ from the normal and view it by eye at a different angle, so if the eye is at normal incidence the magnification is equal to sec $\alpha$.

**Figure 4.150** Oblique light edge profile.

Since this time there have been many other devices and techniques which, in some form or another, use the non-normal technique, although they do not necessarily image a thin spot or line onto the surface as above. One which does is shown in figure 4.151 [116].

The illumination shown in figure 4.151 is that of a laser. It is focused on the surface at an angle of 45° (the spot size is about a few micrometres) and as the surface moves transversely the spot moves about as seen by the microscope, which in turn relays the movement to a TV camera or similar detector. Thus, knowing the angle, the height can be measured assuming that the vertical movement is small—which is not often true. Other oblique methods will be discussed in the section on flaw measurement.



**Figure 4.151** Oblique section topography.

A more accurate device uses two oblique incident rays and gets rid of the problem of surface tilt to some extent. This is shown in figure 4.152 [117].

An image of approximately 1–2 $\mu$m is projected onto the test surface via two separate paths divided off by a beam splitter. Two mirrors deflect the rays of the respective paths into oblique pencils of equal and opposite angles of incidence. The reflected light returns via the opposite path to the beam splitter and then to a lateral



**Figure 4.152** Bottomley technique.

sensor. Variations of surface height are detected as lateral changes of the sensor (which could have been achieved with only one ray!). However, if the surface is tilted a double image is formed but the lateral sensor gives an unchanged reading due to the symmetry of the optical arrangement.

Methods involving oblique incidence and small spots have been used mainly for position sensors rather than surface texture or form measures. Using the principle of triangulation but a mixture of both normal incidence and off normal reflectance can still be useful for both position sensors and texture measuring as shown in figure 4.153 [118]. This type of technique will be discussed more with regard to how the texture can be evaluated from the scatter.



**Figure 4.153**

The advantage is that generally some extra magnification of the feature to be measured is possible due to the obliquity; the disadvantages are that often height and spacing information gets mixed up—the horizontal movement of the part can affect the readings, and it is not easy to economize on the optics, whereas for normal incidence it can often be possible to use the same optics for the incident and reflected beams.

### 4.3.6   Phase detection systems

Many instruments use phase detection as the essential part of the measuring system, in particular interference methods, both homodyne methods in which one wavelength is used and heterodyne techniques when two variables (which can be wavelengths) are used. Before looking at one or two examples of this, the term coherence will be explained because the concept is important in interferometry and speckle.

#### 4.3.6.1   Spatial and temporal coherence

These are two different properties of a wavefront. In conventional interference one wavefront is split into two by partial reflection. Later these are recombined. What determines how well they combine is the 'temporal' coherence of the light. This is equivalent to the inverse width of the spectral bandwidth of the light.

Thus if $B$ is the bandwidth of the signal then $1/B$ is the coherence length of the light—the same as the correlation length concept in surface finish described earlier. If the path difference between two beams in an interferometer is greater than $1/B$ then they will not add constructively or destructively—the phase information is effectively lost (figure 4.154).

Spatial coherence concerns the phase relationship between points on the same wavefront, such as A and B on the spherical wavefront shown in figure 4.155, and determines things such as speckle effects.

In this case, providing that the point or line source is sufficiently fine it is possible to have wavefronts with similar phases emerging in slightly different directions from the source O. The greater the area over which the wavefront must be coherent (spatially) the smaller must be the angle the source subtends at the wavefront. Examples of this type of interference are the Fresnel biprism and mirrors, Lloyd's mirror, the Billet split lens and the Rayleigh interferometer. When the source is too large the fringes disappear. This

disappearance of the fringes was actually used by Michelson in his method of measuring the angular diameter of stars in his stellar interferometer [119] and by Gehreke [120].

The division of the wavefront rather than amplitude can involve more than two parts, the obvious example being that of the diffraction grating (see e.g. [121]. In the case of amplitude division the interference fringes also disappear eventually owing to the temporal coherence (figure 4.156). For a typical laser this coherence length is many metres and this is relied upon when measuring distance interferometrically. However, it is not always a bad thing to have a laser which has a short temporal coherence. In fact it can be used.

In summary, temporal coherence is relevant when the interference is caused by amplitude division. Spatial coherence is relevant when the interference is caused by wavefront division.



**Figure 4.154** Relationship between temporal coherence and bandwidth.



**Figure 4.155** Spatial coherence.



**Figure 4.156** Coherence length of narrow-band spectrum.

### 4.3.6.2 *Interferometry and surface metrology*

In interference methods most systems can take as a starting point the Michelson interferometer (figure 4.158).

Interferometry is one of the most important tools of metrology. Until recently this has usually meant dimensional metrology rather than surface metrology.

Interferometry or, more precisely, the interference between two rays of light, was discovered by Thomas Young in 1802. He observed fringes (figure 4.157).



**Figure 4.157** Young's fringes.

Young found that he could not explain this behaviour with 'straight line' theory. He proposed that light was composed of waves.

It was left to James Clerk Maxwell in 1860 to provide a theoretical justification for the wave theory. He formulated the wave propagation theory based on his work in electricity and magnetism.

The basis of modern electromagnetic theory are Maxwell's equations, which are given below together with their adopted names. X denotes cross product.

$$
\begin{aligned}
&(1)\ \text{Gauss' law} && \nabla.E = 0 && \text{(a)}\\
&(2)\ \text{Magnetic monopoles} && \nabla.B = 0 && \text{(b)}\\
&(3)\ \text{Faraday's law} && \nabla.XE = -\frac{\partial E}{\partial t} && \text{(c)}\\
&(4)\ \text{Ampere's law} && \nabla XB = -\mu_0\,\varepsilon_0\,\frac{\partial E}{\partial t} && \text{(d)}
\end{aligned}
\qquad (4.160)
$$

where E is the electric field, B the magnetic field, $\mu_0$ is the magnetic permeability and $\varepsilon_0$ is the dielectric constant, altogether yielding the beautifully symmetrical wave equations

$$
\left.
\begin{aligned}
&\nabla^2 E - \frac{1}{C^2}\frac{\partial^2 E}{\partial t^2} = 0\\
&\nabla^2 B - \frac{1}{C^2}\frac{\partial^2 B}{\partial t^2} = 0
\end{aligned}
\right\}
\quad
\begin{aligned}
&\text{Where } C \text{ is the velocity of light}\\
&C^2 = \frac{1}{\mu_0\varepsilon_0}
\end{aligned}
$$

(4.161)

The solution to the wave equations reveals the oscillatory nature of light

$$
\begin{aligned}
E &= E_0\, Cos(wt + \theta)\\
B &= B_0\, Cos(wt + \theta)
\end{aligned}
\qquad (4.162)
$$

where $\theta = 2\pi/\lambda k.r$, $w = 2\pi f$ and $f = c/\lambda$ where $\lambda$ is wavelength.

The fact that these waves have sinusoidal character opens up the possibility of measuring 'phase' or relative shift between two waves (*k* and *r* are unit vectors). The question arises of seeing and measuring 'phase'. When the phase changes by $2\pi$ which is about half a micrometre (which corresponds to a time interval of 0.000,000,000,000,002 sec), it is too small to measure directly so other means have to be found.

The fringes are a direct result of the addition of fields which happen to be electrical, so the amplitude is given below as E.

$$\text{Thus } E = A\cos(\theta + wt) + B\cos(\theta_2 + wt) \tag{4.163}$$

$$\text{The intensity is } E^2 = A^2 + B^2 + 2AB\cos(\theta_1 - \theta_2) \tag{4.164}$$

Where the term $2AB\cos(\theta_1 - \theta_2)$ contains the difference in phase $\theta_1 - \theta_2$ and is within the value of the intensity.

Interference is the name given to the cross term and it is sensitive to phase difference.

The Michelson's interferometer of 1882 made use of the fringes to measure distance.

Fringes produced between the reference beam and the surface can be contour fringes if the surface and reference are normal to each other and rows of profiles if one is angled with respect to the other. The shape of each fringe can be used to estimate the roughness and the form of the surface (figure 4.158(a)). The problem is that the fringes are relatively diffuse, being of the form

$$\varphi(\delta) = \varphi(v)\cos(2\pi v\delta). \tag{4.165}$$

$\varphi(v)$ is the radiant flux of the monochromatic spectral line at wavenumber $v = 1/\lambda$. This is a case of homodyne detection because of the single wavelength of light used. The interferogram in equation (4.165) has a cosinusodal shape (figure 4.158(b)). This can be converted to the spectrum by a Fourier transformation.

Young's experiment was revealing but was not put to practical use. It was Michelson's interferometer in 1882 which earmarked interferometry as a tool [122].



**Figure 4.158** (a) Schematic optical system—Michelson interferometer. (b) Intensity pattern of tilted fringes.

The apparatus is shown in figure 4.158(*a*). One of the major uses of the interferometer is in movement of X-Y tables. Figure 4.159 shows such a set-up. Two interferometers are used, one for each axis.



**Figure 4.159** Double interferometer for XY stage.



**Figure 4.160** Practical scheme of moving stage based on polarization.

A detailed scheme of each interferometer can be seen in Figure 4.160.

If the workpiece is curved the reference waveform has to be suitably curved. The plane mirror interferometer shown in figure 4.160 is a good example of the case where the reference wavefront is a plane.

Figure 4.161 shows a simple modification that allows curved parts to be investigated. The lens in the figure changes the plane wavefront from the reference mirror. The only problem here is that curved references have to be concentric so the lens has to have adjustments in the lateral directions as well as on axis.

A slightly different configuration also uses a lens in one of the arms—in this case at the detector. This is the Twyman Green interferometer, which images the surface so that fringes produce a map of the surface. The Twyman Green interferometer has extended collimated illumination of the reference mirror and test object and is most often used for measuring form.

In effect, the image wavefront has interference with the plane wavefront from the reference. Such systems are very practical nowadays with the advent of diode lasers.

Another alternative interferometer is the Fizeau type shown in figure 4.163. Here the reference is semi-transparent so that it can be used in the same arm as the object. It is positioned in a way that allows a big separation of the mirror and reference. This is useful for measuring large parts. Notice that there is nothing in

**Figure 4.161** Curved surface measurement.



**Figure 4.162** Twyman Green interferometer.



**Figure 4.163** Fizeau interferometer.

the path between the reference and the object. Light loss in the semi-transparent references is small because the reference—object configuration constitutes a laser cavity to enable optical reinforcement. Figure 4.164 shows a practical set-up.

It should be emphasized that the basic optics of modern interferometers is the same as in the past. Methods of making the optics, the source and the detectors are where changes have been made.

In Figure 4.164 a rotating diffuser is inserted to reduce speckle. It is customary to place a pinhole in the source path at the back focal plane of the collimator to clean up the source. This is imaged at the pinhole (spatial filter).

To revert to the basic interferometer—the Michelson, this has spawned a number of variants such as the Fabry Perot and Linnik interferometers.



**Figure 4.164** Practical system.

Really sharp profile fringes are obtained by the use of a Fabry-Perot interferometer in which multiple reflections take place—a practical embodiment of surface interferometers.

The interferogram can be converted from a contour view to a series of parallel traces similar to those obtained by the parallel tracking of a stylus instrument simply by lifting up an edge of the optical flat (figure 4.158(b) and 4.166). This produces oblique sections in parallel across the field of view and a very acceptable picture of the surface. The use of higher-power optics (with the higher NA values) can be very useful for examining scratches and flaws as well as some of the microgeometry.

The system devised by Linnik [8] using a variant of the Michelson interferometer is still regarded as one of the best. In this the reference arm is removed from the surface and a comparison obtained using the conventional two arms, as seen in figure 4.165.

The well-known method of reducing the width of the fringes and increasing the sharpness by means of multiple reflections has been highly developed by Tolansky [17]. This method refers to figure 4.166 because the optical flat has to be close to the surface. In this technique both surfaces have to be highly reflective, which is often the case if very fine surfaces are being measured.

**Figure 4.165** Linnik interferometer.



**Figure 4.166** Schematic diagram of surface interferometer (Tolansky): (*a*) reference parallel to test surface; (*b*) reference tilted relative to test surface.

### 4.3.7 Heterodyne methods

#### 4.3.7.1 Phase detection

Phase shifting interferometry is the method of detecting and using the very small distances corresponding to phase by measuring fringe intensity. Before any movement or shift in optical path length of the test arm occurs the starting phase has to be determined. Thus if the intensity across one fringe is measured to give $I_1$, $I_2$, ...... the phase $\theta$ is determined by

$$\theta = \tan^{-1} \frac{I_1 - I_3}{I_2 - I_4}.$$

$$(4.166)$$

An alternative method of measuring movement or shift is to keep the phase the same and to change the wavelength of the source. Its distance is to be measured rather than a shift in the wavelength changed by means of a frequency modulator at the source. A schematic diagram is shown in Figure 4.167.

**Figure 4.167** Interferometer with frequency change.

### 4.3.7.2 *Frequency-splitting method*

Another method similar to Nomarsky, but much more sophisticated, is due to Sommargren [123]. The basic component is a phase detector in effect, although it can equally be placed here because of the commonality of components. It comprises a Wollaston prism, a microscope objective and surface orthogonally, linearly polarized beams [124] with a frequency difference of 2 MHz. This method is a classical example of the heterodyne technique in interferometry in which, in this case, the combination signal of two slightly different signals in frequency shows much lower frequency variations which can give, via phase detection, an accurate measurement of position. (This is sometimes referred to as envelope detection in radio.)

As shown schematically in figure 4.168, these signals are refracted slightly by the Wollaston prism on either side of the normal. The objective brings them to a focus at two points a and b on the surface about 100 $\mu$m apart. Each spot is about 2 $\mu$m. The 2 MHz frequency difference is achieved by Zeeman splitting due to a magnetic field across the path between the laser cavities.

The beam from the laser is split by beam splitter B into a reference path and the test path. The reference path has phase bias adjustment by means of the $\lambda/2$ plate so that the average phase measurement takes place in the centre of the range of the phase detector. Changes in the height $\Delta h$ of the surface between the two points a and b on the surface $\Delta h$ introduce an optical path change $\Delta z$ where

$$\Delta z = \frac{4\Delta h}{(\mathrm{NA}^2 + 4)^{1/2}}$$
(4.167)

where NA is the numerical aperture of the microscope. Thus, since the phase change $\Delta\varphi = (2\pi/\lambda)\Delta z$

$$\Delta h = \frac{\lambda}{8\pi}(\mathrm{NA}^2 + 4)^{1/2}\Delta\varphi$$
(4.168)

so if a is kept constant, for example by situating it at the centre of rotation of the part, a circular profile of the surface around a (or b) is possible. Because of the circular track used by this system its stability can be checked at the start and finish of each cycle.

The computer system extracts the values of $\Delta h$ after first removing the eccentricity signals due to the tilt of the workpiece on the roundness instrument and any second harmonic present due to bowing.

The results show that the out-of-roundness of the rotation needs only to be less than 2 $\mu$m (the spot size). Axial deflection is second order in importance. Sensitivities of the order of 1 Å have been obtained, giving a ratio of range to resolution of about one part in 5000, which is a relatively high value.

Other heterodyne methods have been devised which claim sensitivities of 0.1 Å [125] in which there have been two beams, as in reference [124], but this time the beams are concentric, one being of 2 $\mu$m diameter and

**Figure 4.168** Sommargren heterodyne probe.

the other 50 $\mu$m. As in all heterodyning methods two frequencies of the source are used—in fact in this example the same source is used with one spot being frequency shifted by about 30 MHz by means of a Bragg cell rather than Zeeman splitting as in reference [124].

Because of the concentricity of the two beams any vibration introduced into the system obviously introduces about the same phase change in the optical paths of both beams. Hence measuring differential phase eliminates the effect—again the skid principle!

Assuming one beam to have signal $E_1$ and the other $E_2$

$$E_1 = A\cos(\omega_0 t + \varphi_0 t + \varphi_s(t))$$
$$E_2 = B\cos(\omega_0 t + \omega_a t + \varphi_0 + \overline{\varphi}_s(t) + \varphi_a)$$

(4.169)

where $\varphi_s(t)$ is the phase produced by the surface, $\varphi_0(t)$ is the phase produced by vibration, $\overline{\varphi}_s(t)$ is the averaged phase produced by the surface for a large spot and $\varphi_a$ is the electronic phase produced by $\omega_a$– the shift.

The heterodyne detection is achieved by multiplying $E_1$ and $E_2$ and using a square-law detector. Thus

$$I = E_1 \times E_2$$
$$= \text{ac} + \text{dc}$$

(4.170)

where

$$AC = \tfrac{1}{2} AB \sin(\omega_a t + \overline{\varphi}_s(t) + \varphi_a)$$
$$DC = \tfrac{1}{2} AB = \text{ direct current through detector.}$$

(4.171)

Hence

$$\frac{AC}{DC} = \sin(\omega_a(t)) + \overline{\varphi}_s(t) - \varphi_s(t) + \varphi_a).$$

(4.172)

$\varphi_0$ is therefore eliminated in equation (4.172). Consequently the surface finish $\overline{\varphi}_s(t) - \varphi_s(t)$ is measured in the presence of vibration.

With the progressive development of phase-sensitive detectors it is highly likely that these heterodyne methods will grow in use. The problem is their limited range and relatively slow speed.

Variants on interferometers which do not use moving members yet have a wider range than conventional interferometers have been developed. One such method uses an oblique approach similar to the non interferometric optical probe [126] (figure 4.169).

The idea behind this method is to increase the effective value of $\lambda/2$, the fringe spacing due to the obliquity factor. In fact this works out quite simply from the geometry to be

$$\frac{\lambda'}{2} = \frac{\lambda}{2 \cos \theta}.$$

(4.173)

For $\theta$ up to 75° the effective value of $\lambda'/2$ is $2\lambda$ which is a useful gain of four over other devices. However, it does lose equally on resolution and introduces problems associated with directionality of the surface and non-circular illumination. But it does show, as do all the methods mentioned so far, the fantastic versatility of optical methods using relatively standard components. Recent developments highlight this comment, for example in simple interferometers.



**Figure 4.169** Mirau interferometer

Here the whole sensitivity can be increased by means of optical path difference multipliers [127, 128] (figure 4.170). It is, however, also very sensitive to tilt.

Instead of having fixed plane mirrors, combinations of corner cube prisms (figure 4.170), or corner cube and right-angled prisms, can be used which produce a much larger path difference for a given movement of the movable mirror than the straightforward Michelson interferometer.

**Figure 4.170** Multiple path distance interferometer (from [127]).

So in the case above, if $b$ is the distance between reflected beams side by side, and $a$ is the distance between beams biaxially, if the corner cubes are transversely shifted by $b/2$, then the diameter $D$ of the corner cube required to get a multiplication of a is given by the expression

$$a^2 + (\alpha - 1)^2 b^2 \leqslant (D - d)^2 \qquad (4.174)$$

where $d$ is the diameter of the laser beam.

In the case of simple biaxial reflections

$$\alpha(a + d) \leqslant D. \qquad (4.175)$$

This trick enables fringes to be obtained at mirror displacements of $\lambda/10$ or less.


### 4.3.7.3 Other methods in interferometry comparable with heterodyne methods

Another interferometric method relying on phase measurement for the evaluation of surface height is described below [129]. The basic idea is one which endeavours to remove the problem of fringe interpretation, which has proved to be so difficult in the past in getting interferometry accepted as a simple tool for surface measurement. It consists of producing more than one fringe pattern from one surface by some means or other, storing them and then unravelling the surface deviations on a computer. The problem then condenses from trying to correlate relative phases of points in the 2D field of the surface by eye using the fringe pattern to that of measuring the difference in phase changes produced by the path difference of the rays from the test surface and reference surface point by point.

There are a number of ways of achieving the phase changes. One is step-by-step discrete positioning of the reference relative to the surface. Another is to have the reference moving at constant speed. Take the former.

The intensity at a point $(x, y)$ is

$$I_1 = A + B \cos(\varphi(x, y)). \qquad (4.176)$$

The reference is moved by $\pi/2$ giving

$$I_2 = A + B \cos(\varphi(x, y)) + \pi/2 \qquad (4.177)$$

and by $\pi$

$$I_3 = A + B \cos(\varphi(x, y)) + 3\pi/2). \qquad (4.178)$$

The three positions enable the three variables $A$, $B$ and $\varphi(x,y)$ to be found where $A$ and $B$ are ambient light and the intrinsic amplitude at $(x, y)$. From these equations (See 4.166)

$$\varphi(x, y) = {}^{-1}\left(\frac{I_3 - I_2}{I_1 - I_2}\right).$$ (4.179)

The beauty of this method is that because of the subtractions and additions at each point, local sensitivity problems in the optical receiver or light intensity are cancelled out, a strategy often used in glossmeter techniques described later on.

Knowing $\varphi(x, y)$ the surface height can be found. Thus

$$z(x, y) = \varphi(x, y)\frac{\lambda}{4\pi}.$$ (4.180)

Wyant *et al* [129] have pointed out that this step of $\pi/2$ is not necessarily the only one. Obviously any incremental step can be used but it makes sense to simplify the mathematics. They in fact start off one device with a $\pi/4$ initial shift.

For example, if $I_n$ is the intensity of the interference pattern obtained by stepping the phase by $n2\pi/N$, the phase can be obtained by using the following equation:

$$\varphi = \tan^{-1}\left[\sum_{n=1}^{N} I_n\sin\left(n\left(\frac{2\pi}{N}\right)\right)\Big/\sum I_n\cos\left(n\left(\frac{2\pi}{N}\right)\right)\right].$$ (4.181)

Note that the resulting intensity is multiplied by a cosine and sine. This is equivalent to homodyne detection (synchronous) used in communications to extract the signal from a carrier.

Another method [130], is called the integrating bucket or window method. This is more difficult to compute but has the advantage that the effect of vibration can be reduced.

These methods are restricted to phase shifts of much less than $\lambda/2$ (usually less than $\lambda/4$) otherwise ambiguities occur.

A practical embodiment incorporates a Mireau interference microscope (figure 4.171) described by Breitmeier [131]. The whole assembly of microscope objective interferometer and beam splitter is mounted on a piezoelectric stage driven by a computer. Using such a technique very high-quality results can be achieved using a follower principle.



**Figure 4.171** Mirau surface interferometer (UBM, WYCO).

Such modulation methods are also used with the technique of speckle but have not yet produced comparable answers.

Like all methods involving phase measurement, certain assumptions have to be made about the nature of the surface. However, for fine homogeneous surfaces correlation with stylus methods is high, yet for complicated shapes there is much less agreement because of the phase limitations of the former.

### 4.3.7.4  Relative merits of different nanometre instruments

The question arises as to which of the methods outlined above is the best. In order to answer this the signal-to-noise value for the Sommargren (or frequency heterodyne) type of microscope and the point-to-point instrument using the Mirau interferometer will be compared. Both make use of phase heterodyne interferometry. However, they differ widely because the former constitutes an optical stylus instrument with an output signal of continuously varying phase, whereas the latter derives phases of points of a varying fringe pattern from just four (or three) intensity values issuing from four discrete modulations of the fringe pattern.

Although one type of instrument works with a Zeeman effect laser stabilized in both frequencies and the other—the point-to-point instrument—uses an arc lamp with a bandwidth of about 10 nm, they both are subjected to the same sort of factors affecting performance.

Contrast is the decisive quantity in determining vertical resolution, whether of fringes or visibility. If the contrast is low then there is a loss of light to a standing background illumination. This useless light gives rise to a standing photoelectric current which adds to the noise power but not to the signal.

The factors influencing the noise of the optical systems are the following:

1. Roughness of optical elements. The proportion of light lost is given approximately by $1 - R_q^2 K$, where $R_q$ is the composite root mean square roughness of the elements.
2. Unequal intensities of reference beam and measuring beam. The combined total intensity is given by

$$I_2 = A_1^2 + A_2^2 + 2A_1 A_2 \cos(\psi_2 - \psi_1) \tag{4.182}$$

   where the subscripts indicate reference and measured light intensity. The third term is the interference term. For optimum contrast, $A_1 = A_2$. In the case of the Mirau instrument there is usually a choice of reference surface of different reflectances available so that the best can be chosen for optimum conditions.
3. The finite bandwidth of light accepted from a spectral filter is symptomatic of pure illumination. If there is an optical path difference of A between the beams the phase difference is spread over the available wavenumbers. The loss of contrast is

$$\int_{k_{low}}^{k_{high}} (I_k / I_{mean}) \exp(-jk\Delta) \, dk \tag{4.183}$$

   where $k_{high} - k_{low}$ is the light bandwidth.
4. Inclination of the surface affects the amount of light collected by the objective beam. The loss is dependent on NA and the slope values, and affects the interference signal directly.
5. The longer spatial wave components of the surface or the form of the surface can be of an amplitude of height variations which are comparable with the depth of focus. The effect is to cause a variation of phase within the core of the measured beam. The contrast is reduced by a factor

$$\mathrm{sinc}(\tfrac{1}{2} k l \alpha^2) \tag{4.184}$$

where $l$ represents the path differences between beams and $\alpha$ is the NA. The same phenomenon causes a systematic phase shift in the measured beam emanating from the optical element from its nominal values of $\frac{1}{2} k l \alpha^2$.

### (a) Calculation of signal-to-noise ratio: Mirau instrument, signal acquisition

In the following calculations the data assumed for the CCD photodiode matrix is abstracted from the published data in respect of an early type of $64 \times 64$ elements. Thus the value assumed for the mean charge generated and stored on a detector element during exposure is the optimum of half the maximum at a value of approximately 1 pC. The charging photoelectric current on average per element, with an irradiance of green light on the chip overall of 200 $\mu$W cm$^{-2}$, is estimated from the data at 15 pA. This gives an integration time of 0.07 seconds.

Of the two most familiar types of noise, shot noise and Johnson noise, the former has influence in the acquisition of data, but not the latter, because resistive noise sources are not involved in the storage process.

Ostensibly the shot noise is given as

$$i_n^2 = 2 I_m q B \tag{4.185}$$

where $i$ is the RMS noise current, $q$ the electronic charge, $I_m$ the mean current and $B$ the bandwidth.

The effects of the steady dark current (i.e. one of zero frequency) and the average of the dark-current variations of periods much greater than the integration value are eliminated in the numerical procedures of compensating for integrator output drift which is described above. Thus the frequency band required for the storage process is reckoned to extend from 1 to 40 Hz.

Hence the intrinsic value for shot noise current $i_n^2 = 2 \times 1.6 \times 15 \times 40 \times 10^{-31}$ A$^2$ or $1.4 \times 10^{-14}$ A.

The RMS signal current for a fringe displacement at 100% contrast is 10.5 pA. Hence the signal-to-noise ratio by shot noise alone during acquisition is 700:1.

From available data it seems that the dark current should not exceed 25% of the photoelectric current due to illumination. But in the low-frequency band of the integration process the dark-current variations constitute a flicker noise which varies according to a $1/f$ law. The noise of this component is calculated from the expression. $I_{nd}(\sqrt{f_C} \log(f_H/f_L) - f_C + f_H)$ The corner frequency $f_C$ is estimated at 100 Hz and the dark shot noise component is one-quarter of $I_{nd}$ times the noise current. Hence, the component of $1.12 \times 10^{-14}$ A due to flicker noise has to be added; the signal-to-noise ratio drops to 546:1.

### (b) Signal retrieval

The retrieval of the signal from the photodiode elements during readout is a very fast process. As the scan may allow a discharge time of typically 0.5 microseconds only per element, the output signal current is increased over the charging current by a factor of $10^{10}$ (i.e. to 2 $\mu$A). The corresponding shot noise current of readout, which must be added to the noise component accumulated in the storage process, also increases in inverse relation to the timescale. This is because the noise is proportional both to the square root of bandwidth and to the square root of signal current. Hence another noise current in the ratio of 700:1 is generated.

The output amplifier is of the current type and sinks the signal current into a resistor of a value of 22 k$\Omega$ typically. The Johnson noise current is given (at 30°) by $1.3 \times 10^{-16} \sqrt{B/R}$ A, where $R$ denotes the resistor value. The evaluation assuming a bandwidth of 4 MHz is for a noise current of $1.3 \times 10^{-9}$ A, that is a signal-to-noise ratio of $10^3$.

The amplifier noise is probably in the region of an equivalent noise voltage of 10 nV Hz$^{-1/2}$. Thus at 4 MHz bandwidth the noise current in the resistor is approximately $10^{-9}$ A.

For the data retrieval the signal-to-noise ratio is 525:1; the overall signal-to-noise ratio is 377:1; and the shot and Johnson noises are very roughly in balance, representing efficient timing in the chip. Very fast rates

of output scan meet the barriers of crosstalk and switching transients. Thus performance is limited by these in readout and dark current in storage.

The resolution of the Mirau instrument is 1/377 radians of a fringe at the green-light wavelength or 0.01 nm, which is well within the specification.

### (c) Calculation of typical signal-to-noise ratio: frequency heterodyne instrument

The frequency heterodyne instrument is studied in respect of signal-to-noise ratio. The beam power of 200 $\mu$W is incident upon a single-element detector, related directly to the single optical element of the optical stylus. Assuming a responsivity of 0.2 A W$^{-1}$ a signal current of value 40 $\mu$A is generated.

As a resolution of 0.01° is required, it is probably necessary to change the carrier frequency of 2 MHz to approximately 1 kHz where the electrical angle resolution of that order may be measured digitally. Thus a final bandwidth of 4 kHz is appropriate.

The shot noise current is thus $2.4 \times 10^{-10}$ A and the corresponding signal-to-noise ratio is approximately $10^5$.

Assuming a current amplifier sink resistor of value 22 k$\Omega$, the Johnson noise current is $6 \times 10^{-11}$ A, giving a signal-to-noise ratio of $5 \times 10^5$. The amplifier noise is not very important but it should be possible to fit an amplifier with an equivalent noise voltage of 5 nV Hz$^{-1/2}$, and therefore a value of $1.4 \times 10^{11}$ A of noise current in the output resistor should be attained.

Thus the overall signal-to-noise ratio is controlled by the shot noise at a value of about $10^5$. In this respect the frequency heterodyne instrument is better than that of the Mirau one by over two orders. But it is essentially a profile instrument, whereas the Mirau instrument is designed around an areal capability.

For both instruments a reduction in fringe contrast has a direct effect upon performance, the useless light still contributing to noise while the loss of contrast reduces the signal in direct proportion.

### (d) Comparison

The electro-optics of the frequency heterodyne instrument has a reserve factor of about 250:1 in meeting its resolution of 0.1 Å whereas the Mirau has a reserve factor of about 40:1 in respect of its resolution aim of 3 Å. Both types are clearly adequate to meet the design specification. Indeed it seems from these calculations that both could increase their specification.

### (e) Discussion of specification: Mirau models

The NA for the instrument with $10 \times$ objective is about 0.17, yielding a resolution of 3.8 $\mu$m. A maximum surface slope is taken to be one half fringe divided by one sample interval, which gives 5.9°. The depth of focus is theoretically 19 $\mu$m but is sometimes reduced to preclude the spread of signal from one photodiode element to the next.

### (f) Frequency heterodyne model

One of the most striking features is the circular trace of the measurement. This is perfectly acceptable for flat surfaces but poses questions on shaped ones.

The change of frequency to 1 MHz from 2 MHz adds to the time of data acquisition. The total acquisition time of 30 seconds typically allows about 30 milliseconds per data point, given a millisecond to sample, indicates time averaging to reduce the effect of any vibration in the system.

The NA of the objective indicates an NA of about 0.5. The range of the vertical measurement seems to be restricted to one half wavelength of the light and not more!

In conclusion it seems that the heterodyne frequency method (e.g. Zygo) is very suited to atomic-level measurement of accurately controlled surfaces, whereas the Mirau type (Wyko) is more robust and versatile. It is suited to subnanometre measurement but hardly beyond.

The calculations above give some idea of the factors involved in producing an optical instrument to the high specification required. The figures demonstrate that most manufacturers are sensibly pessimistic in fixing specifications.

### 4.3.7.5 *High-precision non-contacting metrology using short coherence interferometry—white light scanning interferometer*

This technique makes use of what has previously been regarded as a weakness of an optical system, namely a short temporal coherence length. Basically the concept is that if a source is wideband spectrally its coherence length will be short. Furthermore, its coherence function will decay in a well-defined way according to the shape of the spectrum. This known decay of the coherence can be used to measure absolute distance rather than relative distance which is usually the case for interferometers.

In conventional Michelson interferometry one fringe looks exactly the same as the next one to it.



**Figure 4.172** Fringe counting.

As the fringes are usually separated by $\lambda/2$, a distance $L$ in figure 4.172 can be measured simply by counting the number of fringes encountered as the detector moves across the fringe pattern. In the figure $L = {}^{3\lambda}/_{2}$ where $\lambda$ is the wavelength. However, if the power is cut there is no distinguishing feature between fringes when power is restored. The detector is literally lost. The reason for the identical fringes is the long coherence of the light. White light has a very short coherence length. Under these circumstances each fringe has an address!

There is a Fourier transform relationship between the coherence length of the light and its bandwidth (figure 4.173).



**Figure 4.173** Coherence length via Fourier transform.

$$\left.\begin{array}{l} F(w) = 2 \int_{0}^{\infty} f(x) \cos wx \, \mathrm{d}x \\[2mm] F(w) = \dfrac{2}{\pi} \int_{0}^{\infty} f(w) \cos wx \, \mathrm{d}w \end{array}\right\} \qquad (4.186)$$

$L_c$, the coherence length for optics is the same as the correlation length of the autocorrelation (equation (4.181)), where $L_c$ is given by

$$L_c = \int_0^\infty | f(x) | \, dx$$

(4.187)

The sensing system relies on (i) coherence-modulated interferometry and (ii) coherence discrimination.

Coherence modulation of the fringe field at the output of a Michelson interferometer derives from the coherence function of the source (figure 4.174) where the visibility $V$ is given by equation 4.188

$$V = \frac{I_{max} - I_{min}}{I_{max} + I_{min}}.$$

(4.188)

The coherence function of the source is a function of the absolute path difference $(l_1 - l_2)$ in the interferometer. This therefore enables the absolute fringe order number to be determined from the fringe envelope modulation providing that the difference in visibility between different fringes is resolvable above system noise (i.e. $\Delta V$). See figure 4.180 for the various bandwidth options.



**Figure 4.174** Broadband absolute interferometry.

Coherence discrimination (figure 4.175) enables the interference between the reference and measurement arms in the remote sensor to be detected by a transceiver-mounted tracking interferometer.

Coherence discrimination is achieved by adjusting the relative path length of the beams to satisfy equation (4.189) (figure 4.175)

$$| (a + c) - (b + d) | \ll l_c$$
$$| (a + d) - (b + c) | \gg l_c$$

(4.189)

where $l_c$ is the coherence length of approximately 100 $\mu$m.

The probe geometry for surface finish and profile is shown in figure 4.176. Measurement is made by using the reference beam effectively as a stylus instrument uses a skid. The central spot varies as the coherence modulation with roughness height.

An important feature of coherence discrimination is that it allows the remote sensor to be connected to a transceiver by a common path optical link in which the interfering beams are separated only in the time domain. Such a link is inherently immune to random noise. The sensor need not have the Michelson configuration.

**Figure 4.175** Broadband fibre optic interferometer.



**Figure 4.176** Surface finish variant.

The interferometer output can be expressed in spectral or temporal domains:

1. Spectral:

$$I_t = T(x,\lambda)I(\lambda) = \frac{I_i(\sigma)}{2}[1 + \cos(2\pi\sigma\delta)]$$

(4.190)

where $\sigma = 1/\lambda$, $I_t$ is transmitted intensity, $I_i$ is incident intensity, $\delta$ is path length (a function of $x$), $I(\lambda)$ is the source spectrum and $T(x, \lambda)$ is the spectral transmission. The spacing of adjacent spectral peaks in the output is given by $\Delta\sigma$ where $\Delta\sigma = 1/\delta(x)$.

The number of cycles of spectral modulation will depend on $\Delta\sigma$ and, hence $\delta(x)$, and also the extent of the input spectrum.

2. Temporal: the output $I(t)$ is given by

$$I(t) = 2I[1 + \gamma\cos\varphi(t)]$$

(4.191)

where $\varphi(t)$ is the time-varying phase and $I$ is the intensity of the interfering beams. In this form it can be seen how the coherence function $\gamma$ modulates the output. From this the phase (and $x$) can be determined absolutely.

Obviously equations (4.190) and (4.191) are a Fourier transform pair and so are equivalent.

Such a development as the wideband interferometer is very important for the shapes and profiles of objects such as cams and gears, as well as texture. Perhaps, if the spectral shape could be fixed to this technique, there is great potential here.

One way of building an absolute interferometer resembles the way of using echelons for calibrating length standards. Basically sets of coherent envelopes have to be linked together in order to extend the range of absolute measurement. This has been done, as mentioned above. A simple way is to use a CCD camera in such a way that each pixel acts as if it were a contact probe. Then the trick is to record the maximum peak contrast as a function of the scan; this is a way of linking coherence bundles from each pixel to the frame of the image which can then be linked together or so-called 'stitched' together. Frame boundaries have to be constrained to make the connection possible.

'Stitching' is a technique often used to enlarge the range of measurement. One simple example is in the measurement of aspherics [132]. In this the geometry of the aspheric is divided into several diametral zones. Each zone is then specified as a surface of revolution usually in terms of ellipses and polynomials.

The basic idea behind stitching interferometry is to divide the wavefront of the measured surface into several segments so that the fringe density over a sub-interferogram from each measurement remains resolvable. In other words to shift the object with respect to the reference waveform yet retaining positional integrity.

Moving the object through the reference waveform 'sweeps' the fringe pattern in turn throughout the object making sure that at each position the fringes are resolvable. This sequential focusing provides a set of sub-interferograms which are combined after corrections for aberrations.

Range extension can be along the axis or in the lateral direction. In the aspheric case (figure 4.177) the synthesis is an example of normal axis stitching.

The peak fringe contrast is recorded as a function of scan position. This mode of operation is called 'scanning white light interferometer' SWLI'. The length unit becomes the increment between frames compared with the wavelength of light for phase interferometry. The frame scan increment can be associated with the coherence of the source as shown in figure 4.178.



**Figure 4.177** Stitching in the normal mode.

There are two ways of utilizing the white light (broad band) fringes. The one method is to consider the coherence envelope.

In the former case it is considered be possible to get 3 nm positional accuracy by addressing the shape of the envelope. On the other hand use of the fringes within the envelope by measuring phase can be much more sensitive allowing 0.1 nm to be resolved. In some cases it is possible to use both coherence and phase. The coherence is used for 'coarse' positioning and the phase for high precision 'fine' position. There is no ambiguity if the two are used together because of the absolute identification of the individual fringes.

See figure 4.175 for an example of tracking interferometry using broadband (white light) interferometry.

White light interferometers have a very important feature. This is the presence within the instrument of an absolute reference for position. Sometimes this is called the positional datum or zero of the instrument. It occurs when the optical path difference (OPD) between the arms produces a maximum fringe intensity. Measurement of distance or position all relate to this point. There can be a zero point in other interferometers but these are ambiguous. Figure 4.180 shows the options (not to scale).

**Figure 4.178** Lateral association of frame scan.

(*a*) Phase of fringes                    (*b*) Envelope of fringes



**Figure 4.179** Envelope and phase detection.

In (*a*) there is no discrimination between the fringes, although they are easy to count. There is no 'absolute' datum. This is the conventional interferometer. In (*b*) more than one wavelength are mixed. For, say, two frequencies $f_1$ and $f_2$ the relative fringe intensity takes the form $\cos 2\pi(f_1 - f_2)\cos 2\pi(f_1 + f_2)$ where $f_1 - f_2$ is the envelope and $f_1 + f_2$ the carrier.



**Figure 4.180** Bandwidth options.

Interference of two
simple wavefronts

Reference
flat

Part

Speckle images
(rough surface)

**Figure 4.181** Transition fringe to speckle.

The maximum intensity occurs periodically at $\lambda = \left( \dfrac{1}{f_1 - f_2} \right)$.

This option has a wide range with a number of zeros. Option (*a*) has no zeros but unlimited range.

Option (*c*) is the white light alternative having a definite unique zero but a small range that can only be extended by stitching.

In shape and position monitoring using coherent light there is one factor which has to be taken into account. This is the magnitude of the surface texture. The presence of surface roughness can considerably degrade any fringe pattern [122].

Figure 4.181 shows the transition from fringes produced between fine surfaces and which are typical of Twyman Green, Michelson and even sharper with Tolanski. As the surface becomes rough the fringes degenerate into speckle in which local height variations are big enough to destroy coherence. This degradation of the fringes can be used to measure the roughness, as will be seen later.

### 4.3.8  Moiré methods [133]

#### 4.3.8.1  General

Any system which has the superposition of two periodic structures or intensity distributions can be called a moiré system. The name moiré probably originates from the French textile industry and has the meaning of wavy or watered appearance. Moiré fringes have been used in the silk industry since the Middle Ages for quality control.

Moiré fringes can be used as an alternative to interferometry and holography and are becoming increasingly used for measuring form, but not as yet surface texture to any large extent because of the difficulty of making gratings whose width is much less than 10 $\mu$m.

Figure 4.182 shows the principle. If there is a rotational mismatch between two identical gratings of, say, $\delta$, it produces a new system of equidistant fringes which are easily observed. These are called moiré fringes. The pitch of the fringes is $P_{\mathrm{m}}$ where

$$P_{\mathrm{m}} = (P/2)\sin(\delta/2) \tag{4.192}$$

where $P$ is the spacing of the original gratings (figure 4.183).

This equation is similar to that obtained describing the pattern of two interfering plane waves. For this reason the moiré effect is often called 'mechanical interference'. There is a difference, however, because in

**Figure 4.182** Moiré fringes created by two successively positioned periodic structures.



Cens of dark Moiré fringes

Cens of bright Moiré fringes

**Figure 4.183** Moiré fringes created by rotational mismatch. $N$ moiré index, $m$, $n$ number of grating lines, $p$ pitch, $\delta$ intersection angle.

interference the interference term only consists of the difference of the optical path length. The moiré superposition, like acoustic beats, contains two terms, the sum and the difference frequencies.

The sum term in the moiré fringes has nearly half the pitch $P$ so that when using dense gratings (~10 line pairs per millimetre) it is not possible for the eye to observe these fringes. Usually, however, the additive moiré fringes are rarely used and it is the difference fringes which are used. The number of fringes produced is related to the number of lines on both of the other gratings. If they are $m$ and $n$ respectively, then the number of moiré fringes $N$ is given by

$$N = m - n \tag{4.193}$$

Along each moiré fringe the index is constant. Moiré fringes can also be generated by a pitch mismatch, not an angular one. The resulting moiré fringe is then given by

$$P_{\mathrm{m}} = \frac{P_2 P_1}{P_2 - P_1}. \tag{4.194}$$

These pitch mismatch fringes are sometimes called Vernier fringes. Sometimes, even though an optical system cannot resolve the pitch of either of the two generating gratings, it will easily see the moiré fringes.

Gratings utilized in moiré applications usually have a density of between 1 and 100 line pairs (1p) per millimetre [134].

In application the fringes are projected onto the surface under test. Any distortion of the shape or form of the surface shows itself as a deviation from the ideal shape of fringe which is projected. The moiré fringe projection is used mainly in three surface applications:

1. Strain analysis for determining in-plane deformations and strains of a surface.
2. The shadow and projection moiré methods for the determination of the contour of an object or for the comparison of a displaced surface with respect to its original state.
3. The reflection moiré and the moiré deflectometry methods for the measurement of the slope (out-of-plane) distortion of a surface with respect to its initial state or the measurement of waveform distortion produced by specularly reflecting objects or by phase objects.

### 4.3.8.2 Strain measurement [135]

In strain measurement the grating is usually placed in contact with the surface. The orientation of the grating has to be in the direction of the assumed deformation. Prior to deformation the grating with the pitch $P$ is given by

$$y = mP \quad m = 0, \pm 1, \pm 2 \ldots . \tag{4.195}$$

After deformation the shifted surface distorted the grating to

$$y + s[u(x, y), v(x, y)] = nP \quad n = 0, \pm 1, \pm 2 \ldots . \tag{4.196}$$

In this equation the displacement $s[u(x, y)]$ of each point $x$, $y$ with the two independent components $u(x, y)$, $v(x, y)$ can be observed with the moiré effect. Positioning an undistorted grating with the same pitch near the distorted grating gives an optical magnification due to the moiré effect. The subtractive moiré fringes give the fringes of displacement:

$$s[u(x, y), v(x, y)] = (m - n)P = NP \quad N = 0, \pm 1, \pm 2 \ldots . \tag{4.197}$$

This equation gives the displacement only in the $y$ direction. To get it in the $x$ direction both gratings have to be rotated in the $x$ direction.

### 4.3.8.3 Moiré contouring

Out-of-plane methods are triangulation techniques; they are used to measure the form of plane or low converse shapes of mostly diffuse reflecting surfaces.

### 4.3.8.4 Shadow moiré

In this the master grating is positioned over the surface (figure 4.184). It is then illuminated with light at angle $\alpha$. The shadow region is thus viewed at angle $\beta$ through the same grating. Obviously there are $m$ lines in the incident beam and $n \neq m$ in the emergent beam. Hence the moiré fringe number $N = m - n$.

The difference in height between the surface and the master grating is CE:

$$CE = \frac{NP}{(\tan\alpha + \tan\beta)} = \frac{(N + \frac{1}{2})P}{(\tan\alpha + \tan\beta)}$$
$$\text{bright fringe} \qquad \text{dark fringe} \tag{4.198}$$

and the height sensitivity $\delta(CE) = CE/N$. In this case there is a linear relationship between the height difference and the moiré order $N$.

**Figure 4.184** Shadow moiré.

This type of configuration is mostly used for small parts.

Limits to this method are the resolving power of the optics, which can cause a fading of the fringe contrast. This can be controlled by making sure that a point on the surface is not greater than five times the grating pitch from the grating. Another factor is the 'gap effect', which is the distance between the object and gives a lateral displacement of $xCE^2P$ where $x$ is the point on the surface. This effect is reduced by using telecentric projection. Diffraction can also be a problem on high-density gratings. Again, the surface and grating should not be separated much.

### 4.3.8.5 Projection moiré

Unlike shadow moiré the master grating is remote from the surface. Here there are two gratings: the projection grating and the reference grating. The projection grating can be of any sort, such as one etched onto glass and projected with white light, or it can be a Michelson fringe pattern. The pattern made, however, is thrown onto the surface.

The pitch $P_0$ on the object is $P_0 = mP$, where $m$ is any magnification introduced by the projection (figure 4.185).

In this method the fringe is formed by the viewing of the deformed object fringe through the reference grating in front of the camera. If the distorted grating on the surface has $n_s$ line pairs per millimetre the height sensitivity is $1/(n_s \tan \theta)$. This is strictly non-linear. The usual best resolution of this technique is 100 line pairs per millimetre.

A variant of this is not to use a master reference grating before the camera but to use the columns of the CCD array in the camera to provide the reference. This is cheaper but the observer cannot control the moiré fringe density $N$. This is decided by the chip-maker of the CCD.

In the methods outlined diffuse surfaces are needed. In the case of a mirror-like surface, a different strategy is used. The surface is treated as a mirror despite its deformations (figure 4.186).



**Figure 4.185** Projection moiré.

**Figure 4.186** Moiré using surface as mirror.

The specular reflection combined with the deformation, often of the grating on the surface, is viewed via a semi-silvered mirror by means of a camera. This method can be used in a similar way to double-pulse holography in which the 'before and after' scenes of the grating on the surface are recorded on the same film.

*4.3.8.6 Summary*

Moiré techniques have been used up to now for estimating rather crude deformations on surfaces. Typical gratings have had spacings of 25 $\mu$m. They have been regarded as a coarse extension of the other fringe methods, that is interferometry and holography. This role is changing with the advent of much better gratings, better resolving optics and pattern recognition of digital fringe interpolation methods.

*4.3.9 Holographic techniques*

*4.3.9.1 Introduction*

The vital concept of producing a phase record of light which has been diffracted by some object can be attributed to Gabor [136]. The thinking behind the technique which is now called holography was developed over many years by Gabor, in connection with the resolution limits of electron microscopes [137].

It is convenient to regard holography rather as one would a single-plane interferometer. That is to say that by using a recording medium such as a photographic plate, it is possible to record a spatial slice of reflected light from the object in such a way that an amplitude-modulated record is produced which exactly represents the phase structure emanating from the object (see figure 4.187). The amplitude-modulated data recorded on this slice of material can subsequently be used to phase-modulate a plane wavefront to recreate the original complex mapping of intensity information in the object space. Thus the hologram may be thought of as a grating or set of gratings each having some particular spacing and orientation, and it is this which has the ability to reconstitute the wavefront so that the object appears to occupy its original position in space.

Apart from the spectacular three-dimensional characteristics of the reconstructed object field, it is important to recognize that, for the first time since the inception of photography, a data recording technique which completely exploits the information-holding capabilities of the photographic medium is possible. Indeed, the technique goes far beyond the modulation capabilities of silver halide emulsions. In digital

**Figure 4.187** Types of fringe produced by interference.

computing terms it would be possible to get more than $10^{13}$ bits of information onto a perfect amplitude-recording medium which is infinitely thin and a mere 10 x 12 cm² in area. Even greater information storage is possible if a perfect device with some finite thickness is considered.

In realistic terms it is now possible to record about $10^{10}$ bits on present-day emulsions, which goes a long way towards solving the problems of recording and comparing area texture on machined parts.

In many instances it may be found that the full fidelity of the holographic technique is not needed, and in such cases it is possible to reduce the information-recording capacity of the system to suit.

The detail that a hologram has to record is a function of the maximum angle $\theta$ between the reference beam and the signal beam. If $\delta$ is the average spacing of the fringe pattern and $\lambda$ is the wavelength then

$$\delta = \lambda / \theta \tag{4.199}$$

For example, if $\theta = 30°$ then the photographic plate will have to be capable of resolving about 1000 lines per millimetre.

The introduction of the laser was the key to the development of holography [138]. A typical arrangement for the recording and reconstruction of an object is shown in figure 4.188.

How the phase is retained is simply shown, as in figure 4.189. Consider the ordinary scattering of light from an object. The intensity of light hitting the photographic plate or film at, say, P is simply

$$I_P = A_1^2. \tag{4.200}$$

In this signal there is nothing to tell P from which direction the ray $r_1$ is coming.

Now consider the case when, in addition to the illuminating beam reflection $t_1$, there is superimposed another beam $r_2$ which is from the same source but which is directed onto the film and not the object. The intensity at P is now given by

$$I_p = A_1^2 + A_2^2 - 2A_1A_2\cos\theta \tag{4.201}$$

**Figure 4.188** Reconstruction of image.



**Figure 4.189** Difference between hologram and photograph.

Now because the source is coherent the signal received at P contains some information about the direction of the ray $r_1$ because the direction of the ray $r_2$ is known. It is this term which produces the 3D effect. Notice, however, that it only does so in the presence of the reference beam $r_2$. To see the object the configuration is as shown in figure 4.189(*c*).

So it is seen that holography is no more than two-stage photography. At the intermediate stage after the photograph has been taken in the presence of the reflected and reference beams (the hologram) no direct information is visible; the holographic plate looks a mess and no image is apparent. However, once the hologram has been 'demodulated' by the presence of the reference beam then the image is visible and as shown in figure 4.188. Two possibilities exist: the virtual image which cannot be projected and the real image which can.

The most popular view of holography is derived from the work done on the recording and reconstruction of the spatial elements of an object. A great deal of work has gone into the quantification of depth information

over comparatively large areas [139–144]. Probably the first attempt to measure roughness (as opposed to form) was due to Ribbens [145,146] using the system shown in figure 4.190.

In the figure $u_1$ is the reference beam and $u_2$ is the light from the surface. If the texture is small there is a high correlation between $u_1$ and $u_2$ and an interference pattern is formed. The point to make here is that the contrast ratio of the fringes is determined by the ratio of the two light intensities and the coherence between the light amplitude components $u_1$ and $u_2$. The latter is influenced by the spatial and temporal coherence of the incident light beam and by the roughness of the test surface. For a flat surface and a highly correlated source the coherence is determined very largely by the roughness. The roughness value in principle is found by measuring the fringe contrast ratio in the image plane.

Ribbens showed that under certain conditions the contrast ratio $R_c$ is given by

$$R_c = \left( \frac{\rho + \exp(-k^2 R_q^2 / 2)}{\rho - \exp(-k^2 R_q^2 / 2)} \right)^{1/2}$$

(4.202)

where $\rho$ is the ratio of light amplitude of $u_1$ to $u_2$ as determined by the reference beam and the hologram diffraction efficiency, and $k$ is determined by the reference beam intensity, $R_q$ is the RMS surface finish. He subsequently extended the method to the use of two wavelengths which theoretically extends the range of measurement [146]. Equation (4.202) then becomes

$$R_c = \left[ \left( \frac{(1 + \rho)^2 + (2\pi R_q / \lambda_{\text{eff}})^2}{(1 - \rho)^2 + (2\pi R_q / \lambda_{\text{eff}})^2} \right) \right]$$

(4.203)

where

$$\lambda_{\text{eff}} = \frac{\lambda_1 \lambda_2}{\lambda_1 - \lambda_2}.$$

(4.204)

Attractive as these methods may seem, they require many assumptions. The latter technique is a way of extending the range of measurement in theory without the need to use a laser having a very large $\lambda$ and introducing all sorts of complications of detector sensitivity, lens transmission problems and so on, and yet still produce an effective wavelength sufficiently large to allow a small phase modulation approximation of the surface to the wavefront.

Problems with these methods include the facts that a hologram has to be made, that the method is best for plane surfaces and, more importantly, that the system has to have a relatively large bandwidth. As in all small-signal phase modulation, the information lies in the sidebands (diffraction orders); it is necessary that



**Figure 4.190** Roughness measurement due to Ribbens.

the resolution of the hologram is sufficient to record higher diffraction orders from the surface. Ribbens suggests that as the method is meant for rough surfaces, the effective system bandwidth is relatively low, of the order of ~200 lines per millimetre—which means that most holographic plates can record enough diffraction orders to preserve the statistics of the surface at the image plane.

He points out that the system will also work with exposures of the plate for liquids of two different refractive indices rather than two wavelengths: both are entirely rigorous methods of changing the ratio of the wavelength of light to the surface roughness asperities.

In all the methods involving the reflection of light from surfaces some idea of the theory of light scattering is essential. This will be dealt with later.

Whereas the methods outlined above have been suggested as methods for measuring the amplitude of roughness, other possible uses exist. These are concerned with surface classification possibilities. Holograms have been used as cross-correlation filters for the automatic recognition of various spatial parameters [147, 148]. It may be that this approach could be used to recognize various machine-fault conditions. However, there would be enormous practical difficulty in implementing such a technique.

The most important application of holography is not essentially in surface metrology as such, but in strain and vibration measurement. The reason for this is profound. Therefore some explanation will be given here in simple terms.

Conventional interferometry involves fringe formation between optical wavefronts that have been divided either by wavefront or by amplitude. Both methods involve spatial separation of the beams and their subsequent recombination. Interferometry using holograms can be different. In ordinary interferometers the two beams are spatially separated at some point before recombination and one of the beams is modulated either by a movement of a surface or a refractive index change. In holographic interferometry the two beams can be separated in time. The modulation is caused by a change of a single spatial path with time.

Hence, holography allows the possibility of interferometry in time. In figure 4.191 temporal path 2 could be different from path 1 because of a spacing change between the object and the mirror.

The technique for strain measurement or vibration measurement using holographic interference (figure 4.192) is sometimes called double-pulse holography because the two exposures of the interferogram are taken with pulses a short time apart.



**Figure 4.191** Conventional (spatial) interference (*a*) and temporal interference (*b*).

**Figure 4.192** Holography for vibration for strain.

There is an extensive literature on this subject, so it will not be pursued here. However, some very interesting variants have emerged, such as that developed by Abramson [149], who uses double photographic plates for each exposure. With this extra plate a degree of freedom is made available which would not be otherwise. This is the effective directioning of the reference beam after the hologram has been made. But this 'sandwich' holography has problems with exposure times.

Holography has also been used as a means of making lenses and various other optical devices such as space frames in 3D measurement.

Basically holography has not been used effectively as a general purpose tool in surface measurement, but there are other possibilities open with coherent light, such as the use of the speckle phenomenon.

Holographic methods are therefore not specifically useful for surface roughness evaluation. They are useful in the same sense that interferometry is useful but no more. No instrument intended specifically for roughness measurement using holograms has been made commercially, nor is it likely to be. Surface texture is more important for its property of degrading holograms.

### 4.3.10  Speckle methods

One of the problems encountered with holography is the granular effect seen from diffuse objects when illuminated by coherent light. This granularity, called speckle, is now a useful metrological tool rather than the nuisance it was originally branded as.

It has been point out by Erf [150] that speckle phenomena were known long before the laser was invented but it has only been recently that applications have been found for it. There are two basic methods; direct laser speckle photography and speckle interferometry.

In principle there are three functions that control the nature of the observed speckle:

(1) coherence of the source;
(2) nature of the scattering medium;
(3) aperture of the imaging system.

Speckle has two very important properties. These are contrast and number of spots per unit area. The contrast can be shown to have a relationship to the correlation length of the surface, whereas the density is more concerned with the resolving capability of the imaging system.

Surface information has been obtained by using the contrast of speckle patterns produced in the first instance near to the image plane and, second, near to the diffraction or defocus plane.

Polychromatic speckle patterns have also been used, as have the correlation properties of two speckle patterns. From these methods roughness values from 0.01 to 25 $\mu$m have been estimated.

Simply put, the theory is as follows. The way in which speckle is viewed is shown in figure 4.193.



**Figure 4.193** Formation of speckle patterns.

The wavefronts reflected from the surface and picked up by the optical system will interfere with each other to create randomly varied speckles of various size throughout the space covered by the reflected light and projected onto the screen. The pattern is unique to the surface texture of the object, the illumination and viewing direction. The size of the speckles depends only on the angular extent over which the scattered light is received. Thus, if an imaging system is used to record the speckles, the resulting speckle size will be inversely proportional to the resolution of the system. The speckle size $S_p$ is

$$S_p = 1.2\lambda/(d/D) \tag{4.205}$$

for the surface far removed, and

$$S_p = 1.2\lambda F \tag{4.206}$$

where $F$ is

$$\text{aperture ratio} = \frac{\text{distance to image}}{\text{aperture diameter}}$$
$$= f \times (1 + m) \tag{4.207}$$

where $m$ is the magnification and $f$ is the $f$-number of the lens.

This is a gross simplification which demonstrates only the influence of the optical system. To get some idea of the influence of the surface, it is necessary to consider some more complex ideas.

Following here the introduction given by Jones and Wykes [151], figure 4.194 shows in principle what happens.

The path difference of many of the waves hitting P from the surface will vary by more than $\lambda$ because the surface roughness may be larger than this. Also, because the source of the rays could be from anywhere in the region of the illumination, there will be a massive summation of rays at any one point in the viewing

**Figure 4.194** Average speckle.

plane. Because of the randomness of the surface (usually implied) the way in which the summation shows itself as a resultant intensity will vary point by point on the viewing plane. Speckle is this variation in intensity in examining the intensity difference at $P_1$ or $P_2$ or anywhere else on the screen.

Goodman [152] has shown that the probability density of the intensity at a point P lying between $I$ and $I + dI$ is

$$p(I)dI = \frac{1}{\langle I \rangle} \exp \frac{-I}{\langle I \rangle}$$

(4.208)

where $\langle I \rangle$ is the intensity of the speckle pattern averaged over the received field—the variance of which depends on the number of scattering points within the illumination area—which depends critically on the correlation length of the surface. This obviously assumes a Poissonian probability within the interval.

Using this breakdown, the mean square value of the intensity is $2\langle I \rangle^2$ so that the standard deviation of the intensity is

$$\sigma = (\langle I^2 \rangle - \langle I \rangle^2)^{1/2} = \langle I \rangle$$

(4.209)

If the intensities of the scattered light at two points are compared, then if they are close the intensities will be correlated but if far apart they will not. The mean speckle size is related to the autocorrelation function of the intensity distribution $A(r_1, r_2)$.

It has been shown [118] that the correlation function of intensity pattern points separated by $\Delta x$ and $\Delta y$ in the observation plane is

$$A(\Delta x, \Delta y) = \langle I \rangle^2 \left[ 1 + \text{sinc}^2 \left( \frac{L\Delta x}{\lambda z} \right) \sin \left( \frac{L\Delta y}{\lambda z} \right) \right]$$

(4.210)

where $z$ is the distance between the object and viewing plane and $L \times L$ is the size of the illuminated object (assumed to be uniform). The average speckle is therefore $\Delta x$ (or $\Delta y$) for which the sine function becomes zero, yielding

$$Ax = \lambda z / L$$

(4.211)

Physically this corresponds to the explanation given by reference [151] (figure 4.195). The path length is

$$S = P_1 Q - P_2 Q \sim \frac{xL}{z} + \frac{1}{2}\frac{L^2}{z}$$

(4.212)

and the difference to $Q'$ is approximately

$$\frac{xL}{z} + \frac{1}{2}\frac{L^2}{z} + \frac{\Delta xL}{z}$$

(4.213)

The change from Q to Q′ is therefore $\Delta x L/z$ from which (4.211) results.

So from equation (4.211) as the area of illumination increases the 'size' of the speckle increases. This type of speckle is termed 'objective' because its scale depends only on the position in which it is viewed and not the imaging system used to view it. As a screen is moved from the surface, the speckle pattern will continually change.

Some mention of this will be made when fractal behaviour is being considered.



**Figure 4.195** Image plane speckle.

Image plane speckle is shown in figure 4.195. Here

$$QQ' = 1.22\lambda v/a. \tag{4.214}$$

The size of the speckle can be taken to be twice this ($2.4\lambda v/a$). The distance $P_1P_2$, which is the radius $r$ of the element of the object illuminating Q, is

$$\frac{1.22\lambda u}{a} \tag{4.215}$$

where $u$ is the object distance.

Goodman [152] has discussed an expression for the autocorrelation function of the image plane speckle;

$$C(r) = \langle I \rangle^2 \left[ 1 + 2J_1\left(\frac{\pi a r}{\lambda v}\right) \middle/ \left(\frac{\pi a r}{\lambda v}\right) \right] \tag{4.216}$$

from which the speckle size is taken when the separation between the two first minima of the Bessel function is measured.

$$d_{\text{speckle}} = 2.4\lambda v/a \tag{4.217}$$

The maximum spatial frequency is determined by the size of the lens aperture and $v$, the lens—image distance. This is called subjective speckle because it depends on the viewing apparatus.

In principle, some idea of the correlation length of the surface can be found by varying the aperture size and measuring the standard deviation of the intensity. If the aperture is small, there are fewer independent points on the surface which contribute to the speckle.

Following Asakura [153], the basic theory is as follows.

The speckle pattern formed at the image and diffraction planes is shown in figure 4.196. The object at $P_1$ is transmitted via the imaging system by a lens system having a spread function (the inverse transform of the OTF (optical transfer function)) given by $S(P_1P_2)$.

**Figure 4.196** Speckle pattern at image and diffraction planes.

The speckle intensity contrast is given by

$$I_{\mathrm{sp}}(p) = |A(P_2)|^2 = \left|\int S(P_1 P_2) A(P_1) \mathrm{d}p\right|^2 \qquad (4.218)$$

where $A(P_1)$ is the complex amplitude reflection at $P_1$ and $A(P_2)$ the complex amplitude at $P_2$. In general

$$A(P_1) = A\exp(j\varphi) \qquad (4.219)$$

so if the surface is smooth (and assuming that the amplitude factor is more or less constant with angle) the smoother the surface the smaller the variation in phase. This results in a reduction in the contrast $V$ of the speckle defined as

$$V = [\langle I_{\mathrm{sp}}^2(P_2)\rangle - \langle I(P_2)\rangle^2]^{1/2} / I_{sp}(P_2) \qquad (4.220)$$

where $\langle\ \rangle$ indicates an average value of the ensemble of scattering points.

Consider surfaces which are relatively smooth:

$$A(P_2) = a(P_2) + c(P_2) \qquad (4.221)$$

where $A(P_2)$ is the amplitude of the speckle corresponding to a certain point $P_2$ of the observation plane, $a(P_2)$ is the diffuse component of the speckle and $c(P_2)$ is the specular component.

Assume that Gaussian statistics apply, which in turn presupposes that the correlation length of the surface is small compared with the illuminated area (failure to follow this is treated in reference [154]). The feature that determines the contrast variation in the speckle pattern is the number of independent surface features within the illumination area, in much the same way that the variance of a statistical parameter is dependent on the number of observations. Letting

$$\begin{aligned}\mathrm{Re}(A(P_2)) &= a_r + c_r \\ \mathrm{Im}(A(P_2)) &= a_i + c_i\end{aligned} \qquad (4.222)$$

$$\begin{aligned}\sigma_r^2 &= \langle (A_r^2) - (A_r)^2\rangle = \langle a_r^2\rangle \\ \sigma_r^2 &= \langle (A_i^2) - (A_i)^2\rangle = \langle a_i^2\rangle\end{aligned} \qquad (4.223)$$

the average diffuse component $I_D$ and specular component $I_s$ are

$$\begin{aligned}I_D &= \langle |A(P_2)^2| - |A(p)|^2\rangle = \sigma_r^2 + \sigma_i^2 \\ I_s &= |\langle A(P_2)\rangle^2|^2 \qquad\qquad = c_r^2 + c_i^2\end{aligned} \qquad (4.224)$$

The probability density function of the speckle is obtained from the joint probability density of $A_r$ and $A_i$ and is assumed to be Gaussian:

$$p(A_r, A_i) = \frac{1}{2\pi\sigma_i\sigma_r} \exp\left[-\left(\frac{(A_r - C_r)^2}{2\sigma_r^2} + \frac{(A_i - C_i)^2}{2\sigma_i^2}\right)\right] \tag{4.225}$$

assuming the cross-correlation between $A_i$ and $A_r$ is zero, from equation (4.225).

The average contrast $V$ of speckle contrast is given by

$$V = [2(\sigma_r^4 + \sigma_i^4) + 4(c_r^2\sigma_r^2 + c_i^2\sigma_i^2)]^{1/2}\big/(\sigma_r^2 + \sigma_i^2 + c_r^2 + c_i^2). \tag{4.226}$$

Asakura [153] transforms equation (4.225) using $A_r = I^{1/2}\cos\psi$ and $A_i = I^{1/2}\ln\psi$, giving

$$p(I) = \frac{1}{4\pi\sigma_r\sigma_i}\int_0^{2\pi} \exp\left\{-\left[\left(\frac{\cos^2\psi}{2\sigma_r^2} + \frac{\sin^2\psi}{2\sigma_i^2}\right)I \right.\right.$$
$$\left.\left. -\left(\frac{c_r}{\sigma_r^2}\cos\psi + \frac{c_i}{\sigma_i^2}\sin\psi\right)I^{1/2} + \frac{c_r^2}{2\sigma_r^2} + \frac{c_i^2}{2\sigma_i^2}\right]\right\} \tag{4.227}$$

These equations represent those required for speckle investigation.

If the forms for $\sigma_r^2$ and $\sigma_i^2$ and $c_r$ and $c_i$ are known, the variations in speckle and the joint density can be found. As will be seen later the $p(I)$ term is useful in describing the nature of the surface spectral characteristics.

The important question is whether or not it is possible to extract surface finish information from the speckle pattern. There are two basic features that influence the speckle variability: one is the likely phase variations resulting from the amplitude of the roughness heights, and the other is the likely number of scatterers that could contribute to the pattern at any point. It has been highlighted by researchers what a critical role the optimal imaging system has in determining the actual number of independent surface features which make up the intensity at any one point in space [155–157]. It is obvious that the speckle intensity is considerably influenced by the number of scatterers contributing to the intensity. The amount of variation will be lower the larger the number of contributors. The surface roughness parameter which determines the unit scattering element is the correlation length of the surface defined previously. Obviously this is relatively straightforward to define and measure for a random surface but not easy for a periodic surface or, worse, for a mixture of both. The weakness of all the approaches has been the problem of attempting to estimate the roughness parameters without making global assumptions as to the nature of the statistics of the surface under investigation. Beckman and Spizzichino [158] started by suggesting a Gaussian correlation function within the framework of Kirchhoff's laws. This is not a good idea, as will be shown in the section on optical function. Probably the confusing influence of a short-wavelength filtering produced by the instrumental verification of the surface model prompted this. Far more comprehensive models should be used. However, this produces problems in mathematical evaluation [159].

Pedersen [156] used such complex models in a general form and showed that the speckle contrast is highly dependent on the surface height distribution, thereby tending to limit the usefulness of the method in practice.

Partially coherent light has been used with some success, in particular by Sprague [155] who found out that if the surface roughness height was comparable with the coherence length of the light—which implies quite a large bandwidth in fact—then the speckle pattern in the image plane near to the lens did bear some relationship to the roughness. Furthermore, a meaningful correlation was not confined to one process.

This work has been carried further with polychromatic light resulting in

$$P_1^2 = P_1(K_0)^2 / [1 + (2WR_q)^2]^{1/2} \tag{4.228}$$

and

$$P_1(K_0)^2 = 1 - \{1 + N^{-1}[\exp(R_q^2 K_0^2 - 1)]\}^{-2} \tag{4.229}$$

$R_q$ being the RMS roughness of the surface, $W$ the bandwidth of the light, $P_1(K_0)$ the contrast development of the speckle pattern for monochromatic light at mid-band. $N$ corresponds to the number of independent events (facets) within the illuminated area.

Other methods based on Sprague's work [155] have been attempted but so far do not appear to be very successful. This is because they are basically not versatile enough to be generally applicable. It is this that counts in instrumentation!

The results by Sprague have been queried by suggesting that perhaps the method has not been completely verified. Fundamentally the technique depends on plotting the contrast ratio as a function of $R_a$ (or $R_q$) or alternatively the number of effective scatterers (figures 4.197 and 4.198).



**Figure 4.197** Variation in speckle contrast with roughness.



**Figure 4.198** Speckle contrast as a function of reflectors.

The argument is that using these graphs the surface roughness ($R_q$ or $R_a$) can be found. Figure 4.198 will be explained later when very few scatters are considered or when phase effects which are large compared with the wavelength of light are present.

Leonhardt [159] introduces a number of variants, one of which is the use of a certain amount of incoherent background light together with some degree of defocus. He calls the methods white-light phase contrast detection and gives some results for different surfaces with varying degrees of added intensity. This is shown in figure 4.199 in which intensity is plotted on a log scale. $t$ is the amount of white light added. The basic idea behind this is that the ratio of constant light to scattered light is less favourable for a high degree of surface roughness, with the result that the measured contrast is the smallest, so the graph is steeper for larger values of $t$.

**Figure 4.199** Intensity dependence of white-light addition.

The whole theory is complicated but much effort is being used to try to explain the problems encountered in radar and similar transmissions through media which introduce both amplitude and phase variations. One of the best ways of understanding problem phenomena is to build a model for testing. Many workers have done this using either real surfaces in reflection as a model or a phase screen having a rough surface to produce the phase changes: from the observed scattering some conclusion can be expected to be obtained about the scatterer—the so-called 'inverse scatter problem'.

Almost all the theoretical treatments have assumed that the wavefront distortions (or equivalently the surface height fluctuations) constitute a joint Gaussian process, which implies that an autocorrelation function is needed for the evaluation of the complete model. Most workers have assumed a Gaussian model for the autocorrelation although this will be shown to be very questionable later in the book. This corresponds physically to a smoothly varying surface containing fluctuations of about the same size. When the surface height is about the wavelength of light or greater, scattered radiation from such a surface leads to amplitude fluctuations which become dominated by large-scale geometric effects. These effects are called 'caustics' or 'discontinuities' and correspond to 'local focusing' effects.

When such a target is illuminated by a laser beam of variable width the contrast of the scattered intensity pattern has the form shown in figure 4.200 if the contrast is taken far away from the scatterer. This corresponds to the far-field or Fraunhofer region which will be described shortly.

At large spot diameters (figure 4.200) the illuminated region on the surface is large compared with the correlation length of the surface and the contrast is unity (more or less uniform). This is called the Gaussian



**Figure 4.200** (*a*) Transmission; (*b*) reflection.

region. As the spot size is reduced, the number of independent scatterers reduces and more variations occur because the central limit theorem no longer applies and the contrast increases to very high values, sometimes in the region of 10.

A similar behaviour is observed at the other optical extreme of the near-field or Fresnel region where the wavefronts tend to be spherical rather than plane as in the Fraunhofer region. Here the independent parameter is the distance from the surface and not the spot size as shown in figure 4.201. Close to the surface there is very little contrast but further away the contrast increases owing to the caustics. Even further away the contrast relaxes to its Gaussian value of unity.

This is the sort of qualitative view that occurs at surfaces as a function of detector position and size of spot and is to be expected. However, sometimes this behaviour does not occur and a lot of thought has been given to why not!

One suggestion is that surfaces are not necessarily single scale; as indicated earlier they may be multi-scale. The reason behind this is that instead of getting the high contrast peak which is expected, the peak is relatively low which indicates that there is a superposition of different contrast patterns, each corresponding to a scale of size. This would tend to smear out any individual peak. Such surfaces would be characterized by a simple power-law spectrum and have been named 'fractal' surfaces after Mandelbrot [159]. A great deal of work has been carried out to examine the so-called non-Gaussian and Gaussian behaviour of such surfaces as well as for ordinary surfaces, in particular by Jakeman and McWhirter [160], Fuji [161], Levine and Dainty [456], etc.

Such surfaces have different characteristics to the single-scale ones. Single-scale surfaces are continuous and differentiable to all orders, which implies that the autocorrelation function is defined at the origin and can be approximated by Taylor's theorem as indicated earlier. Fractal surfaces are continuous but not differentiable and contain structure down to arbitrary small scales.

The use of simple rays to explain the behaviour of these surfaces is definitely inappropriate. Take as an example a surface used by Jordan [162] having a structure like a multilevel telegraphic signal as shown in figure 4.202. This type of surface, although not strictly a fractal surface, has some of its properties. All the rays from the surface (if a phase screen) are parallel and so can form no geometrical features in the scattered radiation (neither can a fractal surface).

At this stage it should be noted that it is very likely that surfaces could exhibit some degree of fractal behaviour, especially in those formed by growth or deposit mechanisms. This has been examined in Chapter 2.

It has been shown that the surface can be determined by the unit event making up the surface (e.g. the grid), its impact height distribution and the spatial distribution of bits on the surface.

For a square grid of random height, width and position

$$A(\beta) = \sum_{i=1}^{4} P_i(\beta)\exp(\alpha_i\beta)$$

(4.230)

where $P_i(\beta)$ are finite-order polynomials in $\beta$ (the lag) and $\alpha_i$ are positive constants which depend linearly on the width of the impressions $\mu$ and the density of events per unit length $\lambda$.



**Figure 4.201**

**Figure 4.202** Facet model for surface zero-order Markov.

The general form for triangular grids having various heights, positions and widths is more complex, being of the form

$$A(\beta) = g_i(\beta)\beta \exp(\alpha_i \beta) \int_{-a}^{a} \psi(z\beta) z \exp(-z^2) \mathrm{d}z$$

(4.231)

where $g_i$, are finite polynomals in $\beta$, and $\alpha_i$ are constants involving $\lambda$ and $\mu$. $\psi(z\beta)$ is a function of $z$, the surface height.

Both (4.230) and (4.231) show the dependence on the exponential form, which means that most random surfaces exhibit some form of Markov characteristic and so would be expected to be multiscale. The assumptions made by Beckman regarding Gaussian correlation functions for manufactured surfaces have not been held by surface workers not working in optics.

In concluding this section on speckle phenomena of surfaces it should be noted that a lot of work has been carried out on interferograms, the so-called speckle interferometry pioneered by Butters and Leendeutz [163] particularly for vibration and form assessment. Surface roughness has not been a priority.

As a general rule at present it seems that some use of speckle for surfaces could be obtained for monochromatic light of about the wavelength of light or less. This can be extended by using polychromatic light or, for example, by using a type of comparison between 'standard' surfaces and 'generic' surface.

Summarizing the speckle approach [153], the coherent light speckle contrast method is useful for estimating relatively fine surfaces of less than 0.25 $\mu$m $R_q$. Polychromatic speckle pattern methods apply more to surfaces having a roughness of 0.2-5 $\mu$m $R_q$. Speckle pattern correlation is applicable to surfaces which are roughly, say, 1-30 $\mu$m $R_q$.

Despite the agreement of theory and practice obtained using speckle methods it should be pointed out that the speckle method is at best an indirect method. It is one form of the 'inversion' problem, estimating surface properties by means of the properties of the light scattered from the surface. There are many unwanted variables. One is the receiving optics, another is the illumination and a third is the nature of the assumptions about the surface which have to be made before the inversion can be made and the value of the surface found. The fact that there is no widespread use in industry for this method testifies to its limited acceptance.

In this section the speckle behaviour of the light scattered from rough surfaces has been investigated. In particular, probability distributions of intensities and correlation functions have been estimated from the speckle. There is another possibility, namely light scatter itself. This is the angular scattering properties of light incident on a surface. Light scatter has received much attention and in some respects at least it is superior to the speckle method.

### 4.3.11 *Diffraction methods*

Probably the most potentially useful application of the general scattering properties of light are to be found in diffraction methods.

Because of its relevance to all scattering, the basic philosophy of the scattering of light from a surface will be considered. This subject has been reviewed in detail by Ogilvy [164] but not from an engineering point of view. There are two ways of tackling this problem in electromagnetic theory, one scalar and one vector [165] and compared in [166, 164]. Easily the most used has been the scalar theory based on the original work of Beckmann and Spizzichino [158] and modified by Chandley [167] and Welford [168] to suit rough engineering surfaces and not radar reflections from the sea.

The scalar theory is based upon the Helmholtz integral and the vector variational techniques developed originally by Rayleigh.

Consider first the scalar theory as clarified by Chandley and Welford. Scattering in the optical or near optical region presents problems because of relating the properties of the actual scattering surface to that of the wavefront emanating from it. Problems of this nature have been encountered in contact and related phenomena where simplistic models of surfaces have been devised in order to obtain tractable results—which in the event turned out to lack credibility in the mechanical sense.

The problem is twofold:

1. Determining the relationship between the wavefront scattered in the region of the surface and the properties of the surface itself. How do the statistical properties of the wavefront compare with those of the surface?
2. Determining the properties of the intensity of the scattered light in the far field.

The theoretical approach is usually based on the Helmholtz integral. This is normally solved under Kirchhoff boundary conditions, which assume conducting surfaces, no reflections or shadows to give the following expression in the far field:

$$I(K_2) = \frac{-j\varepsilon_\mathrm{o}R}{2\lambda f} \iint_s (K_1 - K_2).n \, \exp[j(K_1 - K_2).K] \, \mathrm{d}s \qquad (4.232)$$

where $I(K_2)$ is the amplitude of the field in the direction $K_2$, $\varepsilon_\mathrm{o}$ is the amplitude of the incident wave, $R$ is a reflection function assumed to be constant (it should be valid for scalar forms but not for vector (polarized) forms)), $\lambda$ is the wavelength, $f$ is the focal length of the lens and $S$ is the area covered. Figure 4.203 shows the system and figure 4.204 shows the angular convention.

For a stationary random process, also assumed to be ergodic, Beckmann's solution of (4.232) is a rather messy formula for the contrast ratio in the Fraunhofer plane:

$$\frac{\langle I^2 \rangle}{\langle I \rangle^2} = A\exp(-g) + \frac{F^2}{2L} \int_{-L}^{L} \exp(jv_x\beta)\exp(-g)\exp[gA(\beta) - 1] \, \mathrm{d}\beta \qquad (4.233)$$



**Figure 4.203** Coordinate system.

**Figure 4.204** Angular system.

$\langle I^2 \rangle$ is averaged between—$L$ and $+ L$ (the part of the surface illuminated) and standardized to the value of the intensity of the field reflected along the specular direction. $A$ is $\sin^2 v_x L$, the spread function of the lens, $g$ is

$$(v_z R_2)^2 = \left( \frac{2\pi}{\lambda} (\cos\theta_1 + \cos\theta_2) R_q \right)^2 = \left( \frac{4\pi R_q}{\lambda} \right)^2$$

$\theta_1$ is the incident angle and $\theta_2$ is the scattering angle

$$F = \sec\theta_1 \left( \frac{1 + \cos(\theta_1 + \theta_2)}{\cos\theta_1 + \cos\theta_2} \right)$$

(4.234)

and $R_q$ is the RMS value of the surface.

The terms of (4.233) are rather critical. The first term, $A \exp(—g)$, represents the specular term of reflection and the latter term is the diffuse light.

From here on many different variants exist as how to work with the formulae: many investigators get into problems with assumptions. Beckmann assumes a Gaussian correlation function which is questionable but he does attack the basic vital problem of trying to relate the real surface to that of the wavefront close to the surface. Chandley [167] and Welford [168] give a review of various optical methods in which a good many of the original papers in the subject are reviewed. Earlier workers (e.g. Bennett and Porteus [169], Davies [170] and Torrence [171]) had similar equations.

Thus the intensity for normal incidence light was derived as

$$R = R_o \{\exp[-(4\pi R_q \lambda)^2]\} + R_o (2^5 \pi^4)/m^2](R_q/\lambda)^4 (\Delta\theta)^2$$

(4.235)

where $R_q$ is the surface RMS value, $m$ is the RMS slope, $R_o$ is the reflectance for a perfect surface and $\Delta\theta$ is the acceptance angle of the detector.

The first part agrees with Beckman and is the specular component; the latter term expressed in rather different form is total scatter around the normal. This is modified if non-normal angles are used. Thus

$$\frac{R}{R_o} = \exp\left[ -\left( \frac{4\pi R_q}{\lambda} \right)^2 \cos^2\theta \right] + \frac{2^5 \pi^3}{m^2} \left( \frac{R_q}{\lambda} \right)^4 \cos^3\theta (\pi(\Delta\theta)^2).$$

(4.236)

Obviously this is a useful expression for deriving the slope and $R_q$ value for the surface. However, other statistical parameters such as the power spectral density and the autocorrelation function are much sought after and methods for getting them need to be employed urgently because of their usefulness in machine tool monitoring. The fact that light scatter provides a tool which is non-contacting and can work at the speed of light is attractive. It is also very important when the workpiece is remote and access is difficult—the surface of the moon, for example, or ocean waves.

It is still probable that estimates of surface finish for values less than the wavelength of light will give the most productive results. Chandley [167] has modified (4.233) and derived equations for the autocorrelation

function of the surface in terms of the inverse Fourier transform of the scattered light spectrum in the Fraunhofer plane $A_F$ the wavefront correlation yielding

$$A(\tau) = 1 + \frac{\ln\{(A_F(\tau)/A_F(0))[1 - \exp(-R_q^2)] + \exp(-R_q^2)\}}{R_q^2}.$$

(4.237)

However, such a technique, although promising on paper, has not been developed because of the large number of imponderables involved. Investigators are still trying! Robust assumptions are needed!

The best method seems to rely on simply using the weak scatterer theory where $R_q \ll \lambda$. With this there is a definite possibility of getting valid information about the surface under certain well-defined conditions.

In general, the wavefront emanating from the surface will be

$$w(x, y) = A(x, y)\exp(j\varphi(x, y))$$

(4.238)

where $A$ is an amplitude term representing local changes in reflectance of the surface (as $R_o$) and $\varphi(x, y)$ represents the changes in phase produced by the surface heights.

There is no conflict between equation (4.235) and (4.236). They are both expressing the same physical phenomena. The equivalence rests upon the fact that there is a relationship between the correlation and the slope for a given distribution. For example, for a Gaussian distribution it has been shown earlier that the relationship is

$$m = R_q \sqrt{\frac{\pi}{1-p}}.$$

(4.239)

Bennett and Porteus use the relationship

$$m = \frac{R_q\sqrt{2}}{\tau_L}$$

which is dangerous because it assumes a well-behaved autocorrelation function. The true mean slope is likely to be very much higher in practice, for example in any type of fractal or multiscale surface. The scale of size under consideration is always fundamental.

Insertion of this type of expression from one to another and remembering that Bennett and Porteus only considered scatter in the immediate region of the specular reflection (which corresponds to a value of $n = 1$ in Beckmann's formula) with a specific surface model (section 5.3 in reference [158]) gives only a crude idea of what happens.

In fact slopes become important in the case when the surface is rough ($> \lambda$ as indicated by the formation of local caustics referred to earlier). Straightforward assumptions give some indication of how this could be. Assuming that the surface height distribution and its correlation are Gaussian (which is not a good assumption) and letting the surface slope be $\psi$ it will have a distribution

$$p(\psi) = \frac{1}{\sqrt{\pi}\tan\beta_o\cos^2\psi}\exp\left(\frac{-\tan^2\psi}{\tan^2\beta_o}\right)$$

(4.240)

where $\tan\beta_o = 2R_q/T$.

Beckmann's equation (54) in section 5.3 of reference [158] reduces to

$$\frac{\langle I^2\rangle}{\langle I\rangle^2} = \frac{kT}{\sqrt{g}}\exp\left(\frac{-v_x^2 T^2}{4g}\right)$$

(4.241)

which can be written as

$$\frac{2k}{v_z \tan\beta_o} \exp\left(\frac{\tan^2\beta}{\tan^2\beta_o}\right)$$

(4.242)

which has a similar form to the variation of probability distribution of the slopes. Replacing the factor $F$ from before gives

$$I_\beta I_o = F p(\beta).$$

(4.243)

Thus the diffusion direction does give a measure of slope and plotting the intensity variation with $\beta$ gives the surface slope distribution directly. This is to be expected because as $g$ increases the theory has to converge to a geometrical solution. As $g$ decreases the angular spread of the diffuse light reduces by the factor $[(g-1)/g]^{1/2}$ and increases the maximum intensity by the inverse ratio. These two factors are essential in the understanding of the operation of scatterometers. Hence there is no basic confusion between the various theoretical equations; it is simply that the slope dependence is more clearly seen for the rougher surfaces [170].

Comparative methods for extending the range of measurements from the fine surfaces to the rougher ones have been attempted using Beckmann as a guide with some success [172], the effective range being extended to 1.8 $\mu$m by means of breaking down the surface types into two components, random and periodic. Unfortunately the usual constraints with scatter methods of not being absolute and needing careful preparation still apply. This type of approach in which some reliance is made on theory as well as practice has its place because many of the conventional industrial problems fall within the roughness range of interest. These include pistons, pins, journals, crank pins, etc, as well as very fine metallic mirrors for high-power lasers.

For normal incidence

$$\varphi(n, y) = \frac{4\pi}{\lambda} z(x, y) + 2\pi n$$

(4.244)

where $n$ is an integer. For unambiguous results obviously (in the weak scattering regime)

$$\varphi(n, y) = \frac{4\pi}{\lambda} z_{max}(x, y) \leqslant 2\pi$$

(4.245)

where $z(x, y)$ are the 2D surface ordinate heights and $z_{max}$ is a maximum value.

Now if the wavefront is bounded by an aperture $S$ and the far field (or Fraunhofer diffraction pattern) is formed by a lens of focal length $f$, by applying Kirchoff diffraction theory the complex amplitude is given by

$$F(v, w) = \frac{1}{S} \int\int_{-\infty}^{\infty} w(x, y) \exp\left(-\frac{2\pi j}{\lambda f}(vx + wy)\right) dx\, dy.$$

(4.246)

This can be physically appreciated better if the simple idea is adopted that the exponential term under the integral represents in effect the collecting of plane waves at constant spatial angles $D/x$ and $w/\lambda$, the contribution at the angle being determined from $W(x, y)$. The notation is similar to that used in electrical filtering for linear phase situations mentioned earlier in data processing. In passing, the Fresnel equivalent for equation (4.246) contains exponential squared terms which, when put under the integral sign, represent the collecting of spherical rather than plane wavefronts.

Even equation (4.246) contains a number of assumptions that are usually satisfied in practice and will not be introduced here.

Equations (4.244) and (4.245) are modified if an oblique angle rather than normal incidence is used. Thus if $\theta$ is the angle from the normal

$$\varphi(x, y) = \frac{4\pi}{\lambda} z(x, y)\cos\theta.$$

(4.247)

Note that the cosine term has the effect of making the path difference between rays hitting the top of a peak and the bottom of an adjacent valley smaller with respect to $\lambda$. Alternatively, the wavelength $\lambda$ could be considered to be increased by sec $\theta$.

In equation (4.246), the wavefront can be expressed as in (4.248). Thus

$$F(v, w) = \iint A(x, y)\exp\left( j\frac{4\pi}{\lambda} z(x, y)\cos\theta \right)\exp\left[ -\left( \frac{2\pi j}{\lambda f}(vx + wy) \right) \right]dx\,dy.$$

(4.248)

A number of simplifying assumptions are usually made here. One of these concerns the nature of $A(x, y)$. This is usually taken as a constant over the surface and sometimes even given the value unity. Obviously if only comparative methods are being used, the constant can be removed.

A second assumption concerns the nature of the exponential phase term. If the maximum phase advance is smaller than 0.5, the function $w(x, y)$ can be approximated by

$$w(x, y) = 1 + j\varphi(x, y, \theta) - \varphi^2(x, y, \theta)/2$$

(4.249)

by taking the first three terms of the expansion. From equation (4.249) it follows that the phase term of the wavefront is closely related to the Fourier transform of the surface heights themselves, $z(x, y)$.

If the angle $\theta = 0$, it enables the possibility of simple practical instrumentation to be made—which is one problem addressed in this book.

Therefore from (4.249) and (4.248) with the assumptions included

$$F(v, w) = \iint_s \exp\left[ -\left( \frac{2\pi j}{\lambda f}(vx + wy)' \right) \right]dx\,dy + j\frac{4\pi}{\lambda}\iint z(x, y)\exp\left[ -\left( \frac{2\pi}{\lambda f}(vx + wy) \right) \right]dx\,dy$$
$$- \frac{8\pi^2}{\lambda^2}\iint z^2(x, y)\exp\left( -\frac{2\pi j}{\lambda f}(vx + wy) \right)dx\,dy$$

(4.250)

Obviously this is even more simple if the last term is ignored, which can only occur if $z(x, y) \ll \lambda$.

In practice it is the intensity of light which is measured in the diffraction plane rather than the complex amplitude. This is given by

$$I(v, w) = F(v, w)F^*(v, w) = |F(v, w)|^2$$

(4.251)

where $F^*(v, w)$ is the complex conjugate of $F(v, w)$.

Putting equation (4.250) into (4.251) yields the equation for the intensity in the Fourier diffraction plane (following Rakel's analysis [173]). Thus simplifying the nomenclature and dropping insignificant terms gives, if $\varphi$ and $\varphi(x, y)$ are different phase points on the wavefront,

$$I(v, w) \propto F\left( \frac{1}{S}\iint_s (1 - \tfrac{1}{2}\varphi^2 - \tfrac{1}{2}\varphi^2(x, y) + \varphi\varphi(x, y))dx\,dy \right)$$

(4.252)

where

$$\varphi(x, y) = \varphi(x - x_1, y - y_1)$$

(4.253)

and $F$ denotes taking the Fourier transform.

The relationships between the intensity and the roughness values can be obtained by inserting different types of surface model in equation (4.252).

The essential problem in the whole exercise is that of knowing how many realistic approximations to make, and also knowing what surfaces will be measurable under these assumptions.

As an example of how the technique can be used in an instrument using a range rather larger than that which would be expected. Rakels has devised algorithms which increase the range of measurement by a factor of 4 for the special case of periodic surfaces [173]. This involves an expansion of the exponential wavefront to up to 10 terms in the series. Thus

$$\exp(j\varphi) = 1 + j\varphi + \sum_{2}^{10} \left(\frac{j\varphi}{n!}\right)^n$$

(4.254)

and

$$z(x) = \sum_{n=0}^{m} a_{2n+1} \sin[(2n + 1)x]$$

(4.255)

where (4.255) is chosen to represent various periodic waves in, for simplicity, one dimension.

Insertion of (4.254) and (4.255) into (4.246) enables the amplitude orders in the intensity pattern to be evaluated. So if the diffraction pattern is $I(\omega)$ it can be presented by

$$I(\omega) = \sum_{n=-m}^{m} A(n)\delta(C\omega + n)$$

(4.256)

where $m$ represent order numbers. $C$ is a constant.

For the simplest case when $z(x)$ is sinusoidal

$$A(m) = CJ_m^2(R_q/\lambda)$$

(4.257)

where $J_m(R_q/\lambda)$ is the Bessel function of the first kind of order $m$. By normalizing the amplitudes $a_m$ by

$$a_m = A_m \Big/ \sum_{n-\infty}^{\infty} A(n)$$

(4.258)

the normalized maxima are obtained.

By moving between approximations represented by (4.249) to (4.255), it is possible to develop various algorithms involving the maxima to enable a considerable increase in the range to be made. This is possible only if the type of process and the amplitude range are restricted. However, it shows that by balancing between optics and computing a gain can be made.

This approach, even for periodic surfaces, is not restricted necessarily to single-point machining because the results work for thin wheel grinding or cases where the dressing marks on the grinding wheel are left on the surface.

There are other ways of extending the range of measurement using diffraction methods. Perhaps the most obvious is to increase the value of the wavelength of radiation thereby reducing the ratio of $R_q/\lambda$. This method will work [174], but there is a limit to the application of such methods. The detectors and the optics become expensive outside the preferred range 0.5-2 $\mu$m.

An alternative yet equivalent method is to increase the wavelength artificially. There are a number of ways of doing this. One has already been suggested—that of using oblique incidence light. TIS and gloss meter instruments rely on the fact that the roughness $R_a$ is very small when compared with the wavelength of

light. The typical acceptable ratio of $R_a/\lambda$ is less than $^1/_8$ although this can be increased by a factor of two or three knowing the form of the surface. See Rakel's analysis [173].

At extreme glancing angles the wavelength $\lambda$ becomes effectively lengthened to $\lambda \operatorname{cosec} \theta$. The path difference between rays hitting peaks and those hitting valleys becomes small compared with the effective wavelength. A very cunning development is shown schematically in figure 4.206.



**Figure 4.205** Glancing angle mode.

Figure 4.207 shows a practical embodiment by Zygo. In this configuration the incoming beam is diffracted by the coarse grating. The zero order (direct beam) is blocked off. There remain two first order beams which are projected by means of the fine grating onto the surface at different angles. These two different angles produce two effective wavelengths at the surface which interfere with each other, producing fringes which are viewed (using reciprocal optics) by a camera..

Using such a system enables effective wavelengths of 12.5 $\mu$m to be produced from a laser having a wavelength of only 0.68 $\mu$m. Any optical instrument utilizing oblique rays has some constraints on the shape of the object, e.g. concave objects are difficult to measure.

The obvious way to measure rough surfaces is to use a longer wavelength source such as infra-red having a wavelength near to 10 $\mu$m rather than 1$\mu$m. the scheme shown in figure 4.208 shows a typical system (Zygo corp.) which incorporates separate interferometers for moving the stage. In such systems it is not always possible to use conventional optics because of absorption of the infra-red by the lenses. Special glass is used. It is interesting to note that systems using grazing incident light to simulate long wavelengths do not have this problem.



**Figure 4.206** Schematic of desensitized interferometer [175].

**Figure 4.207** Zygo embodiment.



**Figure 4.208** Infra-red interferometer.

Infra-red (IR) techniques have also been used in the 'total integrated scattering' mode discussed in section 4.3.12. Because this is an integrated method it is fast and has been reported as being effective for measuring surfaces up to $1.7\mu m$. Also using IR makes the method insensitive to ambient light; vibration effects are minimized, again because of the integration [176]. Another method is to change the refractive index of the medium in between the surface and the detector. Note that oblique incident light and the Fraunhofer diffraction mode have been used not only to increase the range over which average surface characteristics can be estimated but also for examining surfaces for flaws [177]. This method of looking at flaws seems to be questionable in view of the obvious problem of separating out single and multiple flaws of different characters in the same field. Spatial filtering methods have worked in which the reconstructed image gives an indication of where the problem areas are likely to be, but viewing in the image plane at least for the rare flaw seems to be a possibility. This highlights the problem of trying to achieve two seemingly disparate requirements: getting average statistical information over an area and finding very small but significant scratches or flaws. Here the space-frequency functions are a distinct possibility, theoretically lying in between the image and diffraction plane. The Wigner function has been mentioned earlier in the characterization of surfaces. The next step is to find an optical way of achieving it. This does not seem to be an impossible task because the function is at least real.

In the angular method of increasing the effective wavelength multiple reflections can occur and also shadowing. How these affect the estimate of roughness will be seen in Chapter 7 in optical function. The

other method by Whitehouse [178] involves the use of replicas and was the first practical use of optical Fourier spectrum analysis for surfaces, although replicas are still being used. One example is in the use of replicas in holographic contouring [179]. Some reasonable results can be obtained but it is a messy procedure.

In this a replica is taken of the surface using the method with film, not resin. This was shown in figure 4.209.

If the surface height is represented by $z(x, y)$ then the path difference between rays touching the valleys and those the highest peaks is $(n_1 - n_2)z(x, y)$. This produces a phase advance of

$$\frac{2\pi}{\lambda}(n_1 - n_2)z(x, y)$$

(4.259)

Note that the wavelength A has effectively been increased by $1/(n_1 - n_2)$ where $n_1$ is the refractive index of the replication material and $n_2$ that of the 'matching' fluid. This multiplication factor can be made arbitrarily large simply by making $n_2$ near to $n_1$. Notice, however, that a factor of 2 is lost in equation (4.259) because the method is based on transmission and not reflection.

There are obviously a number of major drawbacks to this method. The most severe is that a replica has to be taken. This brings with it the possibility of a loss of fidelity. Also the plane face of the replica has to be smooth and flat (also true for the faces of the cell containing the liquid). Other problems include the homogeneity (or lack of it) of the replica material.

Nevertheless, the method has its uses, especially for measuring spectra in bores or on other inaccessible surfaces.

Diffraction can also be used in another way (figure 4.210). In this the surface is floodlit by UV light projected at an angle, thereby deliberately forming shadows. This is imaged by the objective lens onto a film which is photochromic in the sense that it is sensitive to UV light. The film is opaque where the bright and dark shadow pattern of the surface is imaged on it. This shadow pattern is transformed into the spectrum by the transform lens. This method approximates the shadow pattern to the true surface in much the same way as a zero-crossing model of the surface. At present the relaxation time of the film is rather large so that the method is limited. It does, however, change from a reflected transform to a transparent one and is genuinely useful for a very wide range of surfaces. It, in effect, increases the range of diffraction methods but only at



**Figure 4.209** Power spectrum of surface using replica.

**Figure 4.210** Diffraction method for spectral analysis using ultraviolet light.

the expense of being an indirect method. Care has to be taken with the UV light because it has to be high in energy (of the order of watts) because of the absorption by the surface.

Fraunhofer diffraction methods have been used in many practical situations, for example in turning [180] and rolling [181] by Konczakowski, with useful results. These definitely show that the technique can be used not only for estimating the spectra of the surfaces but also for machine process diagnostics, a point found out earlier by a CIRP team using correlation methods directly [182].

It is not surprising that these methods are becoming popular. This is for a number of reasons as listed below:

1. The information received in the transform plane is not unduly influenced by the movement of the workpiece. This is shown in figure 4.211.



**Figure 4.211** Coordinate invariance of transform system.

Movements in the $x$, $y$ plane do not affect the position of the spectrum (unless the part is moving at the speed of light) and so make the technique very useful for in-progress gauging where the part would normally be moving rapidly during machining. Movement away in the $z$ direction has only a small effect because the included angle of the illuminated area as seen from the detector is only slightly less with the shift, so that higher orders could possibly get blocked off at the limited detector aperture. Obviously as the part moves there will be small variations in the actual magnitude of the signal in the transform plane due to the statistical variations between different parts of the surface. These can be integrated out in time.

The situation is slightly different for shapes other than flat ones, such as the cylindrical workpiece shown in figure 4.211($b$). Here the freedom is less; yaw is not allowed, neither is movement in the $x$ direction. Usually the illumination is changed to suit the part shape; in the cylindrical case a slit would be used thereby axially illuminating the surface. Also the illumination of the part should match its shape if possible; for a spherical object there should be a spherical wavefront hitting the surface preferably concentric to the part.

2. Small features are easier to measure than large ones. This is shown in figure 4.212.
   The relationship between $d$ and $D$ is elementary and yet rather profound. Thus

$$D = \lambda f/d. \tag{4.260}$$

The smaller $d$ is the larger is the value of $D$, and hence in terms of sensitivity on the detector it is better. This is a metrological 'freak', because usually the smaller the feature the harder it is to measure!

This is a fortunate phenomenon for the surface metrologist because there is a definite trend to finer and finer surfaces and tighter shape and dimensional tolerances.

3. The method is non-contacting and therefore has certain advantages over stylus methods.

For in-process gauging there is another incidental advantage that follows directly from the first advantage listed above. This concerns the sensitivity of the method to debris and coolant. In the first case there are positive benefits in using this technique because the value of intensity at any point in the diffraction (transform) plane is the result of all light scattered at one angle over the whole illuminated area. A chip in the scattered field only reduces the overall intensity of the whole pattern because it will obstruct (to about the same extent) the light scattered at all angles. Providing that the chip does not completely block out all the light then the information is hardly affected.

The presence of coolant is a problem because it affects the amount of light hitting the surface. The only way out is to use a laser wavelength which has small losses in coolant. This is acceptable if the transmission of the coolant does not change with time, for example it biodegrades.

Alternatively an air knife (compressed air) can be used to clear the surface just in front of the illuminated area. This is effective providing it does not produce a mist, in which case the signal does get



$$D = \lambda f/d$$

**Figure 4.212** Reciprocal relationship between object and transform spacings.

degraded. Some investigators [183] attempt to correlates diffraction results with conventional methods. Their results have to be queried because high correlation is not significant due to surfaces having a common traverse scale.

There is another problem that is associated with the interpretation of the signal in the transform plane and this is concerned with the presence of waviness. The basic problem is that of stability; the illumination is never sufficiently large to get a sufficient sample of waviness even if it does for roughness. If the waviness is very much larger than the roughness the whole centroid of the transform can be shifted or split, in much the same way that a blazed grating throws the maximum intensity pattern off centre. The effect is shown in figure 4.213.

This fragmenting can be difficult to unravel if there is an overlap. Obviously if there were a large sample of waviness it would produce a well-behaved component of the transform very close to the centre of the transform pattern.

Another effect is the superimposition of roughness marks which are random and those which are periodic. These need to be separated. Algorithms which are designed to give roughness from periodic surfaces often fail in the presence of randomness, a well-known signal-to-noise problem. The extraction of the most useful information under these circumstances is complex and involves a number of tricks such as deconvolution of the spread function corresponding to the aperture.

Various rules have been tried for different types of surface [240]. These are not valid for peak parameters.

There are other ways in which to develop theories about light scatter from surfaces. One obvious way is using vector theory (polarization) rather than scalar theory. However, this method has not been adopted in practice mainly because of its difficulty and because of its sensitivity to surface films and crystalline structure. Vector methods usually have only worked on prepared surfaces. Some notable examples have been tried, for instance by Church *et al* [184] and Elson and Bennett [185].

Basically using only the first term of the scattered field $f_r$ the BRDF (bidirectional-reflectance distribution function) is given by

$$f_r^{ss} = \frac{1}{P_o} \frac{dP^{ss}}{d\Omega} = \frac{1}{\pi^2} \left( \frac{2\pi}{\lambda} \right)^4 F_{ss} F(\omega)$$

$$f_r^{pp} = \frac{1}{P_o} \frac{dP^{pp}}{d\Omega} = \frac{1}{\pi^2} \left( \frac{2\pi}{\lambda} \right)^4 F_{pp} F(\omega)$$

$$(4.261)$$



**Figure 4.213** Identification of form error in transform domain.

where $P_o$ is the incident power, $P$ is the scattered power, $\Omega$ is the solid angle, $\lambda$ is the wavelength of the radiation, $F$ is given below, $F(\omega)$ is the surface power spectral density and $\omega = (2\pi/\lambda)$ (sin $\theta_s$ – sin $\theta_i$) is the surface scattered spatial frequency, where $\theta_i$ is the incident angle; ss means that the source and receiver are both polarized with the electric vector normal to the plane of incidence, and pp means that they are polarized parallel to the plane of incidence.

According to (4.261) the vector theory predicts that

$$P_{ss} = \cos\theta_i\cos\theta_s$$

$$P_{pp} = f_{ss}\left(\frac{1-\sin\theta_i\sin\theta_s}{\cos\theta_i\cos\theta_s}\right)^2$$

(4.262)

whereas scalar theory predicts [143]

$$P_{ss} = (\cos\theta_i\cos\theta_s)^4/\cos\theta_s.$$

(4.263)

The total flux according to reference [186] is given by

$$f_r = \tfrac{1}{2}(f_r^{ss} + f_r^{ss}).$$

(4.264)

The same result is obtained from scalar theory because scalar theory treats all polarizations equally so that the BRDF is the average of the polarizations. For in-plane measurements therefore the result is the same as equation (4.264). For out-of-plane measurements the scalar result is half that of the vector.

Comparisons between the vector and scalar theory and practice have been carried out by Wang and Wolfe [186] with exhaustive tests on a single specimen. Their results indicate that there is still a long way to go before a true correlation between surface characteristics and the scattered polarization effects is properly understood. However, it seems obvious that much more information is potentially available than by using scalar methods. On the other hand this extra information may well be a confusing mixture of surface and material properties such as the dielectric constant. More work in this difficult area is needed. More precise results have been obtained using vector theory for the scattering of light from gratings [187].

Using the same criterion as that used for speckle, have diffraction or scalar methods resulted in practical commercial instruments? The answer is a definite, but belated, yes. There is evidence that machine tool manufacturers are considering incorporating diffractometers into the machine tool system principally because of their in-process capability. Using the diffractometer as the front-line instrument for surface control and the stylus instrument as a more accurate quality control check is one way to go.

Of the two methods, scalar light scatter and speckle, the former is the most practical because it requires fewer assumptions to be made and it is less dependent on the instrument system parameters. So again the choice is between getting less rigorous but useful results rather than complicated and often impractical results!

### 4.3.12 Scatterometers (glossmeters)

There are three relatively simple ways of getting limited surface information. These are: examining the way in which the specularly reflected light behaves as a function of angle, measuring the total integrated scatter (TIS) and measuring at a selected number of angles.

Take first the specular reflection $I_r(s)$ where the angle of the detector is taken to be the same as that of the source as in figure 4.218(a) when $\theta_i = \theta_2$: $K$ is a constant

$$I_r(s) = K\exp\left[-\left(\frac{4\pi R_q\cos\theta}{\lambda}\right)^2\right]$$

(4.265)

Obviously for normal incidence this reduces to

$$I_r(s) = K \exp\left[-\left(\frac{4\pi R_q}{\lambda}\right)^2\right]$$

(4.266)

(making the assumption that the roughness of the surface is much less than $\lambda$). This has been found by Beckmann and others as reported. It is simple and refers to the wavefront of the ray from the surface. Specular methods are in effect a resort to geometric optics and as such represent one stage simpler than even the scalar methods. How geometric scatter and scalar scatter can be related to each other will be shown in the section on optical function.

Plotting equation (4.265) as a function of $R_q$ and $\theta$ yields a family of curves similar to that shown in figure 4.214.



**Figure 4.214** Specular reflection as a function of finish.

If the intensity axis is plotted logarithmically then the curve looks like figure 4.215. For an instrument the best value of $\theta$ is the one in which $(\mathrm{d}I/\mathrm{d}R_q)$ is a maximum.

Many instruments have been devised which work on this principle and are used with some success as comparators but not for measuring absolute values.

It is sometimes possible by taking advantage of the non-linear characteristics of detector devices to get an almost linear response of the output of the instrument as a function of $R_q$ [188], as shown in figure 4.216.

For each type of surface the instrument has to be calibrated to take into account such things as reflectivity, lay, etc.



**Figure 4.215**

**Figure 4.216** Resultant output as a function of angle.

Much more general results are obtainable if more than one detector is used. The normal configuration is one in which two detectors are used. One is placed to pick up the specularly reflected light and the other is placed at an arbitrary angle to it to get an estimate of the diffused light.

The diffusion angle $\beta$ is taken at 20° or so and $\theta$ is usually near to 45°. To compensate for possible changes in the light level the surface quality is often expressed as a ratio $Q$ (for surface quality) where

$$Q = \frac{A - B}{A + B}.$$

(4.267)

The rationale behind this is that for very fine surfaces $B \simeq 0$ and $Q$ becomes 1; for very rough surfaces $A \sim B$ and $Q$ becomes approximately 0. This is illustrated in figure 4.217. The value of $g$ becomes bigger as the digrams in figure 4.218 are viewed from left to right, illustrating the Beckmann philosophy.

There are many variants on this simple theme, for example more than one diffusion angle can be taken—sometimes as many as 10 have been used. The limitation is that of practicality, especially with regard to the shape of the surface. To enhance the range of the device infrared is sometimes used. Also, to resolve the signal-to-noise ratio the signal is sometimes chopped. Characterization based on fitting a model to the specular scattered light for the different angles has been attempted. The model purports to give $R_q$ and $R_{ku}$ for the surface. The method is very specialized having only been applied to sheet steel. It is therefore process sensitive [189].

Instead of a set of detectors a photodiode array can be used, in which case the variance of the distribution is used as a measure of the 'characteristic' roughness $S_N$. Thus

$$S_N = k \sum_{i=1}^{n} (i - \bar{i})^2 P_i$$

(4.268)

where $i$ is the diode number, $P_i$ is the intensity value of the $i$th diode, $\bar{i}$ is the mean diode number as ascertained by the intensity values and $k$ is a scale factor [190] (figure 4.219).



**Figure 4.217** Glossmeter arrangement.

**Figure 4.218** Quality measure of surface using scattered light: (*a*) smooth surface; (*b*) medium surface; (*c*) rough surface. $Q = (A - B)/(A + B)$.

Typical variabilities of these scatter methods when compared with stylus methods are 15–20% for peak-type parameters and 5–10% for average values such as $R_a$ (Note that this is after the mean values of both have been adjusted to be equal.)

Fibre optics have also been used to pick up the scattered light, as shown in figure 4.220.

In this method, half of the fibres are used to illuminate the surface and half to pick up the scattered light. Sometimes the fibres are arranged at random in the tube and sometimes are concentric. Fibre optics may seem to be somewhat irrelevant to the system when considering reflection and diffusion, but it is not necessarily so because, as Thwaite points out [174], the whole optical system can be made more compact simply by using fibre optics. This removes the need for beam splitters and generally makes the optics more flexible. Obviously if coherent light is to be used a monomode fibre has to be included, especially if diffraction effects are being studied. As with all things there is a disadvantage in the fact that quite a lot of light could be lost. This has caused problems in measuring surfaces during grinding. Damage to the end of the fibres by chips etc can also cause spurious scattering.

Other alternative collectors are TV cameras and charged coupled devices which obviously take in much more information but are not really designed for use in such hostile environments.

A great deal of work has been carried out by Takayama *et al* [191] on single and multiple fibre collections, taking various ratios and plotting outputs as a function of angle. The results, as expected, indicate that within processes—especially single-point processes—it is possible to get useful results providing that the tool radius is not changed when it becomes necessary to recalibrate.

One of the practical problems encountered with these scatterometers (sometimes called glossmeters) is that the optics is often positioned close to the surface in order to reduce the effect of ambient light. This results in a need to protect the lenses from debris.

An alternative approach to taking a number of detectors to scan the diffused field is to measure the whole of the light scattered from the surface (total integrated scatter, TIS) (figure 4.221).



**Figure 4.219** Multiple detector arrangement.

**Figure 4.220** Use of fibre optics to measure scatter.

From the earlier expressions derived using scalar theory, the total scatter expressed as a ratio of the diffuse reflection is given by

$$\frac{\text{diffused}}{\text{total}} = \text{TIS} = 1 - \exp\left[-\left(\frac{4\pi R_{\mathrm{q}}}{\lambda}\right)^2\right]$$

(4.269)

which to a reasonable level of accuracy can be expressed as

$$\text{TIS} \sim \left(\frac{4\pi R_q}{\lambda}\right)^2 \text{ for } Rq \ll \lambda$$

(4.270)

following Davies [170]. This assumes that the scatter is mostly confined to a small angle around the specular direction.

Church *et al* [184] have approached the whole problem from a simplified vector theory to evaluate the differential solid angle which, as mentioned, is related to the power spectral density (PSD) and some useful comparison between the two methods have been obtained.

TIS methods are now widely used worldwide for estimating the surface finish of very fine surfaces. It should be pointed out, however, that whereas the TIS method was originally devised and based upon the theory of Davies (which assumed that the surface was Gaussian in its statistical properties) it is not a requirement as far as the height distribution is concerned [192].

TIS methods are now especially being used to estimate the surface roughness of very fine surfaces in the range of 0.1 to 10 nm in applications such as laser mirrors and semiconductor surfaces (e.g. [193, 194]). It seems obvious that for measuring surfaces which are ultimately going to be used as mirrors, the use of light scatter techniques should be adopted, but this has not always been the case. However, it is certainly true that for measuring very fine surfaces which have soft finishes, such as aluminium, germanium, copper or whatever, the TIS method is a useful tool to have when stylus methods could be damaging to the surface.



**Figure 4.221** Total integrating scatter (TIS).

From equation (4.270) it follows that the RMS value of the surface, $R_q$ is given by

$$R_q \simeq \frac{\lambda}{4\pi}(\text{TIS})^{1/2}.$$

(4.271)

Usually the TIS measurement is made by means of measuring the total light scattered into an integrating sphere (figure 4.221). Obviously, tracking a single detector around the space by means of a goniometer is an alternative instrumental arrangement, but great care has to be taken to eliminate the effects of the finite size of the detector [195, 196]. In such a TIS method the surface is usually completely opaque. If not, the incident beam has to be projected at an oblique angle of typically 30° because of bulk scattering or scatter from the back surface of the specimen (remember that A has to be modified by sec $\theta$). TIS has also been used in the infra-red, see section 4.3.11.

Various estimates have been made to establish cross-correlations and how they are affected by polarization and angle of incidence [197]. A comparison between the various methods of measuring smooth surfaces [198] has indicated that heterodyne methods at present tend to give lower values of surface finish than TIS or stylus methods because of the limited spatial resolution (typically 2.5–200 $\mu$m) compared with the much higher resolutions of the stylus (0.1–10 $\mu$m) or TIS (1–10 $\mu$m).

Instruments for measuring the scatter of light from surfaces have been used with varying degrees of success. They are quick and relatively cheap. Unfortunately they are not very reliable for making absolute assessments of surface geometry largely because of the assumptions that have to be made and the wide variety and range of engineering surfaces. However, the scatterometer or glossmeter can be used very effectively as a quick comparator. For this reason it can be used in a quality control loop once the credibility of a process has been verified by using a stylus instrument. The fact that the information offered up by a glossmeter is sometimes garbled should not be a deterrent to use one if a comparison is being made. It should also not be forgotten that it is likely to be much cheaper than a stylus instrument. One general comment reflecting a thread that penetrates all optical methods is as follows.

If there is one factor that has held back the use of optical methods, whether sophisticated or simple, it is the lack of national standards on optical methods. All stylus methods are keyed solidly into the national and international standards systems with the almost psychological result that they are used with confidence. No such feeling exists with optical methods. This neglect is not intentional; it arises because of a number of factors. One is the sometimes conflicting evidence on the usefulness of the optical method, which comes about mainly because of the lack of credibility brought about by having to use mathematical models of surfaces rather than properly prepared surfaces. Another, probably more important, reason is that everyone is attempting to correlate optical methods with stylus methods. This may be a reasonable attitude with rough surfaces where geometric optics could come into play and some agreement could be expected, but for very fine surfaces there is little reason why optical and tactile measurements should agree; they are measuring different things! The only way out of this quandary is to relate optical methods (and standards) directly to function (and to manufacture) by passing the link to tactile methods. Belatedly some efforts are now being made in the USA, Germany and the UK to remedy this situation.

### 4.3.13 Flaw detection

#### 4.3.13.1 General

In this section the problem of flaw measurement will be considered, but before this the concept of a flaw must be defined. Here it is taken to mean a singular or infrequent event that is not encompassed by the statistics of the rest of the surface. Obviously this would mean something like a scratch or pit. It is also taken to mean something that is not wanted and represents either a fault in the manufacture of the part, as in the handling of the part during the process, or the deterioration of the part with storage. It is also usually taken to mean something which is not regular or even numerous. Therefore it is usually difficult to find.

It could be argued that waviness should be included in the definition of a flaw because as a rule it is produced by machine malfunction and, consequently, it is not wanted. However, here in this book waviness is considered to be more like a continuous flaw which tends to underlie the whole of the roughness. A flaw is taken to be a feature of discontinuous nature not present at large in the surface geometry but very localized and usually a singularity. Because of this definition methods of investigating flaws can follow a different strategy from those used for detecting waviness.

Other types of flaw are errors in a structured pattern such as circuit faults on a VLSI mask. Such defects require another strategy.

The basic problem with flaws is that they may or may not be functionally important. Some are disastrous, such as a fatigue crack or corrosion pit, whilst others may be simple discolorations having no important role apart from cosmetic effects, which could be a detraction when selling the part.

Consequently flaw measurement could consist of a methodology which has three stages:

1. Detection of the flaw. This could be a preprocessing exercise to localize the area of investigation.
2. Characterization. This could be a pattern recognition exercise.
3. Measurement or quantification.

Because a flaw is an exception it has to be searched for. Also, it is likely to be small. So the requirements for an instrument in terms of resolution and viewing range are very high. For this reason any technique has to be fast. This consideration rules out conventional stylus methods and so usually optical methods are adopted.

There are three ways of tackling the optical instrumentation: (i) in the image plane or (ii) in the far field or diffraction plane or (iii) in both or neither (figure 4.222).



**Figure 4.222**

It might be considered that the image plane is the obvious place to look for a flaw. This is based on the fact that the appearance of a flaw such as a scratch is different from that of the rest of the surface. The problem is that it could be anywhere and has to be searched for over what could be a large area. It may be that the position as well as the presence of a flaw may be needed, as in the case of fatigue where a crack in a highly stressed part of the surface (such as the root of a turbine blade) could be serious, whereas on the rest of the surface it is not.

If position of a flaw is important then the search can only take place in the image plane. However, if position is not a problem then the situation is different. It should be advantageous to use the diffraction plane.

### 4.3.13.2 Transform plane methods

There is a plausible argument, which says that if the image of a flaw is different from the rest of the surface then the diffraction of the scratch will therefore be different from the diffraction of the surface. If this is the

case then the characteristics of the normal surface can be removed from the scene in the diffraction plane and the presence or absence of a flaw determined. It should also be possible to characterize the flaw using the argument that flaws have different scattering properties. This aspect will be considered shortly.

Other possibilities exist that do not require the high-resolution scan which has to be used in the image plane. One such technique is the wide-field approach. This is one of a set of techniques that artificially enhance non-standard aspects of the surface. This then considerably reduces the problem of detection with less sensitive equipment or alternatively allows a faster scan.

There are other differences in approach between image processing and diffraction plane processing. One of the most important is that, in the case of the diffraction plane approach, there is no need for a scan; the whole surface can be floodlit by the light. The spectrum in the diffraction plane is insensitive to the position of the flaw in the object field. Also, because it is the nominal power spectrum of the surface at the diffraction plane the signal will be more or less stationary in both the statistical sense and the instrumental sense. As a result the signal in the Fourier plane will be steady or more steady than any part of the signal in the image plane.

Another point is that because there are no problems with focus, the diffraction plane is insensitive to the position of the object in the $z$ direction. Obviously image plane flaw detection is very sensitive to object position axially.

One disadvantage of flaw detection in the diffraction plane is that there will be very little signal. This argument results from the fact that the flaw is physically very small relative to the rest of the surface. It need not be an insurmountable problem because a sensitive detector could be used at the high-frequency part of the spectrum to pick up the flaw signal. It is unlikely that any normal surface signal of high intensity will be present at the high-frequency end of the spectrum.

Because the flaws can be at any orientation the flaw detectors in the diffraction plane would have to be either annular, centred around the origin or linearized at the transform plane by means of a cylindrical lens.

If the fault is a defect in a structured surface, such as a microcircuit mask, it is advantageous to consider detection in the Fourier plane and then subsequent reconstruction. The arrangement for detecting such a defect can be the following in figure 4.223.

Prior to investigating the structured surface the transform of a perfect surface (either a real object or a computer-simulated one) is photographed. This negative is left in the transform plane $P_2$. If now a test object is put in the object plane $P_1$ its transform will coincide exactly with the negative of the perfect object, providing the orientation is the same. If the test object is perfect then the reconstituted object at plane $P_3$ will show nothing because it exactly matches the perfect object transform. However, if there is a fault it will show as a bright spot in the reconstructed image plane at $P_3$. If there is more than one fault then more than one spot will show. This technique has the considerable advantage of reducing the information bandwidth of the object by many orders of magnitude. Simple detection schemes can show the position and nature of faults.



**Figure 4.223** Optical filtering to detect flaws.

The method has the disadvantage that knowledge of the perfect part is required. This is usually not a problem in the case of structured or patterned surfaces but it is for more random objects such as moving steel sheets in which the 'perfect' object is random and subject to a certain amount of statistical variation.

The fluctuation makes the use of a mask impractical.

If the transform (diffraction) plane is to be used then there might be another way to remove effectively the presence of the normal surface, such as by using oblique illumination [177] as seen in figure 4.224. By making the viewing angle oblique the normal surface becomes progressively demagnified because of the cos $\theta$ effect—the wavelength of the light $\lambda$ becomes $\lambda/\cos\theta$ and the normal surface marks become more specular.



**Figure 4.224** Separation of scratch scatter by oblique illumination.

If some Beckmann assumptions are made, for example that the phase expression $\exp(\mathrm{j}\psi)$ can be expanded into the form

$$1+\mathrm{j}\psi-\psi^2/2 \quad \text{where } \psi=(4\pi/\lambda)z'(x,y)\cos\theta \tag{4.272}$$

and $z'(x,y) = z(x\cos\theta, y)$ where $z(x,y)$ is the surface texture at normal incidence (of variance $\sigma^2$), then $I(u,v)$, the intensity at $\rho$, is given by

$$I(u,v)=K_{\mathrm{o}}\,|A_{\mathrm{o}}|^2\left\{2\pi\left[1-\left(\frac{4\pi}{\lambda}\cos\theta\right)^2\sigma^2\right]\right\}\delta\left(\frac{2\pi}{\lambda z_{\mathrm{o}}}u,\frac{2\pi}{\lambda z_{\mathrm{o}}}v\right)$$
$$+\left(\frac{4\pi}{\lambda}\cos\theta\right)^2\sigma^2\cos\theta\rho_1\rho_2\exp\left[-\rho_1^2\frac{\cos^2\theta}{2}\left(\frac{2\pi}{\lambda z_{\mathrm{o}}}\right)^2u^2-\frac{\rho_2^2}{2}\left(\frac{2\pi}{\lambda z_{\mathrm{o}}}\right)^2v^2\right]. \tag{4.273}$$

Equation (4.273) is approximate only and assumes $8z'\cos\theta \ll \lambda$.

If the area illuminated is large the DC component of the light is converged into a small spot (expressed as a delta function in equation (4.273)) which can be effectively blocked off by a central stop at $P$.

The intensity then becomes

$$I_{\mathrm{out}}=\frac{K_1}{\alpha_1\alpha_2}\cos^2\theta[\mathrm{erf}(L\alpha_1\cos\theta)-\mathrm{erf}(l\alpha_1\cos\theta)][\mathrm{erf}(L\alpha_2)-\mathrm{erf}(l\alpha_2)] \tag{4.274}$$

where

$$\alpha_1=\frac{\rho_1 2\pi}{\lambda} \qquad \alpha_2=\frac{\rho_2 2\pi}{\lambda} \qquad K_1=2K_{\mathrm{o}}\,|A_{\mathrm{o}}|^2\left(\frac{4\pi}{\lambda}\right)^2\sigma^2\rho_1\rho_2. \tag{4.271}$$

Equation (4.274) represents the intensity pattern in the Fourier plane with 'central dark ground' conditions. If a scratch is present it can be regarded as independent of the general surfaces. Then

$$z(x, y) = z_N(x, y) + z_s(x, y).$$

(4.275)

Thus if the normal surface $z_N(x, y)$ is assumed to be white noise and $z_s(x, y)$, the scratch, is small, and the surface in the $x$ direction is assumed to be independent of that in the $y$ direction

$$I'_{out} = \langle | F_n |^2_{P\text{-}P_o} \rangle + \langle | F_s |^2_{P\text{-}P_o} \rangle$$

(4.276)

where P—$P_0$ is the annulus between the outside restriction of P and central stop $P_0$,

$$I'_{out} = (1 - \varepsilon) | F |^2 + F_s^2 (P - P_0)$$

(4.277)

where

$$\varepsilon = \frac{s_o}{s} = \frac{\text{of scratch}}{\text{total area}}.$$

If the sensitivity to a scratch is defined by

$$\text{sensitivity} = (I'_{out} - I_{out})/I_{out}$$

(4.478)

this is as plotted in figure 4.225

Having plotted the sensitivity, it is obvious that the maximum sensibility to a scratch is for angles above 75°.

As the angle of obliquity increases the normal surface becomes more specular and the capability of scratch detection improves.

The sensitivity does not in practice increase right up to $\pi/2$. This is because at the extremes of angle the diffraction due to the scratches also becomes demagnified, so the effectiveness of the method depends on picking an angle in which the scratches are not demagnified but the normal surface is. For a ground surface finish of $R_q \simeq 0.2$ $\mu$m and scratches of approximately 10 $\mu$m in width, the optimum angle is about 80°. It seems therefore that, in principle, the use of various stops and masks in the Fourier transform plane to get rid of the background surface scatter could be used for flaw and defect detectors. This applies to random as well as structured surfaces. The only objection, especially when using oblique illumination, is that of practicability. Only



**Figure 4.225** Sensitivity of a scratch.

quite flat surfaces can be measured because of the need for glancing angle illumination to demagnify the basic surface-finish scatter.

### 4.3.13.3 Image plane detection

Many items requiring fault detection are flat and moving. Typical examples are in the sheet aluminium and sheet steel industries. In these cases the fact that the strip is moving alleviates the scan problem. Only one degree of freedom has to be engineered. The other at right angles is provided by the movement of the strip which is usually very fast. In the hot mill process and cold reduction mill, speeds can be up to 10 m s$^{-1}$. For direct image processing a TV camera could be used (see [199]).

The output is often differential to enhance defect edges and generally to reduce the bandwidth. Once the defect is picked up it is then classified. It is not usual to classify into as many groups as shown in Chapter 2 on flaws, but even for steel there can be many. Some typical ones are given below. The feature is obvious from the name:

(1) slivers
(2) scale
(3) pits
(4) gauges
(5) blisters
(6) stains
(7) seams
(8) scratches etc.

These defects in the image plane are found by the camera using a matching procedure or enhancement. Then they are classified in a computer using pattern recognition methods. Parameters other than those above might be clustering of defects, distances between clusters, etc, and would be used to detect longer-term changes in the manufacturing process.

Other sheet materials could be photographic film or paper. Examples of objects which are not so well behaved are cans and bottles and other food containers. If the outside is being examined on a can, for instance, it has to revolve as well as move in order to expose all the surface to inspection.

Mechanical scan devices have to be used in many cases and their design is complicated even for a linear scan or circular scan—used for example in flaw detection on roads or in tunnels and pipes. One such system could involve the use of a polygon. An example of a typical schematic view of a scanner using a polygon is shown in figure 4.226.



**Figure 4.226** Laser scan system with polygon.

The lens is so positioned that whatever the angle of light from the polygon, the ray hitting the surface is normal. This cannot be held to better than a degree or two but it is necessary so that realistic information from the surface is obtained and no effects like shadowing occur. The problem is that often the slopes on surfaces have rather a small value. At the same time the laser spot has to be comparable with the stylus tip dimension. This imposes serious restrictions on the optical design. The Seidel aberration that causes most problems is coma. The only practical way to reduce the problem to a solvable one is to expand the beam before it hits the polygon. Then the beam has sufficient width to condense to a small spot. Under these conditions, in order to get the required scan it is necessary that (i) the lens system is twice as wide as the scan, and (ii) the incident beam onto the lens is always equal to half the lens diameter.

With the small spot and the scan comes another problem, which is the depth of focus. This should be comparable with the stylus instrument. It is here that there is a fundamental conflict of physical laws, as shown earlier in the chapter.

Consider a laser beam, for example. At the narrowest part of the waist the disc is the size of the diffraction disc related to $\lambda$ and the NA. The diffraction focal distance can be applied as the axial length over which the energy density is substantially constant. This is approximately

$$\lambda/(NA)^2 \tag{4.279}$$

Sometimes the distance is defined as between the cross-sections of $\sqrt{2}$ multiplied by the minimum disc diameter, as in figure 4.103. As a comparison see table 4.7.

**Table 4.7**

| NA ($\mu$m) | Diameter ($\mu$m) | Length |
|---|---|---|
| 0.08 | 9 | 94 |
| 0.25 | 2.9 | 9.6 |

From table 4.7, it can be seen that if the spot size is to be comparable with the largest acceptable stylus tip used, at present 10 $\mu$m, then the depth of focus is about 0.1 mm or four-thousandths of an inch, which is just about the range of movement used on the very old stylus instruments. Any attempt to reduce the spot will necessarily reduce the depth of focus. This is a serious matter because it restricts the inherent used of a scanner to pick up fine detail of the surface if it has a large form error. Obvious ways out of this, as mentioned in the section on optical styluses, are (i) to move the optical stylus on a fast acting servo, similar to the Dupuy method; (ii) to use the scanning system with the probe null-detecting device controlling a registration of scan displacement (giving as it were points on the contour).

Another point about scanners concerns the effect of the polygon itself. Even if perfect optics could be used, the formation of a spot which is, say, diffraction limited would be impossible because of the non-stationary virtual centre of rotation of the scanning beam. By 'virtual centre of rotation' is meant the point from which the centre line of the reflected beam emerging from the scanner would appear to rotate. Because the reflecting surface of the polygon is a flat plane and the axis of rotation of the surface is not within that plane, the intersection of the surface with the centre line of the incident light beam is not a fixed point in space but moves in position axially along the incident centre line. As a result, the intersection of any two reflected beam centre lines when projected backwards for all positions of the scan occurs not at a single point but over an area (figure 4.227).

The envelope shown in figure 4.227 is typical of a 12-sided polygon of about 20 mm radius. Notice the rather large distances involved.

**Figure 4.227** Location of centre of polygon face with incident beam.

Obviously the alternative solution to mechanical scanners is to use a mirror where the axis of rotation can be very close to the plane mirror itself. This is acceptable in certain cases but it can definitely introduce vibrations if the mirror is flapping and is more difficult to drive if rotated (figure 4.228).

If a rotating mirror is used the entrance pupil has to be placed at the front focal plane of the lens. Arranging that the centre of the input beam coincides with the centre of the pupil ensures that the focused beam always appears like an axial beam when viewed normally to the surface. If the mirror is rotated through $\theta/2$ then the focused spot moves through a distance of $d = f \tan \theta$, where $F$ is the focal length of the lens. This is also the height of the centre of the beam on the lens which accounts for the absence of beam tilt in the object space (figure 4.228).



**Figure 4.228** Scan with constant numerical aperture—flapping mirror.

Axial displacement of the reflector does not affect spot position or size on the surface. A parallel input beam across the lens results in no change in focus position or change in aperture angle. This is because of Langrangian invariance in the object space. Thus $d\alpha$ = constant. Consequently the spot size is unchanged at approximately $2.44\lambda F$, where $F$ is the ratio of the focal length of the lens to the diameter of the entrance beam. The spot size is the airy disc. In general the optical system arrangement follows normal practice:

1. The diffraction-limited spot is obtained by freeing the image from coma and sperical aberration.
2. The spot size is maintained by having flat astigmatism-free correction over as large a field aperture as possible.

These problems for a scanning system have to be solved when the stop is situated at the front focal plane of the lens, as seen in figure 4.229.



**Figure 4.229** Lagrangian invariance of constant NA system.

Scanners which employ electronic means, such as in the scanning electron microscopes, rely on deflecting a beam of light or electrons and will not be considered here.

In the cases discussed the defects have been in some instances enhanced by means of optical filtering in the transform plane or electrical filtering on the signals. It is possible to get some enhancement in the scattering domain itself as seen in the 'D field' method.

### 4.3.13.4 'Whole-field' measurement—flaw detection

In the system shown in figure 4.230 the light from the source is cast onto the surface under test.

The scattered light reflects from a retroreflector made up of a large number of small retroreflecting elements such as beads. Some of the light again reflects back from the surface onto a detector.



**Figure 4.230** Whole field—'D' sight for form or flaws.

The whole system works on the principle of one-to-one imaging—almost like a flying-spot system except that the spot does not move. This is equivalent to a discrete spatial spot imaging system and has all the signal-to-noise improvement of a flying-spot method without the need to scan dynamically. It has been named the 'D' sight. Basically if there is any flaw on the object it throws light in all directions and this is lost to the detector because it does not hit the retroreflector at the right geometric angle. The parts of the object which are smooth simply obey reflection laws and bounce light onto the retroreflector and then back onto the detector. Hence all flaws and scratches appear dark and the normal surface is bright.

This is an optical system for enhancing defects so that their presence or absence can be recognized. Parameters such as a density count can be measured but not the depth of the flaw. This technique has been used for large bodies as well as small [200].

Although most flaw detection schemes rely on a pattern recognition scheme the problem often still arises where the flaw or defect has to be quantified. Because the defects have rather a wide range of properties, mostly including sharp or rough edges such as pits and scratches (but also sometimes having smooth edges like grooves or hollows), it is necessary to quantify sizes and position with a special instrument.

What is remarkable is that there is no standard recognized way to do this. As a result problems arise when components are transferred between companies. The international standards committees have recognized this problem (ISO/TC172/SC1) and are attempting to solve it. One of the principal advocates of standards for measuring flaws is Baker [201, 202].

In particular he proposed the use of a special purpose microscope called a microscope image compositor, which works in the image plane. This microscope, instead of being used to study the fine structure of an object for the purpose of recognition, requires a knowledge of the type of object present, and its function is to relate the size, light scatter or density of the defect to a set of reference standards incorporated within the instrument.

As a first step the comparison has been used only for scratches and digs on an optical surface. The basis of the comparison is to adjust the instrument until the true scratch on the surface matches in visibility the standard scratch held within the instrument [95] (figure 4.231).

The system works on polarized light. Adjustment of polarizer $Z_2$ is made until equal visibility is achieved between the test and the reference scratch. The interesting point about this system is that the reference defect



**Figure 4.231** Flaw comparator system—Baker.

is itself a real object kept in the instrument. It is not a data set in a computer. The use of data sets is more or less dominated by systems in which the flaw detector is only one part of a complete package of software which not only identifies flaws but also the object size and shape [203].

### 4.3.13.5 Comment

Looking for flaws is a good example of the dilemma facing the metrologist. This dilemma is that the direct way of examining the surface flaws is very difficult to implement satisfactorily because of their complexity and variety and also the high speeds at which they have to be examined. The other impediment is the large range to resolution demanded from the instrument.

The result of these problems is that the metrologist uses an indirect method instead (e.g. using the Fourier transform plane rather than the image plane). This move can alleviate some of the above problems and for this reason can be justified in use. However, the indirect method or parametric method only works properly in well-defined cases. If there is any ambiguity or large noises present, then it is more difficult to appeal to experience or common sense if the problem has been posed in an indirect frame of reference. Hence direct methods should be used whenever possible. Only for overwhelming reasons should an indirect method be used over the direct one.

### 4.3.14 Comparison of optical techniques

The question arises as to which optical method is best and how it relates to stylus methods. Using the criterion adopted for transducers, namely the frequency response and the range-to-resolution ratio, it can be seen from figure 4.232 that there is a definite pattern.

At present the optical techniques are faster, yet they have a limited range to resolution; the stylus methods are improving in range to resolution but are not getting much faster. This is due to the amount of processing necessary to unravel different features taken from an integrated measurement. It appears that those optical methods most likely to make a large impact will be the rather crude ones suitable for in-process gauging.

The term optical used here refers to the 'parametric' methods rather than the profile methods or followers. The term parametric is used when there is an indirect way of estimating the roughness parameters. This is usually via an 'inversion' technique where the way in which light interacts with the surface infers some value of the surface geometry. This inferred value is obtained by assuming some statistical model for the surface.



**Figure 4.232** Trends in optical systems versus mechanical probe system.

This method would be used more for comparisons rather than as absolute, except in isolated special cases such as the diffraction device.

The stylus techniques look to be most useful for post-process verification—much the same role as they have at present. It should be pointed out that ideally the preferred way to go should be the vector sum of the two (shown as the broken line in figure 4.232). It is possible that this will be achieved by a combination of mechanical fixturing together with optical (or other) probes incorporating a limited amount of logic. This is now possible using semiconductor materials. More will be discussed on this subject when the issue of measurement in manufacture is taken up in Chapter 5. What is certain is that optical methods offer an impressive range of options, much wider in fact than stylus methods. As yet, however, there has not been one optical method that can cover both the in-process requirement and the integrated measurement requirement. It is still an either/or choice. The same is true for stylus methods, only that they are much more restricted in the in-process requirement, yet more versatile for integrated measurement.

### 4.3.14.1 General optical comparison

Table 4.8 provides a performance comparison for various optical techniques. The values given are only typical—what is more important is the relative values between methods.

**Table 4.8** Optical comparison. (Para = parametric. Prof = direct)

| Technique | Type | Vert. res. | Vert. range | Lat. res. | Lat. range |
|---|---|---|---|---|---|
| Specular | Para | $\sim l$ nm | $\sim 1\ \mu m$ | $\sim 1$ mm | $0.6\ \mu m$ |
| TIS | Para | $\sim l$ nm | $0.1\ \mu m$ | $\sim 1$ mm | $0.6\ \mu m - 0.1$ m |
| Diffuse | Para | $\sim 10$ nm | $\sim 10$ nm | 1 mm | $0.6 - 0.1\ \mu m$ |
| Diffraction | Para | $\sim 10$ nm | $\rightarrow \lambda/8$ | $\sim 10\ \mu m$ | $10\ \mu m - 10$ mm |
| Laser speckle contrast | Para | $\sim 10\%$ of range | $\rightarrow \lambda/5$ | $\sim 10\ \lambda$ | $0.6 - 0.1\ \mu m$ |
| Polychromatic speckle contrast | Para | $\sim 10\%$ of range | $\rightarrow \lambda/5$ | $\sim 1$ mm | $0.6 - 0.1$ m |
| Speckle correlation | Para | $\sim 10\%$ of range | 0.15-6  8-32 | $\sim 1$ mm | $1\ \mu m$ |
| Ellipsometry | Para | 1 nm | | $\sim 1$ mm | $100\ \mu m$ |
| Interferometer | Prof | 1 nm | $0.01 - 10\ \mu m.$ | $2\ \mu m$ | 1 mm |
| Follower focus | Prof | 1 nm | $0.01 - 5\ \mu m$ | $1.6\ \mu m$ | 10 mm |
| Follower heterodyne | Prof | 0.1 nm | $2\ \mu m$ | $1.6\ \mu m$ | 10 mm |
| Oblique | Prof | $0.5\ \mu m$ | $2 - 200\ \mu m$ | $5\ \mu m$ | 10 mm |

## 4.4 Capacitance techniques for measuring surfaces

### 4.4.1 General

Capacitance methods for measuring surfaces seem to have a long history. Perthen [9] initiated the idea at much the same time as stylus and optical instruments were first being investigated. However, historically the method did not catch on, mainly because of the practical problem of having to match the shape of the electrode to that of the surface [204].

Fundamentally the capacitance is known to be a function of a number of things:

(1) the useful area of the conducting plate;
(2) the separation of the plates;

(3) the dielectric constant between them.

Usually capacitance instruments are used for proximity gauges and so it is the gap between the electrodes which is the variable to be measured, although measurement of the capacitance of the actual surface is being explored with scanning microscopes.

The general theory is quite simple for a flat-plate system [205]. The capacitor is a two-plate system. The plates are conducting; one is of area $A$, which corresponds to the probe that is close to the surface being measured. The surface acts as the other electrode and so must also be conductive.

For any small area $\delta x \delta y$ of the specimen separated by $Z$ from the probe the capacitance is given by

$$\delta C \propto \frac{\delta x \delta y}{Z}.$$

(4.280)

Integrating over the area of the probe $L_x L_y = A$

$$C = K \int_0^{Lx} \int_0^{Ly} \frac{\mathrm{d}x \, \mathrm{d}y}{Z}$$

(4.281)

or $C = KA/Z_m$ where $K$ includes the dielectric constant and $Z_m$ is the reciprocal of the mean value of $1/Z$ over area $A$.

Let the surface finish be $z(x,y)$. The gap will be

$$Z = t_0 + z(x, y)$$

(4.282)

so

$$C = KA \int_0^{Lx} \int_0^{Ly} \frac{\mathrm{d}x \, \mathrm{d}y}{t_0 + z(x, y)}.$$

(4.283)

The probe system and capacitance diagram from reference [205] are shown in figure 4.233.

The plane represented by $Z_m$ thus represents the effective position of the workpiece surface in terms of the observed capacitance. Sherwood called this the 'capacitance plane'.



**Figure 4.233** Capacitance method of measuring surface texture.

The capacitance can be measured and the area is known, so in principle it should be possible to determine something about $z(x, y)$. Note that

$$\frac{1}{Z_m} = \frac{1}{L_x L_y} \int_0^{L_x} \int_0^{L_y} \frac{dx\, dy}{t_0 + z(x, y)}.$$

(4.284)

$t_0$ would normally be fixed for a given probe, usually by means of non-conducting feet that rest on the surface to maintain the gap. Problems arise if the surface is curved, in which case the probe shape has to be changed to the opposite form.

The unfortunate fact is that the capacitance is related to the inverse of the roughness, and large peaks are unduly weighted relative to valleys. This may be quite a useful feature in some functions but for a general purpose instrument it is not a sound basis. It follows from this effect that capacitative methods are sensitive to the form of the amplitude distribution of the surface heights.

If the displacement of the capacitance plane from the plane of the peaks is denoted by $R_c$ then

$$Z_m = R_c + t_0.$$

(4.285)

Alternatively, $R_c$ also represents the displacement of the capacitance plane for a given surface from that obtained by nulling onto an optical flat.

Equation (4.283) has been applied to a number of different surfaces. These have been idealized for ease of calculation.

Some numerical solutions assuming random surfaces having a Gaussian distribution indicate that the capacitance is very sensitive to the mean plane separation [206] (figure 4.234). The method therefore has to be viewed with some reservations, especially in view of the fact that it is very temperature sensitive.

Recent attempts to incorporate compound dielectrics as part of the probe system [207] and to make the probe flexible certainly get rid of the probe shape problem for non-flat surfaces. Yet this method depends somewhat on the contact pressure and the problem of wear of the dielectric with repeated use.

As the roughness is basically a function of the air space or levelling depth between the high peaks and valleys the effective levelling depth increases as the roughness increases. In the device shown in figure 4.235 the probe is brought into contact with the workpiece surface. The capacitative sensor is firmly held against the surface by a non-conductive elastomer cushion. The capacitance is read off directly. The surfaces have to



**Figure 4.234** Sensitivity of capacitance to mean plane separation.

be reasonably clean in this sort of technique because the insulating values of air and other materials, especially oil and cutting fluid, are different. Dirt particles on fine finishes tend to force the sensor pad away from the surface giving the effect of a large roughness.



**Figure 4.235**  Portable capacitance probe with self-adjustable sensor.

### 4.4.2    Scanning capacitative microscopes

These are now being used to measure very fine detail [208] on surfaces with limited spatial resolution.

In figure 4.236 the probe itself is an electrode, which is supported and scanned above the workpiece at height $z$ above it. The capacitance between it and the component is measured. This capacitance is a function of the geometry of the electrode.

In the usual method the electrode is mounted on a diamond stylus (about the same size as a mechanical stylus used in conventional instruments). Given that the height $z$ is much less than the width of the electrode (into page, $W$) the sensitivity of the capacitance $C$ to changes in the height $z$ is given by equation (4.286) below and shown in figure 4.234:

$$\frac{\mathrm{d}C}{\mathrm{d}z} = \left(\frac{8}{\pi}\right)\varepsilon W/z$$

(4.286)

where $\varepsilon$ is the dielectric constant for air for $W = 2.6\ \mu$m and

$$z = 20\,\mathrm{nm} \qquad \frac{\mathrm{d}C}{\mathrm{d}z} = 5.6\,\mathrm{aFnm}^{-1}.$$

(4.287)

Taking noise into account a vertical sensitivity of about 0.3 nm can be detected (comparable with existing surface instruments). The problem is of course that the device makes contact with the workpiece.

The capacitance method suffers from the problem that it measures variation in dielectric as well as topography. It is also not so sensitive in the $z$ direction as the tunnelling microscopes because the relationship between $z$ and tunnelling current is power-law related, whereas in capacitance methods it is reciprocal



**Figure 4.236**  Scanning capacitance probe.

only at best. Also the lateral resolution is somewhat larger than the tip dimension. It therefore is not the preferred method if very high lateral resolution is required. Obviously the output format can be made the same as any of the other scanning methods. Its field of measurement is determined by the carriage or probe translation units. The depth of field is quite limited; the plausible variation is determined by the reciprocal relationship between capacitance and $z$. In practice it can be compared with the Nomarsky microscope. It is very much inferior to the scanning electron microscope, yet the SEM has a limited field.

In summary, the capacitance microscope is a useful addition to the available surface instruments. It does seem to have the capability of good sensitivity at the smaller separations between the workpiece and the probe because of the relationship $\delta C/C = -\delta z/z$. At the present time it does not have the versatility of application of any of the 'atomic-scale' microscopes.

Capacitance transducers have been used not only to measure surface capacitance but also to detect distortion in the metology loop.

Electrical properties other than capacitance methods have been attempted in a limited way.

### 4.4.3  *Capacitance as a proximity gauge*

Nevertheless, capacitance methods can be used very effectively for spacing gauging and also for monitoring the radius in holes. Fabricating a set of them together in a given shape can enable the profile of difficult parts to be estimated very quickly indeed. Its use therefore seems less relevant to roughness but more to form measurement and certain dimensional properties.

One point worth noting is that by the very nature of the method an average value over an area is obtained. This is completely different from the stylus methods and some of the optical methods so that the technique does offer the potential for areal assessment albeit of a rather scrambled nature.

## 4.5  Inductance technique for measuring surfaces

One technique has been to measure the roughness of two magnetic surfaces by measuring the inductance between them as a function of separation [209]. Some correlation between the mutual inductance and surface roughness was found but the method is restricted to magnetic surfaces and should only be used as a comparator even for the same manufacturing process.

## 4.6  Impedance technique—skin effect

High-frequency impedance to electrical signals is largely in the skin of the conductor—the so-called 'skin effect'.

By increasing the high frequency it is possible to arrange that the skin thickness is of the order of the surface roughness. In the limit it could be argued that if the frequency is very high the current would follow the 'real length' of profile $R_1 \simeq 1 + \frac{1}{2}\Delta_m^2$. At low frequencies the 'length of conductor' would correspond to unity, so in principle, the increase in impedance could be related to surface slope $\Delta_m$. This assumes that there are no substantial changes in the purely electrical properties of the skin relative to the bulk.

Both of these methods give average results for the surface and as a result it is not possible to obtain any local information.

## 4.7  Other non-standard techniques

### 4.7.1  *General*

There are numerous other techniques to measure roughness. Most of these have been exploratory in nature and made to attack a specific problem. They have therefore not necessarily followed the conventional path

as in stylus and optical methods. Some of the methods have been entirely valid in the sense that the instrument mimics the function being controlled. Such instruments are therefore matched in purpose to just a limited number of applications and cannot be considered to be serious competitors to the stylus and optical methods. Other instruments have been developed to measure roughness on an altogether different scale to that of engineering surfaces. Examples of this are the measurement of roads, runways, ships' hulls, etc. In these cases some completely different constraints may have to be considered such as portability, robustness, waterproofing.

Some examples will be given below. In most cases the theoretical basis of the instrument is not known in detail; its credibility is based on usage.

### 4.7.2   Friction devices

The flexible-blade technique mentioned earlier in section 4.2.4.6, called the mecrin and developed by Rubert, is an attempt to replace the fingernail. A thin blade is pushed along the surface in much the same way as a fingernail test on the surface [84]. As it is pushed the angle of inclination of the blade is changed until at a critical angle the blade buckles. The buckling makes an electrical contact which indicates the angle and by implication some property of the surface such as surface mean slope. The problem here is that the device is just as sensitive to the variations in the adhesive component of friction as it is to the surface. This adhesive component is a function of the material of the probe and that of the surface, so different results are to be expected from different materials having the same finish. This makes the instrument a comparator of limited applicability. It should be used directly as a friction-measuring device rather than indirectly trying to measure geometry! It is not clear whether this has ever been done.

#### 4.7.2.1   The friction dynamometer

The test piece is rubbed by a reference piece attached to a pendulum bob which has been released at an angle of 30°. The damping introduced by the friction between the test and reference pieces eventually brings the pendulum to rest. The time to do this is taken as a measure of roughness. This method is very doubtful in view of the fact that it measures the dynamic friction, not the finish, and this dynamic friction occurs at different speeds.

### 4.7.3   Rolling-ball device

The angle at which a ball of radius $r$ starts to roll down an incline is supposedly related to the roughness [181]:

$$R_p = r(1 - \cos\alpha)$$

(4.288)

where $R_p$ is the average peak measured from the mean line. This is a technique related to the static friction properties. It seems not to depend on material properties, but this is difficult to justify.

Other methods have been used to estimate roughness using balls (see Thomas [206]) but the relationship to the roughness values is tenuous and not commercially viable.

### 4.7.4   Liquid methods—water

Some methods have been used for measuring large-scale roughness which would be completely unsuitable for surface finish measurement on the microscale. Some physical properties do not scale down linearly. Water flow is one. Moore [210] proposed a scheme whereby a cylinder which is open ended and has a compliant annulus at the end contacting the surface is filled with water. The time taken to allow a given volume of water to escape is supposed to be a measure of the surface finish.

Such a method and similar methods involving the filling up of surface voids with various kinds of filler are not really practical.

### 4.7.5 Liquid methods—oils

Bickerman [211] has reported various ways of measuring surfaces using liquids. One uses an oil spot of known volume artificially spread on the surface. The ratio of the spread area to the oil droplet volume is a measure of the roughness. In another method the oil drop is timed as it flows down an inclined plane, and in yet another method the surface is covered with a non-volatile liquid. The surface is held vertical for some time until the liquid is stable and not moving down the plate. At this point it is reckoned that the thickness of the liquid left on the surface is about $1-2 \times R_q$ in value. This, like the others, is interesting but hardly practical.

### 4.7.6 Pneumatic methods

This technique, like the capacitance method, has been in use for many years [10], again at about the same time as serious work began on surfaces. The theory was developed in 1937 and based on that of leakage [212]. Two possible methods could be used (figure 4.237), one dependent on air flow, the other on the back pressure generated in the chamber by means of the reluctance of the air to escape owing to the inhibiting effect of the surface finish (figure 4.237(*b*)).



**Figure 4.237**  Pneumatic method.

A circular nozzle was originally used but there have been all sorts of variants. As with the capacitance probe the method was originally used as a proximity sensor. Under certain conditions, however, the back pressure can be made to have a linear relationship with surface finish—supposedly independent of lay.

Results so far indicate that useful results could be obtained in the range of 1 to 20 $\mu$m, the roughest surfaces where the presence of long wavelength is less likely to upset the readings.

Practical problems of importance include the proper drying of the air. Water condensation can not only rust the part but also give an artificial reading. Spatial cut-off valves are dependent upon the size of the orifice aperture. Various options are available simply by changing the aspect ratio of the probe. Very thin probes have been suggested to simulate the average over a profile in order to make comparisons with the conventional stylus easier.

Flow-type methods have been used for measuring rough surfaces [213]. Providing that the surface had minimum curvature on it and regular lay the change in flow could be related to the texture. However, when the workpiece is moved relative to the probe a dynamic pneumatic effect becomes evident. It was estimated that these effects only became significant at about 30 m s$^{-1}$ for surfaces of about 10 $\mu$m $R_a$ finish.

An exact analogue of the Wheatstone bridge has been developed for measuring surface finish using a bellows, two fixed fluid resistors and a variable resistor which is accurately controllable and measurable [214]. The latter is the key element in the system and is based on a variable orifice made by the movement of a curved lens over the air supply. This is an ingenious technique which definitely extends the range of pneumatics. Typical values range from 0.2 to 6 $\mu$m $R_a$. However, as always the system needs to be calibrated and for any reading has to be 'tweaked' which is not conducive to in-process gauging.

It appears at present that these methods have some potential in the in-process fields, especially at slow speeds, but only as comparators. It will be interesting to see how they ultimately develop relative to optical methods. The indications are that they may specifically have an advantage in very hostile environments, principally because of the fact that the measuring medium (i.e. air) removes debris and coolant in much the same way that the stylus does. Note that displacement measurement at a remote site can be made using sonic orifices (when the flow is independent of downstream pressure).

Also, as in the capacitance method, specially shaped probes can be made to control diameter and cylindricity by means of two opposite measuring nozzles for use in cylindrical bores and shafts.

The pneumatic element in the measurement need not be direct; it can be part of the transducer rather than the sensor. For example, the sensor can be a ball which touches the workpiece, the position of the ball being sensed rather than the surface itself. This technique is used for high sensitivity rather than normal usage. A similar application of the indirect method has already been mentioned when considering transducers with nominally zero damping and constant force versus distance characteristics.

### 4.7.7 Thermal method

Other methods include thermal comparators based upon the amount of heat lost by a preheated ball which is making contact with the surface relative to a similar ball which does not contact. Thermocouples are connected to both and the differential temperature measured. This is related to the surface finish [215].

From the point of view of instrumentation this technique is not likely to be practical. In fact the number of unknown factors involved is so great that the only useful application is likely to be functional and that would be in thermal conductivity tests. One alternative to optical methods is the possibility of using ultrasonics to measure roughness.

### 4.7.8 Ultrasonics

A number of investigators have attempted this approach, principally Berry [216] and various papers by Quentin and colleagues (e.g. [217, 218]). Basically the idea is that of investigating the characteristics of echoes as modified by the surface. Berry's work derives from the investigation of subglacial topography of polar regions.

It is necessary to note how the time of the first echo (first return) varies as the source/receiver point is moved in a plane above the surface. Of course the roughness affects the echo. It is common for echo-sounding experiments to employ quasimonochromatic pulses. Berry considered the whole analysis from a statistical point of view making use of scalar diffraction theory similar to Beckmann and Spizzichino's, making use of the ill-founded Gaussian correlation function. However the whole approach is remarkably parallel to the optical development except for the paint source and the fact that the sound waves are pulsed. The fundamental surface parameters considered were the autocorrelation function and the RMS value of the surface. Other factors that have to be taken into account are the pulse/wavelength ratio and the height of the source from the surface $z$.

The conclusions are based purely upon the theoretical analysis and not upon practical considerations. Some of these conclusions are that half of the echo power arrives after a time $T + \frac{1}{2}P$, where $P$ is half the pulse width. This half-power time $T$ is given by

$$T = 2.77 z R_q^{2} / c L^2$$

(4.289)

or

$$T = \frac{3z\lambda^2}{16\pi^2 c L^2}(1 - \delta)$$

(4.290)

where $c$ is the propagation speed of the waves, $L$ is the spatial extent of the surface finish (to be taken as about equal to the average wavelength derived earlier), $\delta$ is the coherence length of the source and $\lambda$ is the wavelength of the source.

Equation (4.289) relates to surfaces which are Gaussian, whilst (4.290) relates to step surfaces shown with reference to fractals (also flat surfaces having varying reflectivity). Obviously from the point of view of comparison the first equation is more relevant to the investigation of manufactured surfaces. It shows more than anything that some estimate of slope is possible. This is by making crude use of $R_q/L$ as a measure of gross slope.

Also information can be derived in the tail of the echo. However, the assumptions for real surfaces need to be clarified and extended to make direct comparisons possible.

It remains to be seen whether it will be possible to utilize this sort of method in practice. Problems will exist because of the natural absorption of the air, which may be a critical factor in industrial environments, especially with mists and debris floating around. One way around this is to immerse the whole object in a liquid thereby getting a much more controllable constitution [217, 218].

In this method use is again made of the diffraction theory approximations of Kirchhoff, only this time in sound, so again the thread of the use of far-field diffraction is maintained throughout the modern techniques (figure 4.238).



**Figure 4.238** Ultrasonic methods of measurable texture.

Experiments relating to the backscattered response led to estimates of the $R_q$ height in the range of 6 to 50 $\mu$m with a precision claimed to be 1 $\mu$m. The probe has a centre frequency of about 5 MHz and a pulse of about 0.2 $\mu$s. Obviously in this case the roughness values are large compared with most manufactured surfaces. Also the coupling factor is hardly suited to easy use in industry (unless the whole part is immersed in coolant!).

Extending this same theoretical work using normal incidence rather than a set of angles [219] has shown that not just the $R_q$ value could be obtained but the whole amplitude distribution of the rough surface, so that at a given frequency $f$ the ratio A($f$) of the pressure scattered from the rough surface to that of a plane smooth surface is given by

$$A(f) = \int_{-\infty}^{\infty} p(z)\exp\left(\frac{-4\mathrm{j}\pi f z}{c}\right)\mathrm{d}z$$

(4.291)

where $p(z)$ is the amplitude distribution (assuming no multiple scattering is present and the slopes are low). So, in principle, inverting equation (4.291) yields $p(z)$ from which $R_q$ can be obtained.

As in optics the scattering of sound waves by rough surfaces becomes more diffuse as the ratio of the roughness to the wavelength of the source approaches unity. This is exactly the same problem as in optics when the $R_q$ value is the same as the wavelength of light. The difference here is that the wavelength of the ultrasound can be varied more easily over a wide range—in fact it could take over where the optical possibilities fall short in the larger roughness bands.

*Ultrasonic sensors*

Ultrasonic microscopy was originally developed to look at the internal structure of small opaque objects that do not transmit optical rays. A typical example is a ceramic material. Ultrasonics are especially useful

for investigating the elastic properties of materials at the microscopic level. The method is not an 'echo' technique but a 'wave' technique.

Potentially, ultrasonics are a powerful tool for surface analysis because as a 'wave' it can be scattered from rough surfaces in the same way as optical methods (i.e. Kirchhoff scattering). However, it has the advantage that for suitable angles of incidence the sound waves can penetrate the surface and thereby investigate strain underneath the surface. One problem is the relatively low resolution of acoustic lens of 10–100 $\mu$m at 1–5 GHz; in particular a 13 $\mu$m lens and a 4 .4 GHz frequency.

The basic element of ultrasonics in the microscope mode is a sensor for focusing the acoustic beam. This sensor is usually a zinc oxide (ZnO) piezoelectric film transducer together with an acoustic lens. The beam can be focused down to a spot or a line depending on the application. The other element is a liquid coupler. This can be water or, for better results liquid gases such as nitrogen. The problem lies in finding a coupler with little attenuation. In a liquid attenuation is proportional to the square of the frequency, consequently, when the frequency is very high signals are obscured by noise due to attenuation in the coupler. The only way to reduce attenuation and increase the resolution is to decrease the radius of curvature of the acoustic lens, which is often a sapphire, rod-shaped at the end and capped at the other end with a ZnO transducer [220].

The 'signature' approach to controlling rough surfaces may be a possibility for the future.

The question of whether light diffraction extends into the rougher regions or sound diffraction into the smoother is a fight between two instrument sensor technologies.

### 4.7.9   Summary

The last few possible methods of measuring surfaces have been interesting yet speculative. There are many more which are not sufficiently practical or researched to consider here. However, the main issue in many of the proposed techniques is different from one to another. Normally it is proper to try to predict function. Failing that, it is useful to get quantitative information of the surface roughness about one or more parameters; failing that it is important to be able to get a quantitative 'feel' for the surface detail.

It is this last category where instruments such as electron microscopes, tunnelling microscopes and the like have their *forté*. Balancing the lack of geometrical information is the increased capability to explore the chemical and molecular nature of the surfaces. Obviously in a book like this the emphasis is on geometry, but the underlying requirement has to be a combination of all factors that could influence performance. For this reason some aspects of microscopes will be considered. Optical microscopes do not bridge the barrier between geometry and chemical/physical properties and therefore will not be considered *per se*. It must be highlighted in this context that their usefulness is not questioned; what needs to be considered is their potential for the future examination of surfaces.

## 4.8   Electron microscopy

### 4.8.1   General

There are two types of microscope that rely on the atomic properties of materials to explore surface topography. One is the scanning tunnelling microscope, which relies on tunnelling effects and the other is the electron beam microscope in which an electron beam is directed onto the surface. The former has been considered in section 4.2 because of the great similarities of the pick-up mechanism with that of the tactile stylus instrument. The latter is considered below.

The three factors of importance to be considered in electron beam microscopy are the effective resolution, the way in which the electrons impinge on the surface and the methods of scanning. Other features are to be found in many texts.

Take first the definition of resolution as used for optical microscopes. The resolution is

$$1.22\lambda/\mathrm{NA} \tag{4.292}$$

using the well-known theory. What are $\lambda$ and NA for a typical electron microscope? For a $\lambda$ of 0.55 $\mu$m and an NA of 1.6 the best an optical microscope can expect is about 0.4 $\mu$m.

For electronic measure $\lambda$ is related to the momentum of the electron $p$ and Planck's constant $h'$ by the formula

$$\lambda = h'/p. \tag{4.293}$$

Given that the charge of the electron is $e$, mass $m$, $v$ the velocity and $V$ the voltage (in electron volts) governing the energy required by the electron passing through $V$ volts (1 e$V = 1.602 \times 10^{19}$ $J$), then

$$\lambda = h'/(2meV)^{1/2}. \tag{4.294}$$

Therefore for 20 kV $\lambda = 8.6$pm. Note that this is almost $10^5$ times smaller than that of the typical optical wavelength, which would indicate from equation (4.294) that the resolution could be better by the same proportion.

One could argue that there are electromagnetic waves (e.g. x-rays) which have much smaller wavelengths than those of visible light. The response is yes, but how can they be focused to get an effective large numerical aperture? Apart from de Broglie, who formulated equation (4.293), the next important step in the development of the instrument was the notion that magnetic fields, distributed radially, could focus the beam of electrons. The same applies to electrostatic focusing.

The resolutions for the transmission electron microscope (TEM) and the scanning electron microscope (SEM) are complex and depend on the particular configuration, but the NA is of the order of 100 times less than that of a typical microscope. This results in a spot size which is about 0.6 nm for the TEM (the NA would be about $5 \times 10^{-3}$ rad). For the SEM the resolution is usually poorer by a factor of 3 or 4. Nevertheless, the depth of focus of the SEM, which is probably its most attractive feature, is in the millimetre rather than the micrometre range, simply because of the low NA value, and remembering that the depth in focus is given by

$$\lambda/(\mathrm{NA})^2. \tag{4.295}$$

Obviously when compared with practical surface measurement the problem with the TEM is that the specimen has to be translucent to the wavelength involved, otherwise it does not work.

The problem with the SEM and TEM is that they have to be worked where the specimen is in a quite considerable vacuum and, furthermore, for the SEM the output does not necessarily depend on the actual geometry of the workpiece. (For a good introduction into simple electron microscope theory see references [221–223].)

In order to understand the mechanism better, a picture of the penetration of electrons into the surface is useful. Figure 4.239 shows that many things happen as far as scanning electron microscopy is concerned. It is the secondary electrons which are usually picked up (figure 4.240) by a detector situated on the side (figure 4.241).

The secondary electrons result from an interaction of the incident electrons with the bulk atoms releasing electrons of lower energy (typically less than 50 eV). However, to get out of the surface they must have an energy greater than the work function of the material. They tend to be formed at about 100 nm into the surface. The important point to note is that they do not come from the surface but underneath it.

The local geometry of the surface is also very important because it affects the secondary electron yield, $se$ (figure 4.242). Thus

$$se(\theta) \propto k(1 - \cos\theta) \tag{4.296}$$

where $k$ is a constant.

**Figure 4.239** Reaction of electrons with solids.



**Figure 4.240** Incident electron beam normal to surface.



**Figure 4.241** Position of detector in the SEM.

The result is that any point of a surface which is convex tends to produce more electrons than one which is concave or flat. This produces an effective enhancement of the surface feature over what would be seen by a stylus. This visual effect is often misleading to an investigator. Many displays use the $z$ modulation format [222] in which electron yield is presented as a profile across the screen.

Although SEM relies on the secondary electrons to produce the image it is possible to make use of the backscattered primary electrons. This technique is called electron backscatter diffraction (EBSD)

**Figure 4.242** Incident electron beam at angle to surface.

and is an interesting technique because it can be used in examining surface microstructure as in grain boundary, or plastic strain and also crystallographic properties. Microstructure refers to the fine scale external geometry of crystals whereas crystallography refers more to the internal structure and symmetry within each crystal. Care needs to be taken here because sometimes crystallographic detail is referred to as 'texture' which in this book is the marks left by the machining or forming process. EBSD is therefore to be regarded as being useful for surface and mainly subsurface structure, complementing ultrasonics and x-ray methods.

A typical image from EBSD is that of a pattern (formed by diffraction from the lattice). This usually comprises of sets of lines called 'Kikuchi' lines. Each pair represents a plane in the crystal. The spacing between line pairs is inversely proportional to the interplanar spacing given crystallographic direction.[224, 225].

### 4.8.2 *Reaction of electrons with solids*

Large yield looks like a peak, low yield looks like a valley. The problem here is that people think that this intensity profile is a geometric profile of the workpiece. In fact it is not. Experiments have been carried out to get some idea of the correlation between modulation of the SEM and profile methods [178]. Definite correlation exists between the $z$ modulation and surface slopes as would be expected from equation (4.296). To get an accurate estimate of correlation it is fundamental to make sure that the stylus and the electron beam are traversing in the same path. Correlation values of 0.8 have been obtained with surface slope. This can be clearly seen in figure 4.243. Notice how the location spots get differentiated. In fact there is also a correlation with curvature of peaks. Better correlation is possible using a multiplicity of detectors around the primary beam.

Quite small angles on the surface can be detected. The dependence on specimen inclination can be defined as $C$ [226] where

$$C = \tan\theta \, d\theta. \tag{4.297}$$

For an incident angle $\theta$ the number of electrons given off per unit solid angle appearing at $\beta$ is given by

$$N(\theta, \beta) = \frac{N(0)}{\pi} \cos\beta \tag{4.298}$$

letting $\alpha = \theta + \beta$ where $N(0)$ is the number at normal incidence. This becomes

$$N(\theta, \alpha) = \frac{N(0)}{\pi} (\cos\alpha + \sin\alpha \tan\theta). \tag{4.299}$$

Typically the minimum slope that can be detected is about 0.5°.

The scanning is achieved particularly easily by deflecting the electron beam, as shown in figure 4.244.

| Differences within process | | Similarities between different processes | Differences between different processes |
|---|---|---|---|



Face turned fine    Circumferential turned fine    Slab milling smooth    Electrodischarge machining steel

Face turned medium    Circumferential turned medium    Ground fine    Plasma cut monolith AA124

Face turned rough    Circumferential turned rough    Polished surface cast-iron

**Figure 4.243** SEM pictures



Source

Focusing coil

Scan deflection coil

Detector

**Figure 4.244** Schematic layout of the SEM.

### 4.8.3 Scanning electron microscope (SEM)

The scan can be quite accurate and because of this the SEM is often used for measuring length between or of features on the surface. It is not quite as easy to measure quantitatively other geometrical properties, as has been indicated earlier on. Nevertheless there is a potential way of getting quantitative information about the surface and this is by means of stereographic viewing [193]. In this a micrograph is taken of the specimen. The specimen is then titled slightly, say by 10°. Superposition of the two images in a stereo viewer then gives a correct 3D representation of the surface. If the sample has been tilted towards the collector there can be an appreciable difference in working distance, which causes a different magnification from one part of the specimen to the other. This is most noticeable at low magnifications. The actual size of the feature can distort the measurement of lateral distance as indicated in figure 4.245 [227]. The height of the feature is measured from the parallax shift as the tilt is made.

As the object, shown in figure 4.246, is tilted the parallax $P'$ is related to the height of the feature $h'$ and the length $l$ by

$$P = M[l(1 - \cos 2\delta\theta) + h \sin 2\delta\theta] \tag{4.300}$$

or

$$h' = \frac{P'}{M \sin 2\delta\theta} - \frac{l(1 - \cos 2\delta\theta)}{\sin 2\delta\theta}. \tag{4.301}$$



**Figure 4.245** Distortion of lateral distances in the SEM.



**Figure 4.246** Effect of parallax.

Even this formula is not very accurate for measuring the height normal to the plane of the surface [228].

However, taking sufficient care with measurements can enable height information to be obtained, although it can be messy, especially in the two-dimensional and three-dimensional cases.

It seems ironic that the one feature of the SEM that enables such vivid pictures to be taken, namely the large depth of focus, is the cause for ambiguity in height measurement. (In the same way that the resolution is of the order of a nanometre the depth of focus is of the order of a millimetre.)

Various techniques have been introduced to get better pictures with the SEM. One is to use split detectors and to look at the backscattered or reflected electrons rather than the secondaries. In the same instances the apparent image has been improved. The improvement in contrast has misled some investigators into thinking that they have the possibility of stereoscopic viewing without tilting specimens because of two detectors, but in practice this is not true.

One method has been to use two detectors symmetrically placed relative to the input beam (figure 4.247). Using these two electrodes the sum of the signals $A + B$ is even, whereas the difference $A - B$ is odd and linear over $\pm 75°$. These results emerge using a latex particle as calibrator.

For calculations of the gradient on the surface it is found that

$$\tan \theta = \frac{k(A^2 - B^2)}{(A_m + B_m)^2}$$

(4.302)

where $k$ is a constant which can be evaluated by calibration and $A_m$ and $B_m$ are the detector outputs for a flat surface. Although equation (4.302) is empirical it shows an accuracy of better than 2% over $\pm 75°$.

The profile is best found by integration of the derived slope, the result of which shows that providing the slopes on the surface are not bigger than $\pm 15°$, the profile compares well with the profile obtained by the stylus. Obviously therefore the SEM cannot measure profile directly, but by having a multiplicity of detectors very creditable results could be achieved providing the specimen does not charge up, in which case it has to have a thin film of gold evaporated onto it.

These problems of getting a profile should not detract from the usefulness of the SEM. Over quite a wide range, providing that the deflection coils for charging the electron path do not distort or saturate, vivid



**Figure 4.247** Arrangement of Suganuma.

pictures of surface structure can be obtained. Even when the object charges up, the resultant edge enhancement can magnify the visual impact to good effect.

### 4.8.4 Transmission electron microscope (TEM)

Transmission electron microscopes are not as important in surface geometry investigations as SEMs because most engineering surfaces are opaque. However, it is possible to get some information using replication techniques [229, 230]. In most applications the surface is shadowed with a platinum/carbon mixture and a carbon replica is made. The whole technique depends for contrast on the shadowing with the metal atoms. Large angles are used to throw the asperities into relief; a typical value would be 75° from the normal. The height resolution is approximately proportional to the lateral resolution by the tangent of the shadowing angle which can be as low as 5–10° from the surface and would in theory be 0.3–1 nm. It has also been pointed out that an even higher resolution is possible if tantalum/tungsten shadowing is employed [230].

Behind the purely quantitative information that can always be obtained from replicas in the TEM, some estimates of surface roughness parameters have been made. Claims about getting RMS and autocorrelation function information have been made [230] based on the use of microdensitometers.

The real issue with this and other replication methods is whether they can be regarded as practical in the industrial sense. It should probably be a general rule that replicas should not be used unless necessary. One example is that of styluses. Early SEMs proved to be unsuitable for measuring styluses because, just where the information was needed, that is at the tip, it became charged up. This had the effect that the picture became obscured. The only solution was to use the TEM [229]. In this case the tip was pressed into glass, in fact a microscope slide. The slide was coated with carbon and then shadowed with gold or platinum. Then the carbon replica was floated off and examined under the TEM. This method achieved results far behind those possible using the SEM. For example, simply lowering the accelerating voltage to stop the charge-up only resulted in a poorer signal-to-noise ratio and so defeated the object.

There are also two-stage replication methods for use with the TEM in which some plastic material is softened by heat or by partial dissolution (Formvar, cellulose acetate, polystyrene, etc) and pressed onto the part. This is then removed, and carbon or silicon monoxide is deposited on it. The plastic is dissolved leaving a thin-film replica which is shadowed on the face carrying the structure. This is done if the component material will not retain the thin carbon or silicon film.

### 4.8.5 Photon Tunnelling Microscopy (PTM)

This subject is placed here because of the direct comparisons with SEM. Tunnelling of electrons through a potential barrier is the basis of the now well-known scanning tunnelling microscope and its derivatives, first used by Binnig *et al* [48, 49].

The optical analogue to the STM employing photon tunnelling followed in the form of the scanning near field optical microscope (SNOM). It has been pointed out by Guerra [231] that photon tunnelling effects have been observed for many years in spectroscopy, optical wave guide coupling and holography. Photon tunnelling in microscopy predates the STM by almost two decades in the frustrated total internal reflection microscope [232].

In the absence of the probe used in SNOM the lateral resolution reverts to the usual diffraction limit but high vertical resolution remains. The 'advantage' is that the whole field is preserved and scanning is not needed. The compromise is between 'high lateral resolution probe + scan' and 'diffraction limited lateral resolution + full field.'

Briefly the theory as presented simply by Guerra is as follows. Consider figure 4.248.

When a photon in medium 1, with refractive index $n_1$ is incident to medium 2 $n_2$, (lower value index) at a critical angle given by

**Figure 4 .248** Evanescent field, associated with photons tunnelling into medium 2, as a continuation of the parent field in medium 1.

$$\theta = \sin^{-2}\left(\frac{n_2}{n_1}\right)$$

(4 .303)

or greater, it will tunnel into the rarer medium producing an electric field. As no energy is transferred the photon experiences total internal reflection and the field in medium 2 is called evanescent. The amplitude of the evanescent field decays exponentially normal to the boundary.

$$E_{ev} = E_0 \, \exp\left(-\frac{z}{d_p}\right)$$

(4 .304)

$E_0$ is the amplitude of the electric field in medium 1 at the boundary, $d_p$ is the depth of penetration at $\frac{1}{e}$

$$d_p = \frac{\lambda_1}{2\pi\left(\sin^2\alpha - \left(\frac{n_2}{n_1}\right)^2\right)^{\frac{1}{2}}}$$

(4 .305)

where $\lambda_l$ is the wavelength in the denser medium $1 = \frac{\lambda}{n_1}$, $\alpha$ is the incident angle.

In the presence of a third medium 3 in the figure of refractive index $n_3 > n_2$ the tunnelling photons will couple, thus frustrating the total internal reflection. A linear change in the separation $z$ causes an inverse and exponential change in the energy transfer so that if the surface of medium 3 is a rough surface the light reflected at the boundary in medium 1 is modulated according to the roughness. Hence the microtopography is transformed into a greyscale image. This greyscale image has to be translated into height information by means of a calibration surface of known geometry; usually a hemisphere [233].

Figure 4.249 shows the usual configuration. Contact or near contact between transducer plane and specimen causes frustrated total internal reflection by photon tunnelling. This shows up as black spots on the image. This method has been used to ascertain the contact between gears by means of perspex two disc machines.

**Figure 4 .249** TIR mode of microscope.

While the vertical resolution of photon tunnelling microscopy is about the same as the scanning electron microscope (SEM) the shorter wavelength of the electron yields superior lateral resolution. The photon tunnelling microscope does, however, have one advantage. This is that dielectrics can be viewed directly without coating or shadowing. Samples which release gases and destroy the vacuum in SEMs can be viewed at atmospheric pressure with photon scanning methods. Also there is negligible sample damage from the low energy photons.

Although there is reasonable correlation of PTM images with Nomarski, Wyco and stylus methods it does not appear to have the same versatility.

More recent work corroborates the correlation between PSM and the other techniques. There is a fundamental problem however, in that most of the correlation is in lateral positioning and spacings of detail but not convincingly in the normal direction. Many scanning methods presuppose that the important information is in structure and boundaries.

## 4.9 Merit of transducers

### 4.9.1 Comparison of transducer systems

In this chapter methods of quantifying surface geometry by a number of different types of pick-up have been discussed. In what follows there is a comparative survey of various systems for converting the pick-up information into a usable signal. Before this it is informative to place surface metrology instruments in perspective relative to other displacement instruments.

The basic system is shown in figure 4.250. Each one of the blocks represents a source of noise and has its inherent limitations. Often the mechanical noise can be reduced by good mechanical design. This is less true for electronic noise.

A few representative examples of transducer systems will be compared including that of the well-known stylus method. So that the field is not too wide only systems capable of measuring up to 50 mm will be considered. The basic sensitivity, the noise and the resolution will be considered in detail for a few cases at the end of the comparison to indicate how the values are arrived at.

Because of the continuing improvements in performance the quantitative data is likely to change in time. Therefore the data should be used for comparisons between transducers rather than absolute values [234].

Table 4.9 shows some of the applications of such transducers. It can be seen that surface metrology occupies quite a large proportion of the total applications. Values of the range to resolution shown are only typical and do not necessarily represent the best possible values. Systems such as the STM and AFM have been omitted as they are in such a period of flux at the present time.



**Figure 4.250**   Basic instrument system.

Figure 4.251 shows how the key comparison parameter range-to-resolution can be built up from the sources of error in the system. There are obviously other considerations such as cost, actual size of device, frequency response, etc, which are not shown for clarity.

There are many considerations to be taken into account when specifying transducer properties such as linearity, stability, hysteresis, etc. The relative importance depends on the application. In creep, for example, long-term stability would be paramount.

There are a few basic rules concerned with transducers. One point is that each of the components of the transducer has its own range-to-resolution and frequency response and all contribute to the overall performance of the transducer. Ideally each block should be matched with its neighbours in such a way that the overall performance can be optimized. Shortcomings in one element can often be remedied in another but it is not usually recommended. This process is sometimes called equalization. Most problems usually occur in the mechanical coupling and in the conversion.

For the mechanical registration of surfaces the mechanics commonly consist of a pivoted lever or beam in a 'side'-acting mode where the transducer is at the opposite end of the beam from the sensor. The whole arrangement—the force on the sensor, the return spring on the beam and the force in the transducer—is kept in near balance during the operation. The light construction and weak spring force necessary to resolve microtopographic data to the value of angstroms render this configuration unsuitable for large displacements, for example of values of greater than a millimetre.

The limitations of deformation, hysteresis, friction, etc, only allow a certain range-to-resolution value which, even for good instrumentation design, rarely exceeds 40000:1. Given this figure and accepting the resolution value of, say, 2 Å gives a range of 80 $\mu$m.

**Table 4.9** Application of displacement transducers below 50 mm.

| | Maximum range-to-resolution ratio | Range (mm) | Resolution ($\mu$m) | Frequency response (Hz) |
|---|---|---|---|---|
| **Direct** | | | | |
| Dimensional metrology: | | | | |
|   Length | $5 \times 10^5$ | 50 | $1–10^{-1}$ | <1 |
|     height | | | | |
|     width | | | | |
|   Diameter | $5 \times 10^5$ | 2-50 | $10^{-1}$ | <1 |
|   Thickness (thin films) | $10^5$ | 0.1 | $10^{-3}$ | <1 |
|   Position | $5 \times 10^6$ | 50 | $10^{-2}$ | 10 |
|   Angle | $2 \times 10^4$ | 2 | $10^{-1}$ | <1 |
|   Level | $5 \times 10^3$ | 0.5 | $5 \times 10^{-1}$ | <1 |
| Surface metrology: | | | | |
|   Texture | $4 \times 10^4$ | 0.2 | $5 \times 10^{-3}$ | 100–300 |
|   Roundness | $2 \times 10^4$ | 0.5 | $2.5 \times 10^{-2}$ | 50 |
|   Straightness | $2 \times 10^4$ | 2 | $10^{-1}$ | 10 |
|   Flatness | $1 \times 10^4$ | 1 | $10^{-1}$ | 10 |
|   Contours and surface Deformations | $2 \times 10^4$ | 20 | 1 | 10 |
| Temporal metrology: | | | | |
|   Thermal expansion | $1 \times 10^6$ | 1 | $10^{-3}$ | <1 |
|   Vibrations | $1 \times 10^5$ | 1 | $10^{-2}$ | $10^2–10^6$ |
|   Creep | $1 \times 10^6$ | 1 | $10^{-3}$ | <1 |
|   Tilt | $1 \times 10^5$ | 0.1 | $10^{-3}$ | 10 |
| **Indirect** | | | | |
| Position: | | | | |
|   Force, velocity | $5 \times 10^5$ | | | |
|   Acceleration | | 50 | $10^{-1}$ | $10^2–10^6$ |
| Strain: | | | | |
|   Displacement, velocity | | | | |
|   Acceleration, force | $10^4$ | 5 millistrain | 0.5 microstrain | $10^2–10^6$ |

The ordinary surface instrument that has been developed over the past decade has typically a resolution which is much poorer than angstroms and is more likely to be in the region of nanometres, giving a corresponding range of more like 200 $\mu$m. This is usually sufficient to allow for the measurement of roughness in the presence of form.

In some exceptional circumstances the range-to-resolution ratio has been extended so much that a 4 mm range with a resolution of 5–10 nm has been achieved: a range-to-resolution value of more like $10^6$. Under these circumstances, form as well as roughness measurement is routine. The form is worked out from the whole data by means of fitting curves. An example of this has been given in Chapter 2. A realization of this enhanced type of coupling has only been possible by using carbon fibre beams, a true ball pivot and gravity stylus force with no return spring.

In the case where there is no contact, such as in an optical sensor, the transfer is made by electromagnetic radiation. In order that the dimensions of the illuminated area meet the smallness required, it is necessary, as in the case of all 'followers', to control the position of the objective lens by means of a closed-loop servo system. Again, the typical well-designed servo position control can achieve 40000:1. The complication of the closed-loop servo in effect replaces that of the lever, spring and pivot mechanisms. The transfer of energy optically

**Figure 4.251** Performance features of metrology instruments.

raises the problem of optical path stability in air if resolutions better than a few nanometres are being sought, in which case a heterodyne system is used in which the reference path has the same optical path.

### 4.9.2   Types of conversion and transducer noise

#### 4.9.2.1   Limitations

The transducers of interest are those in which an energy flow either to 'sinks' or receiving elements is achieved or, sometimes, a potential energy difference is set up. The greater the energy gradient the better the possibility of achieving a high signal-to-noise ratio in the transducer.

In general the limitations are systematic and random errors. The former are produced especially in wide-range transducers when the displacement-energy field change law becomes non-linear, due perhaps to the gap between the stator and armature being non-linear or suffering end effects.

What have been termed microsystematic errors can often limit resolution. In electromagnetic or capacitative transducers these errors are constituted in imperfect geometry of the magnetic circuit elements or electrodes. This would include flaws, irregular edges and surfaces which, for this type of transducer, can be avoided by careful manufacture. In optical transducers the lenses and prisms are necessarily of a very high standard and such problems rarely occur. However, any electro-optical devices converting light energy into electrical energy are subject to pinholes and pits of the order of micrometres in dimension. Also the conversion efficiency is uneven due to inhomogeneities. A method of obviating this problem is to avoid scanning light across these photodetectors. Instead the light pattern is maintained stationary by collecting the light at pinhead-sized photodetectors (for which the uniformity is excellent) by placing them at the foci of segmented collecting lenses. By this means the light pattern is caused to vary over glass lenses of which the transmission is virtually free from microsystematic error. It seems possible to compensate for the systematic errors of a

transducer by means of a look-up table in the subsequent data processing unit. This linearizing is sometimes carried out, but often complications of drift, the arbitrary inputs of offsets for operational convenience and the changes of transducer, parameters by servicing or usage limit its scope. Microsystematic errors are more susceptible to instability and complicated noise equalization requires major computation. In fact the boundary between microsystematic errors and random errors becomes blurred. The microsystematic errors give a noise-like pattern to the law of the transducer, affecting its scale value at points over the measurement. A transducer subjected to a ramp input converts the scale noise into displacement noise. It is a basic requirement that the instrument does not develop such noise when displacements of this type occur with tilt and so microsystematic errors have to be eliminated by design wherever possible and then by the careful control of manufacture.

### 4.9.2.2 Random noise limitation

Associated with the physical quantities and structure of the transducer are time-varying noise signals producing error. Some depend on manufacture but generally these noises are inherent manifestations of nature.

(*a*) *Shot noise*
This noise results from variations in standing electric currents flowing through devices such as resistors, photodiodes, etc, and derives from the finite value of charge carriers. It is akin to Brownian motion in molecules and has a white-noise spectrum.

(*b*) *Johnson noise*
This is a property of resistors and is thermal energy. It is thermodynamic in nature and has again a white-noise spectrum.

(*c*) *Device low-frequency noise*
This results from 'spontaneous' events of breakdown or dislocation at sensitive boundary surfaces, possibly including depletion effects in devices. The spectrum of this noise is usually given as a $1/f$ law indicating step-like changes in energy. Commonly the noise is classified as 'flicker noise', '$1/f$ noise' and 'popcorn noise'. Usually it is specified for the baseband below 20 Hz.

(*d*) *Device broadband noise*
This noise comprises shot noise and Johnson noise internal to a device, but being bound up within a package it is determined by manufacture.

(*e*) *Source noise*
Most transducers require excitation by a source of energy (figure 4.251). This source in generating energy is liable to be affected by extraneous noise. The useful line or carrier frequency of the source spectrum is then accompanied by noise sidebands. For some sources the noise sidebands are inherently low, for example less than 1 part in $10^6$. It is easy to suppress this level with filters. Another redeeming factor is that high-performance transducers are usually operated in a differential mode so that noise like the essential carrier component is suppressed as a common mode rejection [235].

### 4.9.3 Types of conversion and types of transducer

### 4.9.3.1 General

In the above section some of the general limitations of transducers and conversion elements have been discussed from the noise point of view. In what follows some specific types will be identified. After this some will be chosen to examine their sensitivities, resolutions and noise levels.

There are four types of mechanical to electrical conversions:

(1) variable resistance
(2) variable inductance

(3) variable capacitance

(4) voltage or current generation.

In the first three, the displacement is converted into another physical quantity. An excitation voltage is required to measure this quantity. In (4) the displacement is converted directly into current or voltage. Signal processing is then required to get the required accuracy.

There are also four types of mechanical to optical conversions in use:

(5) variable light intensity

(6) variable periodic light intensity (moiré)

(7) variable periodic light intensity using an interferometric system

(8) variable light intensity using holographic and speckle methods.

All these have been discussed earlier as a means of directly estimating the value of the surface geometry. Here they are used instead as an integral part of a transducer system. In (5), (6) and (7) the intensity variation is usually analogue. This variation is used to generate electrical currents in one or more photodetectors.

These listed types above do not include some of the more exotic types such as the magnetostricture effect, the Hall effect, etc. A summary of the advantages and disadvantages have been discussed by Sydenham [236] and Neubert [237].

A few general comments will now be given on these conversions (see table 4.10). In surface metrology the conversion type is usually inductive, capacitive or optical. Resistive or generative are also used.

### 4.9.3.2 Variable resistance

Resistive conversion is used most often in some form of strain gauge in a bridge circuit and rarely for the direct measurement of surfaces. For displacement as such, a plunger type of mechanical coupling is used and the transducer works on large-resistance charges. Usually the gauge part is of a wire-wound type and a continuous track of small conductive carbon particles held in a plastic matrix. The big problem with this type of transducer is the effect of the wire diameter and temperature sensitivity of the carbon. In cases where there is a small-resistance charge a bridge circuit is used. Because of the temperature sensitivity often a differential mode is employed. The spatial separation of the gauge is made as small as possible so that this limits the temperature gradient between them.

**Table 4.10** Advantages and disadvantages of different types of conversion.

| Conversion | Advantages | Disadvantages |
|---|---|---|
| Resistive:<br>  Wire wound or conductive strip potentiometers | Economical; minimum of processing electronics; large range; high electrical output; small size | Track wiper interface causing: additional frictional force on mechanical coupling, decrease in overall accuracy, increase in heat dissipation; variations in track producing noise and non-linearity errors |
| Strain gauge | Fairly economical; versatile; small size; minimum of processing electronics; adaptable to conveniently sensing displacements in more than one axis | Requires matching of gauges to minimize: temperature effects, general drift; requires adhesion and positioning on mechanical coupling |
| Inductive:<br>  Variable reluctance and coupling of two coils | Small size; high resolution possible; unaffected by changes of environment; positioning of coupling to conversion not critical | Mass of core positioned on coupling increases with range; inaccuracies increase proportionally with range due to complex construction of the coil; electrical frequency response limited to $\frac{1}{4}$ of modulation frequency |

**Table 4.10** *(Continued)*

| Conversion | Advantages | Disadvantages |
| --- | --- | --- |
| Inductive:<br>Variable reluctance and coupling of a coil and a conducting surface, proximity transducer | Non-contact; unaffected by changes of environment, between component and coil; small size | Non-linear law—requires linearizing; calibration can be difficult when used in proximity with different metallic components requiring a target; small range, depending on size; electrical frequency response limited to $\frac{1}{4}$ modulation frequency |
| Capacitative:<br>Variable area | Easy mechanical design, allowing choice of electrode materials to give high stability and low temperature coefficients; high resolution possible; use at high temperatures; small size; large range | Mass of electrode on coupling increases with range; can be susceptible to stray capacitance, 'shielding' being required; accuracy affected by changes in geometry and condition of electrodes; electrical frequency response limited to $\frac{1}{4}$ modulation frequency |
| (Additional for) variable gap proximity transducer | Independent of 'type' of metal' component in proximity is made of; small size; non-contact | Accuracy affected by dielectric changes between electrodes; requires a metallic component in proximity; small range; nonlinear law, requires linearizing |
| Self-generating:<br>Piezoelectric | High output; high upper frequency; fairly simple processing electronics | Non-linear law; high operating force; temperature and humidity sensitive; low-frequency limit, small range, requires high-impedance-matching interface |
| Light intensity:<br>Use of mask between source and detectors | Small and compact | Limited range; mechanical design contains a variety of optical components affecting thermal stability |
| Segmented cell | Very high frequency response; measurement in two axes; high resolution; simple mechanical design | Bulky; possible thermal instabilities; range limited by light-source beam diameter |
| Lateral effect cell | Simple mechanical design; long range; high resolution; high frequency response; measurement in two axes | Bulky; possible thermal instabilities |
| Periodic light intensity:<br>Self-scanned diode array | Digital output; high frequency response; use for in-process non-contact inspection measurement in one or two axes; resolution independent of range | Expensive because custom built; optics and electronics designed to suit application are required |
| Grating system | Digital output; high accuracy over range; electrical zero can be adjusted; resolution independent of range; hybrid or absolute systems possible | Scale grating size increases with range; thermal instabilities due to variety of optical components used; can be bulky; possible count loss if a maximum velocity is exceeded; gratings need protection from environment |

**Table 4.10** *(Continued)*

| Conversion | Advantages | Disadvantages |
|---|---|---|
| Periodic light intensity:<br>  Interferometric system | Digital output; extremely high accuracy; electrical zero can be adjusted; high resolution; resolution independent of range; end reflector size independent of range | Thermal instabilities due to variety of optical components used; can be bulky; possible count loss if a maximum velocity is exceeded; expensive due to coherent light source and precision optics used |
|   Holographic or speckle system | Displacement of the whole surface of a component is shown | Bulky and expensive; susceptible to setting-up errors and external vibrations; displacement has to be interpreted by operator; very limited range |

### 4.9.3.3  Variable inductance

Variable inductance is often used. Because of this an excitation has to be used which, although bringing in an extra variable, can be advantageous because thermal drifts are reduced and, if the frequency is made reasonably high (2–15 kHz), there is a considerable gain in signal-to-noise ratio. Voltage levels of about 10 V RMS are used. Because the range of movement is small, frequency modulation can be used instead of amplitude modulation, but mainly it is AM which is used because of the ease of getting components that are linear with amplitude.

### 4.9.3.4  Variable reluctance

Variable reluctance transducers utilize a ferromagnetic case which moves axially inside a coil that is often specially shaped. The shaping is used when the gauge is built to act in a side-acting mode with a lever or beam movement rather than the plunger movement. In another sort there is a variable air gap between a ferromagnetic plate and a coil. Charges in the gap change the reluctance. In the former type a differential mode is most often used.

In general linearities down to 0.1% of full scale can be achieved with clever design of the core and coil assemblies. The core moving in a coil is often used in surface metrology as the key element in the transducer. The variable gap technique is used more often in displacement transducers as proximity gauges. They are non-linear and have to be calibrated. The core and coil method is also being used to provide a constant force for moving translation stage linearity, especially if the stage is a flexible hinge type. One example is in the x-ray interferometer technique.

Differences in dielectric constant between the transducer and the surface in proximity gauge applications such as oil films and humidity do not often affect the measurement accuracy. One active and one passive coil are used to minimize temperature effects.

Single-coil transducers can be used with frequency modulation to obtain direct digital information from the gauge. An example would be to obtain a 200 $\mu$m range with 0.1% linearity and a resolution of 10 nm.

### 4.9.3.5  Voltage or current generators

When two coils are used, this arrangement becomes the basic LVDT (linear voltage differential transformer). The output performance is similar to that of the variable reluctance 'push-pull' transducer. The oscillator and demodulator circuitry can be coupled integrally with the conversion and mechanical coupling, usually at the expense of accuracy and temperature dependence.

The coupling can be between a coil and a conducting surface instead of between two coils. In this form magnetic flux lines from the transducer pass into the conductive surface producing eddy currents that react back on the coil impedance if the modulation frequency is high enough. Hence the proximity of the surface can be derived from an impedance measurement. Eddy currents only penetrate the surface skin to about 50 $\mu$m typically. This method is non-linear and therefore requires calibrating.

### 4.9.3.6    Variable capacitance

Variable capacitance methods require excitation. In general the voltage is a factor higher than for the inductive type to get a good signal-to-noise ratio. As in the case of inductance there is more than one form. The *variable area* types are invariably 'push-pull' in operation and often use plunger-type couplings. They are inherently linear over a wide range due to the simplicity of the mechanical design. Linearities of 0.1% are often achieved. This method lends itself well to the null-seeking method of detection used in 'follower'-type instruments and incorporating closed-loop control.

The *variable gap types*—in differential mode—also give linear response but are used more for their greater sensitivity [238, 239]. Jones used this type of transducer in tiltmeters where stabilities of the order of $10^{-6}$ mm per day were required and high resolutions of $5 \times 10^{-12}$ mm with a 1 Hz bandwidth. In all capacitive methods variations in the dielectric have undesirable effects on the performance.

### 4.9.3.7    Self-generation voltage and current

This is achieved using piezoelectric crystals and moving-coil elements. Both are used dynamically because of leakage effects. Their chief usage is in small-range surface metrology instruments. The fact that the transducer itself is leaky can be beneficial because it can be used, in effect, as a high-pass filter of single stage. This obviates the need for eliminating tilts on the workpiece. The problem is that it is not possible to get a true profile of the surface unless some form of integration is used and, even then, the exact cut-off frequency has to be known.

### 4.9.3.8    Photodetectors

Light is often used in the conversion. In its simplest form a photodetector only is used with a slit attached to the mechanical coupling. This slit hinders the passage of light from a collimated light source, perhaps a photodiode, onto the photodetectors of which there are usually two. The difference in light falling onto the photodetectors is a measure of the position of the slit which in turn is a measure of the position of the plunger or beam. The whole arrangement can be made quite linear (~ 0.3%). From the sum total of light falling on both detectors a control voltage is derived to ensure constant source brightness.

It is possible to dispense with the slit by using a special purpose segmented photodetector in the form of a dual (single axis) or quadrant (two axis). Also, photodiode arrays are now common incorporating 512, 1024, etc, diodes in line spaced by about 10 $\mu$m or even less. Arrays or matrices of photodiodes to cover line profiles or spectra and two-dimensional spectra are being used as well as cameras.

Light methods using gratings as in a moiré system with scale and index gratings of 10 to 40 $\mu$m are in use with transducers using the plunger type of mechanical coupling.

The displacement is incremental in terms of a count recorded by movement of the scale grating. Loss of count can occur if the light is interrupted or the maximum count rate exceeded and the digital circuitry gets 'locked up'. Overall accuracy is usually given but linearity is inherent.

Systems containing three gratings have been used. These give several advantages over conventional systems and double sensitivity as well as relaxed tolerances of the gap between gratings.

Interferometric elements have been used effectively in transducers as the conversion from mechanical to electrical properties. This usually involves the design of a miniature interferometer to get the overall size to about that of a conventional analogue transducer and acting in a side-acting mode. In this case, the interferometer has to be insensitive to angular movements of the test arm. Resolutions of nanometres are common, with ranges of millimetres yielding a range-to-resolution ratio of the order of $10^6$ together with a frequency response of about 100 Hz. Heterodyne transducers have been discussed in which two paths are used, replacing the mechanical coupling altogether.

Speckle systems like the interferometers, have been used in the non-contact mode to bypass the mechanical coupling to the surface. So far there has been no method in which the optical conversion using speckle has been connected with the mechanical link to the surface. Some technical data for typical types is given in table 4.11.

**Table 4.11** Examples of typical transducers with respect to performance.

| Conversion | Type of conversion and mechanical coupling | Range (mm) | Linearity % FS | Repeatability (μm) | Temp. coeff. (°C$^{-1}$) | Temp. range (°C) |
|---|---|---|---|---|---|---|
| Resistive | Wire-wound potentiometer, plunger with and without spring loading | 50 | 0,3 (0.1) | 40 | 0.005 | −15 to 80 |
| | Conductive strip potentiometer; plunger with or without spring loading | 50 | 1.0 (0.1) | 5 | 0.01 | −15 to 80 |
| | Strain gauge; spring-loaded plunger acting on cantilever | 10 | 0.4 | 40 | 0.1 | −10 to 60 |
| Inductive | Var. rel. or LVDT + { spring-loaded plunger | 100 | } < 0.5 { | 0.1 | } 0.005 | } −40 to 80 { |
| | spring-loaded beam | 4 | } (0.1) { | 0.1 | } 0.005 | |
| | | 0.2 | 0.3 | 0.1 | 0.005 | −5 to 55 |
| | LVDT with internal oscillator and demodulator, spring-loaded plunger | 50 | 1 | 0.2 | 0.01 | −20 to 60 |
| | Variable reluctance type in proximity to ferromagnetic component | 5 | 5* | – | 0.2 | −40 to 150 |
| | Eddy current type in proximity to ferromagnetic component | 50 | 2* | 5 | } 0.2 | } −20 to 90 { |
| | | 1.8 | 0.4 | 0.2 | | |
| Capacitive | Variable area; spring-loaded plunger | 50 | } 0.1 { | 0.1 | 0.001 | −40 to 100 |
| | | 5 | } (0.01) { | 0.1 | 0.001 | |
| | Variable area with internal oscillator and demodulator, spring-loaded plunger | 50 | < 0.2 | 0.1 | 0.005 | −20 to 70 |
| | Variable gap, one electrode in proximity to metallic component | 5 | } 1 | – | −0.1% | 0 to 40 |
| | | 0.05 } | | | | |
| | Variable gap, fringing capacitance of two electrodes in proximity to metallic component | 1.5 | 1* | – | 0.001 | −40 to 600 |
| Generating | Piezoelectric element attached to beam | 0.15 | – | – | – | 0 to 50 |
| Light intensity | Mask between source and detectors attached to spring-loaded beam | 0.05 | 0.3 | 0.1 after warm up | – | −5 to 55 |
| | Laser, lateral effect cell for two-axis measurement attached to component† | 105 | ±5 | – | 0.5 | } 0 to 50 |
| | | 5 | ±5 | – | 0.5 | |
| | Miscellaneous technique, laser, diffraction of a slit between reference edge and component | 0.5 | < 1 | 2.5 | - | 0 to 60 |
| | | | | *Overall system accuracy (μm)* | | |
| Periodic light intensity | Light source, line scan camera system with self-scanned diode array‡ | 10 | ±20 | | | } 0 to 50 |
| | | 5 | ±20 | | | |
| | Light source, grating system; spring-loaded plunger | 10 | 2.5 | | | 0 to 50 |
| | Light source, grating system, spring-loaded plunger (a digital micrometer) | 25 | 5 | | | 0 to 50 |
| | Dual-frequency laser, interferometric system; end reflector attached to a component | 50 or more | ±0.02 | | | 0 to 50 |
| | Dual-frequency laser, interferometric system; light beam reflected from component's surface | up to 38 | ±0.02 | | | 0 to 50 |

\* Requires the use of a linearizer.

† The gap between the cell and laser is suited for measurement applications, limiting resolution to 2.5 μm.

‡ One particular application of a line scan camera system for non-contact, continuous height measurement with maximum sampling frequency of 2kHz.

| Size (mm) | Static force (N) | Relative cost | Whole system Relative cost | Resolution (μm) | Elect. freq. response (Hz) | Range-to-resolution ratio given by output section | Type of processing and output |
|---|---|---|---|---|---|---|---|
| Ø25 | 3 | 1 | 6 | 40 (wiredia.) | - | $2\times10^3$ | Power supply and $3\frac{1}{2}$ digit readout |
| Ø25 | 3 | 1.2 | 6 | 5 (contact noise) | - | $2\times10^3$ | Power supply and $3\frac{1}{2}$ digit readout |
| Ø20 | 2 | 1.3 | 6 | 10 | - | $2\times10^3$ | Power supply and $3\frac{1}{2}$ digit readout |
| Ø19 | 1 | 2 | 20 | 2 | 500 | $2\times10^4$ | Oscillator, demodulator, range switching $3\frac{1}{2}$ digit readout |
| Ø6 | 0.3 | 1.4 | 20 | 0.1 | 500 | $2\times10^4$ | |
| Ø25 | 0.001 | 10 | 100 | 0.005 | 500 | $2\times10^4$ | Oscillator, demodulator, range switching, gearbox, chart recorder |
| Ø25 | 0.5 | 2.4 | – | – | – | – | Requires stabilized power supply and a digital or analogue readout |
| Ø12 | None | 4.5 | 22.5 | 2 | 500 | $2.5\times10^3$ | Oscillator, demodulator, range switching $3\frac{1}{2}$ digit readout |
| Ø74 | None | 2 | 13 | 5 | 104 | $2\times10^3$ | Oscillator, demodulator $3\frac{1}{2}$ digit readout |
| Ø7 | None | 1.2 | 13 | 0.2 | 104 | $2\times10^3$ | |
| Ø30 | 1 | 2 | 50 | 0.1 | 10 | $5\times10^5$ | Oscillator, demodulator, automatic D/A tracking six-digit readout |
| Ø10 | 0.3 | 1.4 | 50 | 0.01 | 10 | $5\times10^5$ | |
| Ø40 | 2 | 1.6 | – | – | – | – | Requires stabilized power supply and a digital or analogue readoutt |
| Ø18 | None | 2 | 20 | 12.5 / 0.025 | $3\times10^3$ | $2\times10^3$ | Oscillator, demodulator, meter readout |
| Ø15 | None | 3 | 16 | 0.1 | 300 | $1.5\times10^4$ | Oscillator, demodulator, automatic D/A tracking six-digit readout analogue readout |
| Ø6 | 0.01 | – | 16 | 0.1 | 3 to $10^3$ | $1.5\times10^3$ | Change amplifier, gearbox, meter readout of average displacement |
| Ø20 | 0.001 | 8 | 60 | 0.015 | 500 | $3.5\times10^3$ | Amplifier, range switching, compensation for intensity variations, gearbox, chart recorder |
| – | None | – | 30 | 2.5 | $1.5\times10^3$ | $4\times10^3$ | Amplifier, range switching, compensation for intensity variations, two-axis meter readout |
| – | None | – | 80 | 0.025 | 105 | $2\times10^4$ | Amplifier, range switching, compensation for intensity variations, analogue output |
| Camera 140 × 300× 200 | None / None | – / – | 65 / 65 | 20 / 10 | $2\times10^3$ | 500 / 500 | Line scan camera provides scanning circuitry for diode array; a processing unit provides drive and logic circuiry, amplifier for chart recorder |

| | | | | | Maximum velocity (mm s$^{-1}$) | | |
|---|---|---|---|---|---|---|---|
| 100 × 30× 12 | 0.5 | 6 | 13 | 1 | 500 | 104 | Amplifier, counter, interpolator five-digit readout |
| 150 × 50 × 30 | 8 | – | 3.2 | 1 | 50 | $2.5\times10^4$ | Amplifier, counter, interpolator five-digit readout, integral with conversion and coupling |
| 80 × 50× 50 | None | – | 250 | 0.01 / 0.001 | 300 / 25 | $5\times10^6$ / $5\times10^7$ | Amplifier, counter, interpolator laser stabilizer, seven-digit readout |
| – | None | – | 250 | 0.01 / 0.001 | 300 / 25 | $4\times10^6$ / $4\times10^7$ | Amplifier, counter, interpolator laser stabilizer, seven-digit readout |

### 4.9.4 Merit and cost

It is very difficult to give any idea of the comparative cost of transducers because it ultimately depends on the application. Overall the cost depends on the four sections of transfer mechanism (mechanical coupling), energy conversion, processing and output, each one having its own characteristics. Basically the transfer mechanism determines the frequency response of the system, the processing being designed accordingly.

It seems that the cost of the transducer is largely dependent on the range-to-resolution ratio and, to a much larger extent, on the frequency characteristic assuming a particular accuracy. An increase in accuracy requires an increase in resolution.

Resistive systems are cheap but slow. Resolutions of 5 $\mu$m are obtainable with a range-to-resolution of $10^4$. The frequency response, however, is low when the system is coupled to a plunger. Values of less than a few hertz are common.

Capacitance methods generally need good electronics to detect the small capacitance charges which occur. They also need good shielding to reject stray signals. However, capacitance methods have been used very effectively at the very highest resolutions to measure the gap between the slides of very high-accuracy instruments. Resolutions of nanometres are possible in such circumstances. A linearity of 0.01 % is achievable using high-performance D/A tracking methods. Frequency responses tend to be low.

Inductive systems have a greater energy gradient than capacitative methods and suffer less from strong fields, which is why inductive methods are often used in metrology. Frequency responses of up to 500 Hz can be achieved with the beam type of mechanical coupler and 15 Hz with the plunger. The range-to-resolution ration is generally less than the capacitative method, being 5 x $10^4$ typically, mainly because of the non-linearities inherent in the complicated field geometry in the core and coils.

Generating methods, although fairly economical and easy to use, require coupling to very high input impedance since they cannot easily deliver currents. This limits the resolution to about 100 nm with a range-to-resolution of $10^3$. Cantilever couplings provide a fast response of the order of kilohertz.

The optical methods in use when utilizing fringes and which are in the non-contacting zone can have very fine resolutions in the range of nanometres or even less using heterodyne methods. They suffer from being relative in the sense that a count of fringes is easily lost in practice. The use of wideband interference methods is making this disadvantage less of a problem.

### 4.9.5 Examples of transducer properties

A few typical systems will be given to illustrate the sensitivity resolution and noise values often obtained [235].

#### 4.9.5.1 Inductive

An inductive bridge transducer is often used. It may have a circuit like that shown in figure 4.252.



**Figure 4.252** Inductive transducer—typical circuit.

(*a*) *Sensitivity*

Typical input current is $0.4 \times 10^{-6}$ A $\mu m^{-1}$ displacement as sensed at the point D, the virtual earth of the pre-amplifier. To achieve this level of sensitivity the gaps of the magnetic circuit between the stator and armature are reduced from the usual 2 mm by a factor of 10.

The sensitivity of the transducer depends on the energy gradient set up in the gaps. These energy gradients can be enhanced by the use of new materials.

(*b*) *Noise*

These are as input currents referred to point D:

1. Transducer circuit shot noise $\sim 3.4 \times 10^{-11}$ A. This value is large and important because of the magnitude of the current supplied to the inductive transducer coils.
2. Transducer circuit Johnson noise $\sim 2.1 \times 10^{-11}$ A. This originates in the feed resistor connection between B and D.
3. Amplifier low-frequency noise $\sim 0$ A. The carrier system enables the noise to be eliminated by conventional bandpass filtering.
4. Amplified broadband noise $-3.4 \times 10^{-11}$ A. This assumes an amplifier of moderate noise performance, that is 10 nVHz$^{-1/2}$.

(*c*) *Resolution*

The combined random noise input is of the order of $5.3 \times 10^{-11}$ A. This is conventionally *half* that of the displacement equivalent of the resolution, that is $2 \times 5.3/4 \times 10^4$ $\mu m$ (i.e. 2.7 Å).

To improve the resolution the excitation could be increased until it gets up to the maximum allowable heating level. An alternative is to pick a better amplifier with superior noise performance and to insert a step-down transformer in parallel with the points $\alpha$, $\beta$.

Note also that the finish of the transducer parts should be very high ($<0.1$ $\mu m$ $R_a$) and with about $10^4$ asperities over the length of the armature. These precautions help to eliminate any microsystematic errors present.

(*d*) *Application*

The calculations show that this type of transducer is adaptable to the measurement of microtopography at high resolution although with limited range.

By scaling the magnetic reluctance in the design (i.e. increasing the gap) the resolvable displacement and measuring range can be increased to 5nm and 200 $\mu m$ respectively. Non-linearity errors grow in proportion to the range so that it is not possible to use such a transducer to measure both form and microtopography.

### 4.9.5.2 *Capacitative*

A typical capacitative transducer is shown in figure 4.253.

(*a*) *Sensitivity*

For a small transducer and according to figure 4.250 a sensitivity at the preamplifier output is deduced as being about 0.5 Vmm$^{-1}$. A transducer of active length 15 mm and diameter of 9.5 mm is assumed with an interelectrode gap of about 1 mm and an output impedance of 0.5 pF. In this set-up there is no capability for increasing the energy gradient, but a very adequate range with linear operation is possible.

(*b*) *Noise*

This is calculated as output voltages of the preamplifier:

1. Transducer circuit shot noise $\sim 2.5 \times 10^{-8}$. This assumes an amplifier bias current of 5pA.
2. Transducer circuit Johnson noise $\sim 6.5 \times 10^{-7}$ V. The noise originates in the 50 MΩ resistor of feedback resistance.

**Figure 4.253** Typical capacitance transducer

3. Amplifier low-frequency noise ~0 A, as before.
4. Amplitude broadband noise $\sim 1.4 \times 10^{-7}$ V. The amplification associated with this noise depends on the feedback ratio of the preamplifier. For the values of the figure the amplification is $2 \times$. The amplifier noise is assumed as ll mVHz$^{-1/2}$.

(*c*) *Resolution*

The combined noise is estimated as about $7.1 \times 10^{-7}$ V. The signal for 1 $\mu$m displacement is $0.5 \times 10^{-3}$ V. The resolution is given conventionally as $2 \times$ the displacement equivalent of the noise $14.2 \times 10^{-7} / 0.5 \times 10^{-3} =$ 3 nm.

To improve the resolution the 50M$\Omega$ feedback resistor could be increased to 500 M$\Omega$. In principle the resistor can be eliminated altogether by using a voltage pick-off acting as a source follower, A resolution of 1 nm is then possible. The finish of the electrodes has to be high ($<0.1$ $\mu$m $R_a$) and without flaws. The measuring range is limited by the fringe field at the outer ends of the stator (~2 mm).

The energy gradient is not concentrated as in other types. Transducers of this type are feasible for the general purpose moderate-resolution moderate-range instrument. The practical difficulty of fitting the pick-off electrode into a side-acting pick has resulted in its absence from instruments measuring microtopography but not necessarily form with a range of 1 mm.

### 4.9.5.3 Optical lateral positional sensor

This type is shown in figure 4.254.

(*a*) *Sensitivity*

Assuming equal optical elements and a consequential optical magnification of $3 \times$, a sensor responsivity of 0.4 (*A*/$\omega$) and a received power of 100 $\mu$m, the sensitivity of differential input current between those at A and C is of value 64 $\mu$A mm$^{-1}$ (also the dimension of the sensor is assumed to be about 2.5 mm).

The specification includes a limiting value for photocurrent and thus for measured power. The calculation contains an assumption of a sensor microcurrent of 40 $\mu$A. (This can be increased by two or three times if required.)

This type of sensor is akin to an electro-optic potentiometer so that the range of linearity does not depend on the spot size and is typically one-third of the sensor dimension.

**Figure 4.254** Optical lateral positional sensor.

(*b*) *Noise*

An input to preamplifier A and C the noise is evaluated as:

1. Shot noise ~ $1.8 \times 10^{-11}$ A.
2. Transducer Johnson noise ~$0.9 \times 10^{-11}$ A.
3. Amplifier low-frequency noise ~0 A.
4. Amplifier broadband noise ~$0.5 \times 10^{-11}$ A.

The noise is relatively large because of the limitation of large output resistance. The output resistance is typically 5kΩ. The resolution is approximately $7 \times 10^{-11}$ /$6 \times 10^{-9}$ $\mu$m = 1.1 nm.

It is doubtful if the resolution can be made much better. There are limitations to increasing the received power and photocurrent. (The resistance values in the sensor are substantial.) The best method of increasing the resolution lies in optical magnification. The spot size is increased simultaneously but this can be advantageous for the reduction of spatial microsystematic error.

### 4.9.5.4   *Optical area pattern photodiodes, transducer for lateral displacement*

This type is shown in figure 4.255.

(*a*) *Sensitivity*

Typically a uniform beam of 3 mm dimension in cross-section is assumed with a power of 1 mW. A polarizing beam arrangement results in a high level of transmission to the photodetector. In accordance with figure 4.250 the return beam is displaced in the same direction as the corner cube and magnified by 2×. Using these assumptions the sensitivity is evaluated at 0.66 mA mm$^{-1}$.

(*b*) *Noise*

1. Shot noise ~$5.6 \times 10^{-11}$ A. This noise is usually rated as the important one for this type of photodetector.
2. Johnson noise <$9 \times 10^{-12}$A. A value of $R_f$ >1 MΩ is assumed.
3. Amplifier low-frequency noise ~$10^{-13}$ A. This assumes an amplifier of moderate noise performance of 10 mVHz$^{-1/2}$. The relative insignificance is the reason why the shot noise is the significant source of noise.

(*c*) *Resolution*

The total noise is about $1.6 \times 10^{-10}$ A. The resolution is $2 \times 1.6 \times 10^{-10}$ /$0.66 \times 10^{-3}$ mm=5 Å.

**Figure 4.255** Optical area pattern photodiode transducer—lateral displacement.

It is possible to improve the resolution by means of optical magnification as the beam can be confined largely by tubes and sealed by glass. The collimation should be of high standard. The introduction of magnifying optics implies also that the photodiodes must be larger (figures 4.256 and 4.257).

### 4.9.6 Talystep

As an example take a typical instrument for measuring fine surfaces (Talystep). Consider some of its noise characteristics, especially the electrical limitations. This is only meant to be concerned with typical values. Better performance in individual cases can always be expected (references [6–11]).

*(a) General*

The electrical noise is generated in the transducer and the input stage of the electrical amplifier. The transducer is typically of the balanced inductive bridge type. At the centre top of the transducer the sensitivity is about 1 mV RMS per micrometre and thus the current sensitivity for supplying current to the input transpedance is about $0.5$ RMS $\times 10^{-6}$ A$\mu$m$^{-1}$.

Inherent noise due to sine wave excitation can be ignored. Also, any noise sidebands that might cross over from the sine wave driver only have a modulating effect to the extent that the transducer bridge is out of balance.



**Figure 4.256** Magnified receiver optics.

**Figure 4.257** Pin diode.

(*b*) *Transducer circuit shot noise*

The current fluctuation given by the shot effect (Brownian motion) is

$$I_n^2 = 2I_{av}eB \quad (\text{A}^2)$$

(4.306)

where $B$ is the bandwidth in hertz, assumed here to be 25, and $e$ is electronic charge ($1.6 \times 10^{-19}$ C). $I_{av}$ is $0.5 \times 10^{-3}$ A (the carrier current).

Only a limited fraction (say, 10%) of the noise in equation (4.306) gets to the amplifier because of the high mismatching of the pick-off. Also, the bandwidth has to be multiplied by $\sqrt{2}$ because the transducer system is one of amplitude modulation of a carrier, so equation (4.306) becomes

$$I_n^2 = \frac{2\sqrt{2}}{10} I_{av}eB \quad (\text{A}^2).$$

(4.307)

Hence ($I_n$) RMS = $2.4 \times 10^{-11}$ A.

(*c*) *Johnson noise of transducer circuit*

Because of the relatively low value of the transducer bridge output impedance, the main Johnson noise contribution comes from the coupling resistor and results in

$$I_n^2 = \frac{4kT}{R} B \quad (\text{A}^2)$$

(4.308)

where $k$ is Boltzmann's constant (JK$^{-1}$) ($1.4 \times 10^{-23}$), $T$ is temperature (K) (300 K), $B$ is the bandwidth allowing for the carrier system and $R$ is 2.5 k$\Omega$, from which ($I_n$)$_{RMS}$ is $1.5 \times 10^{-11}$ A.

(*d*) *Amplifier noise*

This originates mostly at the input and is expressed as low-frequency reciprocal 'popcorn' noise, which is mostly removed by a suitable choice of carrier, and broadband noise, which is virtually related to shot noise. It is often quoted as a voltage of magnitude $V_n$ volts RMS per $\sqrt{\text{Hz}}$ bandwidth across the input port.

The equivalent amplifier noise output for a typical system is

$$V_n \sqrt{B} R_f/R$$

(4.309)

where $R$ is the coupling resistance (2.5 k$\Omega$), $R_f$ is the amplifier feedback resistor and $V_n$ is typically 10 nVHz$^{-1/2}$. The amplifier noise current as an RMS value is therefore about $2.4 \times 10^{-11}$ A.

### (e) Total electrical noise

The total noise input current is the square root of the sum of the noise powers and yields about $3.7 \times 10^{-11}$ A.

However, as the input current for a 1 nm displacement of the sensor is $0.5 \times 10^{-6}$ A, the resolution limit in equivalent displacement terms is about 0.2 nm, which is entirely consistent with measured values as shown by Franks who, in fact, claims even lower results.

### (f) Possible improvements

An obvious improvement is to increase the transducer excitation. This doubles the sensitivity (yet increases shot noise by $\sqrt{2}$; see equation (4.307)) but may well increase the heating effect within the transducer. A trade-off between resolution and bandwidth is possible but also doubtful as the resolution increases only as $Hz^{-1/2}$. It seems therefore that only marginal improvements are likely on the limit quoted above. Thus a calibration system that can resolve 0.02 nm should be the aim. The use of capacitative transducers is also a possibility but they do suffer more from stray fields and there is less energy change per unit displacement than in inductive devices.

### (g) Mechanical possibilities

Great care must be taken to eliminate mechanical vibration, in effect by floating the whole device relative to the external environment. The whole apparatus should also be within two nested levels of a temperature-controlled room. Low-frequency thermal drift need not always be disastrous, provided that its period is long relative to the measurement time for other features of interest.

Most of the important high-frequency noise is a result of motor vibration or sometimes simply acoustic noise impinging on the sensitive transducer-drive-column mechanical loop. Typical loop sizes are about 100mm. The effective mechanical loop can be considerably reduced for step height measurement by the use of differential sensors, since halving the effective loop size reduces noise by a quarter using even the crudest Rayleigh principle. However, surface finish measurement using this technique would require some computation. Temporal integration of noise when measuring steps can be achieved providing the drift is low and AC coupling used.

Another mechanical point which has to be taken into account if surface texture is to be measured is the physical size and shape of the tactile sensor; it is no good trying to measure subnanometre surfaces if the stylus is unknown. Calibration of the stylus is therefore one of the most important aspects of the whole exercise. At the present time spatial resolutions of 50 mm have been achieved.

Noise in lasers [240], is an important subject in its own right and is concerned with the consideration of gain (power) central wavelength or bandwidth.

From the point of view of surface metrology these are not as demanding as might be expected. Usually it is the actual wavelength that is important. Spectral bandwidth varies the temporal coherence and hence fringe contrast reduction, but conventional laser systems which operate at very low powers of the order of milliwatts are adequate for normal surface metrology. The problem of noise is very much more serious for high-power lasers.

The effects of bandwidth under these circumstances can be the generation either of white noise or of coherent beat frequencies. In lasers such as the helium-neon laser the bandwidth is determined primarily by the Doppler profile of the spectral emission line, which is usually Gaussian and highly predictable in width. The properties of ion lasers are somewhat different since the velocities are determined mainly by electric fields in the plasma rather than purely thermal effects. As a rough rule of thumb, laser noise is proportional to emission linewidth, whilst detector noise such as shot noise or Johnson noise is proportional to system bandwidth, as seen in the earlier considerations of transducers.

It is well known that large-area photodiodes are liable to have inhomogeneities and pinhole-like flaws over the photosensitive layers as a result of manufacture. These flaws cause microsystematic errors. For this reason the region on the left in figure 4.257 could include alternative optics. Pin diodes are very uniform and

used at the focus of a segmented lens over which the light beams are subjected to a stabilized light pattern. Transducers like this are suitable for general purpose topographic use, having moderate resolution and range. These few examples show how the electronic noise is determined for typical applications. Basically there are two sources: the transducer itself as seen at the output of the preamplifier, and the amplifier unit which follows it. In many cases it is the shot noise that dominates.

The exercise illustrates that in most cases it is possible to improve on the resolution or range but usually not both, underlining the usefulness of the range-to-resolution ratio as a merit indicator for transducers.

#### (h) Mechanical noise and associated effects

This subject has been discussed in some great detail by Jones [238, 239].

The Brownian movement of internal particles will cause fluctuations in the dimensions of any mechanical object. This effect is very small indeed. Jones has considered this and has estimated that the RMS length change in a 10mm brass block of 1 mm$^2$ area is likely to be $10^{-6}$ m when averaged over a 1 second interval. Other inherent mechanical instabilities are material creep, in which the defects of the structural lattice slowly rearrange under stress to produce dimensional changes. Also, in dynamic applications mechanical hysteresis can give differing strain values for a given stress.

Joints are also important and can be a source of instability. On clamped joints creep and hysteresis have been observed as significant. Bolted joints appear to be satisfactory for nanometre resolutions but hardly for subnanometre work.

The most significant mechanical effect is thermal expansion [237, 236]. Absolute size changes with temperature in the range of 0.1 to 20 ports in $10^6$ k depending on the material. Differential temperature effects may also lead to errors as they produce geometric shape changes in the framework of an instrument. This can be very important in instruments not properly sited. Temperature changes need not be gradual. The effects of the sun sporadically falling onto an instrument can be devastating to its operation. Also, discontinuous air-conditioning can cause a problem. Even the presence of the human body in the vicinity of an instrument can disturb it, especially at high accuracy. Compensation with electric light bulbs has been used to good effect in the past.

A good design aims to reduce thermal effects by the combination of different thermal coefficient materials arranged in re-entrant configurations or used to provide counter stress. Baird [241] used steel wire wound around an Invar rod. A rise in temperature decreases the circumferential stress provided by the steel on the Invar which, in turning, invokes Poisson's ratio and decreases the length of the Invar (obviously only within a relatively small range).



**Figure 4.258** Mechanical loop with extraneous noise

**Figure 4.259** Reduction of noise by double-probe system.

Another factor is barometric pressure. The magnitude of the effect obviously depends on the elastic modulus and Poisson's ratio for the material. It is of the order of 1 part in $10^8$ for a typical diurnal pressure variation. Constant-pressure environments can if necessary avoid this problem, usually by providing a vacuum.

The size of the mechanical loop between the transducer and the specimen can be critical. Using Rayleigh's criterion for the cantilever is a crude rule [242], but the vibrational effects from sources such as motors, gearboxes, drives, etc, can be critical. Presented simply, the situation is as shown in figures 4.258 and 4.259. (See the section on design shown earlier.)

The shorter the distance around the mechanical loop the more stable the situation. The cantilever approximation is an upper bound in the sense that it overemphasizes the possible departure from the geometry shown in figure 4.259 as a simple diametral representation.

By carefully picking the mechanical configuration the mechanical loop can be closed on the workpiece itself rather than in the instrument. Such an example has been given before in the multiprobe methods of looking at roundness, straightness, etc. Another example is the accurate measurement of step height and surface texture using two probes.

### 4.9.7 Comparison of techniques—general summary

Take as an example a typical comparison of the measurement of very fine surfaces [244] (figure 4.260).

From this comparison it is clear that, in most cases, the stylus method still stands out as the general method. Others such as TIS and diffraction have certain advantages in some instances but not often enough to depose the stylus method. Perhaps in the future they will. Optical followers are challenging fast.

What has emerged is that many comparisons are doomed to failure unless like is compared with like. The spatial bandwidth of the instruments must be about the same. Never mind what happens outside the comparison bandwidth—there are enough problems within it. Church *et al* [243] demonstrate this when comparing stylus methods with optical phase methods. In most cases, not only do the bandwidths have to be compatible for the comparison but also the transfer characteristic of the equipment within this band.

Comparisons have been made with diffraction optics and electron optics [244] and the point is made, quite rightly, that no single technique or instrument can measure the entire range of amplitude/spatial wavelengths. The other point is that the ultimate sensitivity will depend on electronic and mechanical noise. The methods are basically complementary, as seen in table 4.12.

**Figure 4.260** Wavelength vs amplitude graph (after Stedman): P, profile; NP, profile; SIP, slope integration; N/T, Talystep; EM, electron microscope.

**Table 4.12**

| Method | Spatial resolution | $z$ resolution | Range $z$ | Frequency | Comments |
|---|---|---|---|---|---|
| Stylus | $0.1\ \mu$m to 1 mm | 0.3 nm | $50\ \mu$m | 20Hz | Contacts workpiece; easy use; traceable |
| Optical probe | $0.5\ \mu$m | $0.1\ \mu$m | $20\ \mu$m | 10 Hz-30−kHz | Non-contacting less traceable; with servo drive; range extended to $50\ \mu$m |
| Heterodyne | $2.5−200\ \mu$m | $0.2\ \mu$m | $0.5\ \mu$m | 10 Hz | Requires computer unravelling |
| TIS | $0.6−12\ \mu$m | 1 nm | $0.1\ \mu$m | seconds | |
| Scatterometer diffraction | $\sim10\ \mu$m | 1 nm | $\lambda/8$ 100 nm | seconds | Resolution depends on aperture; insensitive to movement |
| TEM | 2 nm to $1\ \mu$m | $2\ \mu$m | 100 nm | minutes | Replication needed can destroy surface |
| SEM | 10 nm | 2 nm | $2\ \mu$m | minutes | Vacuum needed |
| STM | 2.5 nm | 0.2 nm | 100 nm | minutes | Vibration-free mounting |
| Normarsky | $>0.5\ \mu$m | | | minutes | Needs interpretation and certain minimum reflectivity |
| Capacitance | $2\mu$m | 1 nm | 50 nm | 2 kHz | Needs conductors |
| Interferometry | $2\mu$m | 1 nm | $1\ \mu$m | minutes | |

Note that ultrasonic pneumatics are suitable for rougher surfaces only.

The most important trend concerning instrumentation is discussed in chapter 9. The mechanism in a philosophical sense will change from the purely deterministic, as in engineering metrology, down to statistical and even probabilistic in nanotechnology. Furthermore, rather than developing instruments that measure specific properties such as surface texture, more integrated measurement will be preferred because at the small dimensions there are no differences. In the atomic and nanometre domains, as well as changes in the fundamental theory behind the measuring technique mentioned above there will have to be a reappraisal of what is actually being measured and how this relates to conventional features.

Possibly the most important trend is concerned with the disciplines. All the earlier surface metrology instruments were designed by and for engineers. This is changing rapidly. Because the applications drive the specifications chemists, physicists and biologists are now far more involved. Unfortunately this brings communication problems. Also trying to marry the top-down approach of the engineer with the bottom-up approach of the physicist is proving to be very difficult. It seems that a new type of person is needed: one equally versed in physics, engineering and chemistry would be a good start!

The introduction of SPM instruments has rocked the comfortable instrument development routines that have existed over the past decade or so. Questions of traceability of quantities other than geometry now have to be addressed. Some of these issues will be examined in chapter 5.

## References

[1]   Schmalz G 1929 *Z. VDI* **73** 144–61
[2]   Berndt G 1924 Die Oberflachenbeschaffen heiger bei verschiedenen Bearbeitungs methoden Loewe *Notizen* **9** 26
[3]   Andrews P 1928 The contouring of smooth surfaces *J. Sci. Instrum.* V 209
[4]   Harrison R E W 1931 A survey of surface quality standards and tolerance costs based on 1929–30 precision grinding practice, *Annual Meeting ASME 1930, Trans. ASME* **53** 25
[5]   Abbott J and Firestone A 1933 A new profilograph measures roughness of finely finished and ground surfaces *Autom. Ind.* p 204
[6]   Clayton D 1935 An apparatus for the measurement of roughness *Mech. Eng.* **27** 321
[7]   Zeiss C 1934 Methods and instruments to illuminate surfaces *Jena Mach.* **78** 502
[8]   Linnik W 1930 Ein Apparat zur Messung von Verschiebungen in der Sehrichtung *Z. Instrum.* **50** 192
[9]   Perthen J 1936 Ein neues Verfahren zum Messen der Oberflachengute durch die Kapazitat eines Kondensators *Maschinenbau Betr.* **15** 669
[10]  Nicolau P 1939 Quelque recent progrès de la microgeometrie des surfaces usinées et de l'integration pneumatique des rugosités superficielles *Mecanique* **23** 152
[11]  Von Weingraber H 1942 Pneumatische Mess und Prufgerate *Maschinenbau* **21** 505
[12]  Schlesinger G 1940 Surface finish *Machinery* **55** 721
[13]  Schlesinger G 1942 Research on surface finish *J. Inst. Prod. Eng.* **29** 343
[14]  Schlesinger G 1942 Surface finish *Rep. Res. Inst. Prod. Eng.* (London)
[15]  Guild A 1940 An optical smoothness meter for evaluating the finish of metals *J. Sci. Instrum.* **17** 178
[16]  Tuplin V 1942 Autocollimator test for flatness *Machinery* **61** 729
[17]  Tolansky S 1947 Application of new precision interference methods to the study of topography of crystal and metal surfaces *Vertrag Tag. Dsche Phys.* Ges. 5–7 Gottingen
[17b] Whitehouse D J 1990 Biographical memoirs of R.E. Reason FRS *Mem. Proc. R. Soc.* **36** 437–62
[18]  Jost P 1944 A critical survey of surface finish parameters *Machinery* **66** 690
[19]  Timms C and Schole S 1951 Surface finish measurement *Metal Treatment* **18** 450
[20]  Timms C 1945 Measurement of surface waviness *Proc. IMechE* **153** 337
[21]  Ransome M 1931 Standards for surface finish *Am. Mach.* **74** 581
[22]  Deale J 1931 Standardisation of machined finishes *Mech. Eng.* **53** 723
[23]  Bodart E 1937 Les etats de surfaces standards HI, p 1–18
[24]  Broadston J A 1944 Standards for surface quality and machine finish description *Prod Eng* 627
[25]  Smith S T and Chetwynd D G 1992 *Foundations of Ultra Precision Mechanical Design* (London: Gordon and Breach)
[26]  Moore 1970 *Foundations of mechanical accuracy* Moore Special Tools
[27]  Pollard A F C 1929 *The Kinematic Design of Couplings in Instrument Design* (London: Hilger and Watts)
[28]  Jones R V 1962 Some uses of elasticity in instrument design *J. Sci. Instrum.* **39** 193–203
[29]  Zhang G X 1989 A study on the Abbé principle and Abbé error *Ann. CIRP* **38** 525

[30] Bryan J B 1979 The Abbé principle revised—an updated interpretation *Precis. Eng.* no 3

[31] Whitehouse D J 1976 Some error separation techniques in surface metrology *Proc. Inst. Phys. J. Sci. Instrum.* 9531–6

[32] Chetwynd D G 1987 Selection of structural materials for precision devices *Precis. Eng.* **9** 3–6

[33] Chetwynd D G 1989 Material selection for fine mechanics *Precis. Eng.* **11** 203–9

[34] Ashby M F 1989 On the engineering properties of materials *Acta Metall.* **37** 1273–93

[35] Ashby M F 1991 On materials and shapes *Acta Metall.* **39** 1025–39

[36] Whitehouse D J 1996 Optical methods in Surface Metrology Milestone 129 *SPIE*

[37] Hicks T R and Atherton P D 2000 The micropositioning book (London: Penton Press)

[38] Agullo J B and Pages-Fita 1970 Performance analysis of the stylus technique of surface roughness assessment—a random field approach *Proc. 15th MTDR Conf.* p349

[39] Whitehouse D J 1979 A theoretical investigation of stylus integration *Ann. CIRP* **23** 181

[40] Radhakrishnan V 1970 Effect of stylus radius on the roughness values measured with a stylus instrument *Wear* **16** 325–35

[41] Walton J 1961 Gramophone record deformation *Wireless World* **7** 353

[42] Church E 1978 *Proc. Opt. Soc. Am.* (*San Francisco*)

[43] Scarr A J and Spence A R 1977 *Microtechnic* **XXII** 130

[44] Frampton R C 1974 A theoretical study of the dynamics of pick ups for the measurement of surface finish and roundness *Tech. Rep.* T55, Rank Organisation

[45] Whitehouse D J 1990 Dynamic aspects of scanning surface instruments and microscopes *Nanotechnology* **2** 93–102

[45b] Levine B M and Dainty J C 1983 Non Gaussian image plane speckle measurement from diffusers of known statistics *Opt. Commun.* **45** 252

[46] Whitehouse DT 2000 Stylus damage prevention index *Proc. Inst. Mech. Engrs.* **214** 975–80

[47] Whitehouse D J 1988 A revised philosophy of surface measuring systems *Proc. IMechE., J. Mech. Eng. Sci.* **202** 169

[48] Binnig G, Rohrer H, Gerber C and Weibel E 1982 Surface studies by scanning tunnelling microscope *Phys. Rev. Lett.* **49** 57

[49] Binnig G and Rohrer H 1985 The scanning tunnelling microscope *Surf. Sci.* **152/153** 17–26

[50] Martin Y, Williams C G and Wickramasingh H K 1987 Atomic force microscope—force mapping and profiling on a sub 100 A scale *J. Appl. Phys.* **61** 4723–9

[51] Wickramasinghe H K 1989 Scanned probe microscopes *Sci. Am.* October 74–81

[52] Vorburger T V *et al* 1997 Industrial uses of STM and AFM

[53] Sullivan NT 1995 Current and Future Applications of in-line AFM for semiconductor water processing 2nd Workshop on ASPM Gaithersburg MD NISTIR 5752 NIST.

[54] Windt D L, Waskiewilz W K and Griffithe J E 1994 Surface finish requirements for soft x-ray mirrors *Appl. Optics* **33** 2025

[55] Jungles J and Whitehouse D J 1970 *Inst. Phys. E. Sci. Intst.* **3** 170

[56] Wang W L and Whitehouse D J 1995 Application of neural networks to the reconstruction of SPM images by finite tips *Nanotechnology* **6** 45–51

[57] Musselman I H, Peterson P A and Russell P E 1990 Fabrication of tips with controlled geometry for scanning tunneling microscopy *Precision Engineering* **12** 3–6

[58] Stevens R M D, Frederick N A, Smith B L, Morse D E, Stucky G D and Hansma P K 2000 *Nanotechnology* **11** 1–5

[59] Nguyen C V *et al* 2001 carbon nanotube tip probes, stability and lateral resolution in scanning probe microscopy and application to surface science in semi conductors *Nanotechnology* **12** 363–367

[60] Tie L *et al* 1998 *Science* **280** 1253–6

[61] Haffner J H, Cheung C L and Leiber C M 1999 *Nature* **398** 701

[62] Wong E W *et al* 1997 *Science* **277** 1971–5

[63] Walter D A *et al* 1999 *Appl. Phys. Lett.* **74** 3803–5

[64] Kislov V, Kolesov I, Taranov J and Saskovets A 1997 Mechanical feautures of the SPM microprobe and nanoscale mass detector *Nanotechnology* **8** 126–131

[65] Oh Tsu M 1998 Near field/atom optics and technology (Springer: Tokyo)

[66] Pilevar S, Edinger F, Atia W, Smolyaninov I and Davis C *Appl. Phys. Lett.* **72** 3133–5

[67] Lambelet P, Sayah A, Pfeffer M, Phillipona C, and Marquis-Weible F 1998 Chemically etched fibre tips for near field optical microscopy, a process for smoother tips *Appl. Optics* **37** 7289–92

[68] Peterson C A *et al* 1998 *Nanotechnology* **9** 331–38

[69] Butt H J and Jaschke M 1995 Calculation of thermal noise in atomic force microscopy *Nanotechnology* **6** 1–7

[70] Marti O *et al* 1990 Topography and friction measurement on mica *Nanotechnology* **2** 141

[71] Burnham N A *et al* 1997 How does a tip tap? 67–75

[72] Gibson C T, Watson G S and Myhra S 1996 Determination of H spring constant of probes for force microscopy/spectroscopy *Nanotechnology* **7** 1996 259–2

[73] Xu Y and Smith S T 1995 Determination of squeeze film damping in capacetance based cantileler force probes, *Prec. Engrng* **17** 94–100

[74] Whitehouse D J 1990 Dynamic aspects of scanning surface instruments and microscopes *Nanotechnology* **1** 93–102

[75] Brigs G A D 1992 *Acoustic Microscopy* (Oxford: Oxford University Press)

[76] Rabe U, Janser K and Arnold W 1996 *Rev. Sci. Inst.* **67** 328

[77] Xu Y, Smith S T and Atherton P D 1995 A metrological scanning force microscope's *Precision Engineering* **19** 46–55

[78] Peters J, van Herk P and Sastrodsnoto M 1978 Analysis of the kinematics of stylus motion *CIRP, Aacis Coop. Doc.* no R9, September

[79] Reason R E, Hopkins and Garrott 1944 *Rep.* Rank Organisation

[80] Shunmugam M S and Radhakrishnan V 1975 2 and 3D analysis of surfaces according to the E system *Proc. IMechE*

[81] Hunter A G M and Smith E A 1980 Measurement of surface roughness *Wear* **59** 383–6

[82] Whitehouse D J 1982 Assessment errors of finishing processes caused by skid distortion *J. Phys. E: Sci. Instrum.* **15** 1337

[83] Von Weingraber H 1956 Zur Definition der oberflachen Rauheit *Werkstattstech. Maschonenbau*

[84] Rubert 1967/68 *Proc. IMechE* **182** 350

[85] Deutschke S I, Wu S M and Stralkowski C M 1973 A new irregular surface measuring system *Int. J. Mach. Tool. Des. Res.* **13** 29–42

[86] Fugelso M and Wu S M 1977 Digital oscillating stylus profile measuring device *Int J. Mach. Tool. Des. Res.* **17** 191

[87] George A F 1979 A comparative study of surface replicas *Proc. Conf. on Metrology and Properties of Surfaces, Leicester* (Lausanne: Elsevier)

[88] Narayanasamy K *et al* 1979 Analysis of surface roughness characteristics of different replica materials *Proc. Conf. on Metrology and Properties of Surfaces Leicester* (Lausanne: Elsevier)

[89] Williamson J B P 1967/68 The microtopography of solid surfaces *Proc. IMechE 182*

[90] Sayles R C and Thomas T R 1976 Mapping a small area of a surface *J. Phys. E: Sci. Instrum.* **9** 855

[91] Tsukada T and Sasajima K 1981 3D measuring technique for surface asperities *Wear* **71** 1–14

[92] Idrus N 1981 Integrated digital system for 3D surface scanning *Precis. Eng.* **37**

[93] Whitehouse D J and Phillips M J 1985 Sampling in a 2D plane *J. Phys. A: Math. Gen.* **18** 2465–77

[94] Li M, Phillips M J and Whitehouse D J 1989 Extension of two dimensional sampling theory *J. Phys. A: Math. Gen.* **22** 5053–63

[95] Baker L R and Singh J 1985 Comparison of visibility of standard scratches *Proc. SPIE* **525** 64–8

[96] Morrison E 1996 Thedevelopment of a prototype high speed stylus profilometer and its application to rapid 3D Surface measurement *Nanotechnology* 37–42

[97] Young R D, Vorburger T V and Teague E C 1980 In process and on line measurement of surface finish *Ann. CIRP* **29** 435

[98] Young R D 1973 Light techniques for the optical measurement of surface roughness *NBS IR* 73–219, (Gaithersburg, MD: National Bureau of Standards)

[99] Dupuy 1967/68 *Proc IMechE 180* pt 3k

[100] Granger E M 1983 Wavefront measurement from a knife edge test *Proc SPIE* **429** 174

[101] Whitehouse D J 1983 Optical coding of focus positioning—Hadamard method *VNIIMS Conf Surface Metrology,* (*Moscow*) paper E2

[102] Simon J 1970 *Appl. Optics* **9** 2337

[103] Tolansky S 1948 *Multiple Beam Interferometry of Surfaces and Films* (Oxford: Oxford University Press)

[104] Milano E and Rasello F 1981 An optical method for on line evaluation of machined surface quality option *ActaW* 111–23

[105] Bennett H E, Bennet J M and Stanford J L 1973 Surface corregulanties and scattering *Proc Infrared Symp,* (*Huntsville, Alabama, January 1973*)

[106] Hamilton D K and Wilson T 1982 3D surface measurement using confocal scanning microscopes *J. Appl. Phys.* B **27** 211

[107] Bennett S D, Lindow J T and Smith I R 1984 Integrated circuit metrology using confocal optical microscopy *R. Soc. Conf. on NDT, July*

[108] Sheppard C J R and Wilson T 1978 *Opt. Lett.* **3** 115

[109] Young J Z and Roberts T 1951 The flying spot microscope *Nature* **167** 231

[110] Anamalay R V, Kirk T B and Pahzera D 1995 Numerical descriptions for the analysis of wear surfaces using laser scanning confocal microscopy wear **181-183** 771–776

[111] Hodgkinson I J 1970 *J. Phys. E. Sci. Instrum.* **3** 300

[112] Bennett H E and Bennet J M 1967 *Physics of Thin Films* Vol 4 (New York: Academic) pp 1–96

[113] Church E L, Vorburger T V and Wyant J C 1985 *Proc. SPIE* 508–13

[114] King R J, Downes M J, Clapham P B, Raine K W and Talm S P 1972 *J. Phys. E: Sci. Instrum.* **5** 449

[115] Schmalz G 1936 *Technische Oberflochenkude* (Berlin: Springer) pp 5–8

[116] Mitsuik Sato 1978 H frequency characteristics of the cutting process *Am. CIRP* **27** 67

[117] Bottomley S C 1967 *Hilger J. XI* (1)

[118] Remplir device—Opto electronics Stenkullen, Sweden.

[119] Michelson A A 1990 *Philos. Mag.* **5** 30 1

[120] Gehreke G 1906 *Anwendung der Interferenzen* p 120

[121] Jenkins F A and White H E 1951 *Fundamentals of Optics* (McGraw-Hill) p328

[122] de Groot P 2001 Zygo Seminar 2D/3D Surface measurement

[123] Sommargren G E 1981 *Appl. Opt.* **20** 610

[124] Sommargren G E 1981 *Precis. Engl.* **131**
[125] Huang C C 1983 *Proc. SPIE* **429** 65
[126] Wade H 1967 Oblique incident micro interferometer *Bull. Jpn Soc. Precis. Eng.* **4** 234
[127] Chapman G D 1974 Beam path length multiplication *Appl. Opt.* **13** 679
[128] Tanhnura Y 1983 *Am. CIRP* **32** 449
[129] Wyant J C, Koliopoulos C L, Bushan B and George 0 E 1983 An optical profilometer for surface characterisation of magnetic media *38th Meeting ASLE,* print 83-AM 6A-1
[130] Wyant J C 1975 *Appl. Opt.* **11** 2622
[131] Breitmeier U 1990 Optical follower optimised
[132] Fanyk, Struik K G, Mulders P C and Veizel C H F 1997 Stitching interferometry for the measurement of aspheric surfaces *Ann. CIRP* **46** 459
[133] Weinheim V C H, Wagner E, Dandliker R and Spenner K 1992 *Sensors* vol 6 (Weinheim: Springer)
[134] Hofler H and Seib M 1992 *Sensors* vol 6 (Weinheim: Springer) p570
[135] Hofler H and Seib M 1992 *Sensors* vol 6 (Weinheim: Springer) p574
[136] Gabor D 1948 A new microscope principle *Nature* **161** 40–98
[137] Gabor D 1948 *The Electron Microscope* (Electronic Engineering Monograph)
[138] Leith E N and Upatnieks J 1963 *J. Opt. Soc. Am.* **53** 1377
[139] Abramson N H and Bjelkhagen H 1973 *Appl. Opt.* **12** (12)
[140] Brooks R E 1917 *Electronics* May
[141] Fryer P A 1970 Vibration analysis by holography *Rep. Prog. Phys.* **3** 489
[142] Henderson G 1971 Applications of holography in industry *Electro. Opt.* December
[143] Ennos A E and Archbold E 1968 *Laser Focus* October
[144] Kersch L A 1971 *Materials evaluation* June
[145] Ribbens W B 1972 *Appl. Opt.* **11** 4
[146] Ribbens W B 1974 Surface roughness measurement by two wavelength holographic interferometry *Appl. Opt.* **13** 10859
[147] Balabois J, Caron A and Vienot J 1969 *Appl. Opt. Tech.* August
[148] Vander Lugt A 1968 *Opt. Acta* **15** 1
[149] Abramson N H 1975 *Ann. CIRP* **24** 379–82
[150] Erf R K 1978 *Speckle Metrology* (New York: Academic) p 2
[151] Jones R and Wykes C 1983 *Holographic and Speckle Interferometry* (Cambridge: Cambridge University Press)
[152] Goodman J W 1975 *Laser Speckle and Related Phenomena* ed J C Dainty (Springer) ch 2
[153] Asakura T 1978 *Speckle Metrology* (New York: Academic) p 11
[154] Jackeman E and Pussey P N 1975 Non-Gaussian fluctuation in electromagnetic radiation scattering by a random phase screen. *J. Phys. A: Math. Gen.* **8** 392
[155] Sprague R W 1972 *Appl. Opt.* **11** 2811
[156] Pederson H W 1976 *Opt. Commun.* **16** 63
[157] Fuji H, Asakura T and Shindo Y 1976 Measurement of surface roughness properties by means of laser speckle *Opt. Commun.* **16** 68
[158] Beckmann P and Spizzichino 1963 *The Scattering of Electromagnetic Waves from Rough Surfaces* (Oxford: Pergamon)
[159] Mandelbrot B B 1977 *The Fractal Geometry of Nature* (New York: Freeman)
[160] Jakeman E and McWhirter J G 1981 *Appl. Phys. B* **26** 125
[161] Fuji H 1980 *Opt. Acta* **27** 409
[162] Jordan D and Jakeman E 1985 *Proc. SPIE Conf.*
[163] Butters J N and Leendertz J A 1974 *Proc. Electro Optics Conf. Brighton 1972* (Kiver Communications)
[164] Ogilvy J A 1991 *The Theory of Wave Scattering from Random Rough Surfaces* (Bristol: Hilger)
[165] Maystre D 1988 Scattering by random surfaces in electromagnetic theory *Proc. SPIE* **1029** 123
[166] Bass G G and Fuchs I M 1963 *Wave Scattering of Electromagnetic Waves from Rough Surfaces* (New York: Macmillan)
[167] Chandley P J 1976 Determination of the autocorrelation function of surface finish from coherent light scattering *Opt. Quantum Mech.* **8** 329
[168] Welford W T 1977 Optical estimation of the statistics of surface roughness from light scattering measurements *Opt. Quantum Electron.* **9** 269
[169] Bennett H E and Porteus J 0 1961 *J. Opt. Soc. Am.* **51** 123
[170] Davies H 1954 *Proc. IEE* **101** 209
[171] Torrence K E 1964 *MS Thesis* University of Minnesota
[172] Milana E and Rosella F 1981 *Opt. Acta* **28** 111
[173] Rakels J H 1986 *J. Phys. E: Sci. Instrum.* **19** 76
[174] Thwaite E G 1980 *Ann. CIRP* **29** 419

[175] de Groot P, Delega X C and Stephenson D 2000 *Opt. Eng.* **39** 8

[176] Bjuggren M, Krummenocher L and Mattson L 1997 Non contact surface roughness measurement of engineering surfaces by total integrated infrared scattering *Precision Engineering* **20** 33-34

[177] Sawatari T 1972 Surface flaw detection using oblique angle detection *Appl. Opt.* **11** 1337

[178] Whitehouse D J 1972 Some modem methods of evaluating surfaces *Proc. SME, Chicago*

[179] Pancewicz T and Mruk I 1996 Holographic controlling for determination of 3D description of surface roughness *Wear* **199** 127

[180] Konczakowski A L 1983 *Int. J. Mach. Tool Des. Res.* **23** 161

[181] Konczakowski A L and Przyblsy W 1983 *Wear* **92** 1–12

[182] Whitehouse D J, Vanhert P, De Bruin W and van Luttervelt C 1974 Assessment of surface topology techniques in turning *Am. CIRP* **23** 265–82

[183] Ramesh S and Rammamoorthy B 1996 Measurement of surface finish using optical diffration techniques *Wear* **195** 148–151

[184] Church E L, Jenkinson H A and Zavada J M 1977 Measurement of the finish of diamond turned metal surfaces by differential light scattering *Opt. Eng.* **16** 360

[185] Elson J E and Bennett J M 1979 Vector scattering theory *Opt. Eng.* **18** 116

[186] Wang Y and Wolfe W L 1983 Scattering from micro rough surfaces: comparison, theory and experiment *J. Opt. Soc. Am.* **73** 1596

[187] Bonzel H P and Gjostein N A 1968 *J. Appl Phys.* **39** 3480

[188] Jungles J 1979 Private communication

[189] Bosse J C, Hansoli G, Lopez J and Mathina T 1997 *Wear* **209** 328–337

[190] *US Patent* 4092068

[191] Takayama H, Sekiguchi H and Murate R 1976 In process detection of surface roughness in machining *Ann. CIRP* **25** 467

[192] Tanner L H and Fahoum M 1976 *Wear* **36** 299–316

[193] Stover J C, Serati S A and Gillespie C H 1984 Calculation of surface statistics from light scatter *Opt. Eng.* **23** 406

[194] Guenther K H and Wierer P G 1983 Surface roughness assessment of ultra smooth mirrors and substrates *SPIE 401*

[195] Masurenko M M, Skrlin A L and Toporetts A C 1979 *Optico* **11** 1931

[196] Bennett J M, Guenther K H and Wierer P G 1983 Surface finish measurements on low scatter laser mirrors and roughness standards *15th Symp. on Optical Materials for High Power Lasers* ed H E Bennett (Washington, DC: NBS)

[197] Bennett J M, Guenther K H and Wierer P G 1984 *Appl. Opt.* **23** 3820

[198] Elson J M, Rahn J P and Bennett J M 1983 *Appl. Opt.* **22** 3207

[199] Sorides G N and Brandin D M 1979 An automatic surface inspection system for flat rolled steel *Automatika* **15** 505

[200] Pryer T R, Reynolds R, Pastorious N and Sightetc D 1988 *4th Conf. on Met. and Props of Surfaces* (Washington, DC: NBS)

[201] Baker L R, Kirkham A J, Martin S, Labb D R and Rourke C P *US Patent* 3892492

[202] Baker L R 1986 Standards for surface flaws *Opt. Laser Technol.* February, pp 19–22

[203] Baker L R 1984 Microscope image comparator *Opt. Acta* **31** 611–14

[204] Cerni R H 1962 Analogue transducers for dimensional metrology *Proc. Int. Prod. Eng. Res. Conf.* (*Pittsburgh*) p 607

[205] Sherwood K F and Crookall J R 1967/68 Surface finish instrument by electrical capacitance method *Proc.* IMechE **182** pt 3k

[206] Thomas T R (ed) 1982 Roughness measurement alternatives to the stylus *Rough Surfaces* (Harbour: Longman)

[207] Franson R E, Brecker J N and Shum L Y 1976 A universal surface texture measuring system *Soc. Mgn. Eng.* paper IQ 76-597

[208] Matey J R and Bland J 1985 Scanning capacitance microscanning *J. Appl. Phys.* **57** 1437

[209] Bhusham B 1984 Prediction of surface parameters in magnetic media *Wear* **95** 19–27

[210] Moore D F 1965 Drainage criteria for runway surface roughness *J. R. Aeronaut. Soc.* **69** 337

[211] Bickerman J J 1970 *Physical Surfaces* (New York: Academic)

[212] Graneck M and Wunsch H L 1952 Application of pneumatic gauging to the measurement of surface finish *Machinery* **81** 707

[213] Radhakrishnan V and Sagar V 1970 Surface roughness by means of pneumatic measurement *Proc. 4th Indian MTDR Conf.* (Madras: Indian Institute of Technology)

[214] Tanner L H 1978 A pneumatic Wheatstone bridge for surface roughness measurement *J. Phys. E: Sci. Instrum.* **12** 957

[215] Wager J G 1967 *Int. J. Prod. Eng. Res.* **7** 1

[216] Berry M V 1973 The statistical properties of echoes diffracted from rough surfaces *Philos. Trans. R. Soc. A* **273** 611

[217] de Bily M, Cohen Tanoudji F, Jungman A and Quentin J G 1976 *IEEE Trans. Sonics Ultrasound SU-23* 356

[218] Cohen Tanoudji F and Quentin G 1982 *J. Appl. Phys.* **53** 4060

[219] Clay C S and Medwin H 1970 *J. Acoust. Soc. Am.* **47** 1412

[220] Chubachi N, Kanai H, Sannomiya T and Wakahara T 1991 Acoustic Microscopy *Acoustic Imaging* **19** d Ermert M and Harjies H P, eds (London: Plenum Press)

[221] Chang T H P and Nixon W C 1966 *J. R. Microsc. Soc.* **88** 143

[222] Grundy P J and Jones G A 1976 *Electron Microscopy in the Study of Material* (London: Edward Arnold)

[223] Bowen D K and Hall R C 1971 *Microscopy of Material* (London: Macmillan)

[224] 1992 *Microtexture Determination and its Applications* (Inst. of Materials: London)

[225] Oxford Instruments—Guide book 1996 Electron Backscatter Diffraction (High Wycombe, UK)

[226] Hearle J W S, Sparrow J T and Cross P M 1972 *Scanning Electron Microscope* (Oxford: Pergamon)

[227] Howell P G T and Bayda A 1972 *Proc. 5th SEM Symp.* (*Chicago, April, 1972*)

[228] Chistenhauss R and Pfeffekorn G 1968 *Bietr. electronen microskop Direcktabb Oberflachen, Munich* **1** 129

[229] Jungles J and Whitehouse D J 1970 An investigation into the measurement of diamond styli *J. Phys. E: Sci. Instrum.* **3** 437–40

[230] Rasigni M, Rasigni G, Palmari J P and Llebaria A 1981 I surface profiles, II autocorrelation functions *J. Opt. Soc. Am.* **71** 1124 and 1230

[231] Guerra J M 1990 Photon tunneling microscopy *Appl. Optics* **29** 3741–3752

[232] McCutchen C W 1964 Optical systems for observing surface topography by frustrated total internal reflection and by interference *Rev. Sci. Instrum.* **35** 1340–1345

[233] Takahashi S *et al* 1997 High resolution photon scanning tunneling microscope *Nanotechnology* **8** 445

[234] Jungles J 2001 Private Communication

[235] Bellwood P R 1983, 1987 Private communications

[236] Sydenham P H 1972 *J. Phys. E: Sci. Instrum.* **5** 83, 733

[237] Neubert N K P 1963 *Instrument Transducers: An Introduction to their Performance and Design* (Oxford: Oxford University Press)

[238] Jones R V 1957 Instruments and the advancement of learning *Trans. Soc. Inst. Tech.* 3

[239] Jones R V 1961 *Proc. R. Soc. A* **260** 47

[240] Bloom A L 1965 *Spectra Phys. Tech. Bull.* **4**

[241] Baird K M 1968 *Metrologia* **4** 145

[242] Bishop R E D and Johnson D C 1960 *The Mechanics of Vibration* (Cambridge: Cambridge University Press)

[243] Church E L, Varburger T V and Wyant J C 1985 *Proc. SPIE* paper 508–13

[244] Stedman M Metrological evaluation of grazing incidence mirrors

[245] Lech M, Mruk I and Stupnicki I J 1984 *Wear* **93** 167

# Chapter 5
# Traceability—standardization—variability

## 5.1  Introduction

Traceability and standardization enable results to be compared nationally and internationally using common units. Failure to ensure that the results of a measurement or experiment conform to a standard is not in itself serious. It may be, for example, that a particular firm is investigating the function of a part it makes, so it does some test experiments. The measurement system may be needed to determine changes which occur during the experiment. It may not be necessary to communicate such results to the outside world and often in a proprietary situation not even desirable. Under these conditions all that is required is that the system within the firm is consistent, not necessarily traceable to national or international standards. However, the day of the family company working in isolation is over. Consequently, results have to be compared and instruments have to be verified according to mutually acceptable rules. Easily the safest way to do this is to relate all measurements to the international system of units and procedures. Then there is no doubt and no ambiguity. The problem is that it can introduce extra costs because of the need for training and the need to buy in test equipment. However, there can be no question that in the long term this is the most cost-effective way to proceed.

The whole subject of metrology is concerned with identifying and controlling errors and with putting confidence limits to measured values. In surface metrology there are two basic aspects to be considered: one concerns the errors in the parameters due to the measuring instrument, and the other is concerned with the intrinsic variability of parameters due to the surface itself. Both will be discussed here as well as aspects of calibration and traceability. Because the errors are often random in nature some treatment of the basic statistical tests will be given. This should not be taken as the basis for a student course because it is not taken to any depth. It is included merely to bring to the forefront some of the statistical tools that are available to the surface metrologist. To start the chapter the nature of errors will be considered.

In what follows, the way in which an instrument is calibrated will be considered, starting at the input and continuing through the system to the display. Obviously not every instrument can be considered, so as a starting point the stylus method will be used. The messages conveyed will be the same for all instruments. These will be considered when appropriate. Scanning probe microscopes will be included in this chapter.

Part of the question of traceability is that of variability, not the fidelity of the mean measurement of a parameter but the variation that arises when measuring it (whatever it may be). Many investigators think that a deviation or systematic error in an instrument can be acceptable but that variability cannot. This is only true if the systematic error is known or at least known about. The danger is that sometimes it is not.

## 5.2  Nature of errors

This is a subject which is very extensively reported. Basically errors can be split up into two: systematic and random errors. In practice the distinction is usually a question of timescale. The level of random errors may

be predicted and, at the same time, systematic or repeatable errors change their character with time. The essential difference is the time within which accurate or traceable readings need to be taken in order to guarantee consistent results.

### 5.2.1 *Systematic errors*

Systematic errors are sometimes called bias and can be caused by imperfect devices or measurement procedures and also by the use of false standards. Not excluded in the definition are the effect of (unwittingly) biased people and the environment. As mentioned above, systematic errors can be measured (if suspected) and can be corrected.

### 5.2.2 *Random errors*

Random errors produce the scatter of readings and are caused by non-constant sources of error either in time or in space. They cannot be eliminated because they have no known or deterministic value. They can often be minimized by the use of averaging techniques.

Other terms sometimes encountered include [1,2] gross error, which is due to a mistake, additive errors, which are usually zero point displacements, and multiplicative errors, which are characterized by their property of being multiplicatively superimposed on the measured value. These do depend on the numerical value of the measured quantity. They are based on the deviations of the measurement system from its desired value.

Absolute measurement errors are defined as the difference between the actual measured value and the true value (which is not known).

Relative errors are absolute errors divided by a reference quantity.

## 5.3 Deterministic or systematic error model

The absolute measurement error $\Delta x$ is the difference between the measured value $x_1$ and the true value $x_s$.
If the errors are time dependent

$$\Delta x(t) = x_1(t) - x_s(t). \tag{5.1}$$

In system theory the transfer property of a system is designated by its transfer function $S(p)$. It is the quotient of the Laplace transform of the output quantity to that of the input. Thus the error-free transfer function of the measurement system is

$$S_1(p) = \frac{x_o(p)}{x_s(p)}. \tag{5.2}$$

The actual output is obtained from

$$S_F(p) = \frac{x_1(p)}{x_s(p)}. \tag{5.3}$$

This leads to the concept of the measurement error transfer function, $S_F(p)$:

$$S_F(p) = \frac{S_1(p)}{S_s(p)} - 1 = \frac{\Delta x(p)}{x_s(p)}. \tag{5.4}$$

This concept is useful in thinking about the propagation of errors through a system.

*(a) Sensitivity*

This is the ability of the measuring device to detect small differences in the quantity being measured.

*(b) Readability*

This is the susceptibility of the measuring device to have the indications converted to a meaningful number.

*(c) Calibration*

This is derived from the word calibre used to describe the size of gun bores. It refers to the disciplines necessary to control measuring systems to assure their functioning within prescribed accuracy objectives.

## 5.4 Basic components of accuracy evaluation

Factors affecting the calibration standard:

    traceability—interchangeability
    geometric compatibility
    thermal properties
    stability
    elastic properties
    position of use.

Factors affecting the workpiece:

    geometric truth
    related characteristics—surface texture
    elastic properties
    cleanliness, thermal effects, etc
    definition of what is to be measured.

Factors affecting the instrument:

    amplification check, inter-range switching
    filters
    effects of friction, backlash, etc
    contact of workpiece adequate
    slideways adequate
    readability adequate.

Factors affecting the person:

    training
    skill
    sense of precision
    cost appreciation
    hygiene—rusty fingers.

Factors affecting the environment:

    thermal (20°C)
    sunlight, draughts, cycles of temperature control
    manual handling
    cleanliness
    adequate lighting.

Although surface roughness calibration standards are nothing like as formalized as length standards it is useful to list the hierarchy of standards for any given country and, in particular, length standards.

Classification of calibration standards:

international
national
national reference standards
working standards and laboratory reference standards (interlaboratory standards).

Calibration chain:

initial procedure
determine and set product tolerance—go-no-go
calibrate product measuring system—calibrate gauges.

Time checks to carry out procedure:

1–3 weeks   Calibrate measuring system used to calibrate product measuring system.
1 year      Ensure working standards to laboratory standards.
1–5 years   Reference laboratory standards to national standards.

## 5.5   Basic error theory for a system

In general it is the task of the measurement to relate a number of output variables $x_o$ to a number of input variables, and this in the presence of a number of disturbances or noise terms $x_r$. This is the general format with which error signals will be determined as shown in figure 5.1. The measurement can be considered to be a mathematical operation which is not necessarily linear.



**Figure 5.1** System parameters in the presence of noise.

There are two important parameters [3]:

1. The number of measurement values $n$ that may be obtained in unit time at the output of the measuring instrument. This is related to the time $t_e$ required by the measuring instrument to reach the correct value from a transient input:

$$n = 1/t_e. \tag{5.5}$$

According to the Nyquist sampling theorem the following correlation exists with the highest frequency indicated by the measuring constraint $f_g$:

$$f_g = 1/2t_e. \tag{5.6}$$

2. The error $\varepsilon$ is a measure of the difference between the correct undisturbed output variable $x_o$ of an assumed ideal measuring system and the actual output $x_{o\,real}$.

So according to equation (5.3) a sensible way to describe errors is to relate the difference of $x_o$ and $x_{o\,real}$. This is shown in figure 5.2. Thus

$$\varepsilon = x_o - x_{o\,real} \qquad (5.7)$$

so

$$\bar{\varepsilon}^2 = \sum_{i=1}^{m} \left| x_{oi} - x_{oi\,real} \right|^2 \Big/ m \qquad (5.8)$$

to take into account a number $m$ of inputs and outputs.



**Figure 5.2** Error system.

There is a purely geometrical interpretation of this which has been used to describe errors in instrument systems. However, the purely Euclidean method will not be pursued here.

## 5.6 Propagation of errors

Measuring surface roughness or roundness in themselves rarely involves a knowledge of the propagation of errors because they are in effect 'stand-alone' measurements. But there are many parameters in surface metrology, which are derived from a number of others, such as cylindricity, squareness, etc. For this reason an elementary guide to error manipulation will be given here.

### 5.6.1 Deterministic errors

Consider the superposition of small errors. Take first the case of a single variable $y=f(x)$, say; for a small change in $x$, $\delta x$, $y$ will change by $\delta y$, so

$$\delta y = \delta x f'(x)$$

or

$$\frac{\delta y}{y} = \frac{f'(x)}{f(x)} \delta x \qquad (5.9)$$

where $f'(x)$ is the differential coefficient of $f(x)$.

For, say, three variables, equation (5.9) can be expanded to

$$u = f(x, y, z)$$

$$\delta u = \frac{\partial f}{\partial x} \partial x + \frac{\partial f}{\partial y} \partial y + \frac{\partial f}{\partial z} \partial z \tag{5.10}$$

where $\delta x$, $\delta y$ and $\delta z$ are small. Equation (5.10) is no more than the first term of Taylor's theorem.

On the other hand, if the values of $\delta x$, $\delta y$, $\delta z$ are not insignificant, equation (5.10) is not sufficiently accurate and other terms have to be incorporated into the form for $\delta u$.

Thus for two variables $u = f(x, y)$

$$\partial u = \frac{\partial f}{\partial x} \partial x + \frac{\partial f}{\partial y} \partial y + \frac{1}{2!} \left( (\delta x)^2 \frac{\partial^2 f}{\partial x^2} + 2\delta x \delta y \frac{\partial^2 f}{\partial x \partial y} + (\delta x)^2 \frac{\delta^2 f}{\delta y^2} \right)$$
$$+ \ldots + \frac{1}{m!} \left( (\delta x)^m \frac{\partial^m f}{\partial x^m} + m\delta x^{m-1} \frac{\partial^m f}{\partial x^{m-1} \partial y} + \ldots + (\delta y)^m \frac{\partial^m f}{\partial y^m} \right) \tag{5.11}$$

written in symbolic notation as

$$\sum_{i=1}^{\infty} \left( \delta x \frac{\partial}{\partial x} + \delta y \frac{\partial}{\partial y} \right)^i f. \tag{5.12}$$

Note that, for large errors, the constituent errors $\delta x$, $\delta y$, $\delta z$ are not additive as in equation (5.10) because of the presence of the cross-terms.

It is rarely necessary to go past two terms in the series, so for large errors and three variables $u = f(x, y, z)$

$$\delta u = \frac{\partial f}{\partial x} \delta x + \frac{\partial f}{\partial y} \delta y + \frac{\partial f}{\partial z} \delta z + \frac{1}{2!} \left( \delta x \frac{\partial}{\partial x} + \delta y \frac{\partial}{\partial y} + \delta z \frac{\partial}{\partial z} \right)^2 f. \tag{5.13}$$

In the case for multiplied or divided variables, such as for the volume of a box,

$$P = abc. \tag{5.14}$$

Taking logarithms of both sides of equations (5.14) gives

$$\log P = \log a + \log b + \log c \tag{5.15}$$

and for changes in $P$, $dP = d (\log P) = d (\log a) + d (\log b) + d (\log c)$

$$\frac{dP}{P} = \frac{\delta a}{a} + \frac{\delta b}{b} + \frac{\delta c}{c}. \tag{5.16}$$

In the general case $P = a^l b^m c^n$

$$\frac{dP}{P} = l \frac{\delta a}{a} + m \frac{\delta b}{b} + n \frac{\delta c}{c}. \tag{5.17}$$

This relationship is an example of differential logarithm errors. The important point is that it shows that the fractional errors add up (or subtract), because usually the worst case to consider is when the errors are considered to be additive.

### 5.6.2 Random errors

The question arises as to how to deal with random errors. In fact they can be treated in a similar way.

So, if the errors are small

$$\delta u = \frac{\partial f}{\partial x}\delta x + \frac{\partial f}{\partial y}\delta y + \dots . \tag{5.18}$$

If $\delta x$ and $\delta y$ are random variables they are represented in terms of a spread of values rather than by individual ones. The spread is usually given in terms of the standard deviation or variance. Thus for two variables

$$\text{variance}\,[\delta u] = E\left[\frac{\partial f}{\partial x}\delta x + \frac{\partial f}{\partial y}\delta y\right]^2 \tag{5.19}$$

where $E$ is the expectation (mean value). Equation (5.19) is approximately

$$S_u^2 = \left(\frac{\partial f}{\partial x}\right)^2 S_x^2 + \left(\frac{\partial f}{\partial y}\right)^2 S_y^2 + 2\left(\frac{\partial f}{\partial x}\right)\left(\frac{\partial f}{\partial y}\right)\rho S_x S_y . \tag{5.20}$$

Note that $S_u$, $S_x$, etc, are the standard deviations of $u$, $x$, etc, if $\partial f/\partial x$, $\partial f/\partial y$ are evaluated at the mean values $\bar{x}$, $\bar{y}$. $\rho$ is the correlation coefficient.

For $\rho = 0$, $\rho$ and $y$ are independent and then equation (5.20) becomes

$$S_u^2 = \left(\frac{\partial f}{\partial x}\right)^2 S_x^2 + \left(\frac{\partial f}{\partial y}\right)^2 S_y^2 \tag{5.21}$$

so if $u = x^a y^b$

$$S_u^2 = \bar{u}^2\left(a^2\frac{S_x^2}{\bar{x}^2} + b^2\frac{S_y^2}{\bar{y}^2}\right)$$

and if $r = (\bar{x}^2 + \bar{y}^2)^{1/2}$

$$S_r^2 = \frac{1}{(\bar{x}^2 + \bar{y}^2)}(\bar{x}^2 S_x^2 + \bar{y}^2 S_y^2)$$

where the bar indicates average value.

## 5.7 Some useful statistical tests for surface metrology

Because no workpiece is perfect, nor instrument for that matter, measured values are always subject to variations. The mean value of surface roughness or the spread of values is not realistically obtained with just one reading, despite the tendency to do so, because of the need for speed.

Hence it is necessary to qualify any value obtained within certain confidence limits. It is quite straightforward and the rules are easy to apply. In this section a few of the most useful tests and their possible applications are given. Furthermore, because the results are invariably obtained from a very limited set of readings, tests for small numbers should be employed. These can be found in any statistical textbook.

### 5.7.1 Confidence intervals for any parameter

The ideal situation is shown below. Suppose that limits are being set on a parameter. It should be possible to express the limits in the form

$$\text{prob}[L_1 \leqslant \text{statistical parameter} < L_u] = P = 1 - \alpha \qquad (5.22)$$

where $L_1$ is the lower limit, $L_u$ is the upper limit, $P$ is a given probability, say 0.95, and $\alpha$ is the signifance level.

The distribution of values of the statistical parameter would look like figure 5.3.

This is called a two-sided test because upper and lower limits are specified. If, for the sake of argument, only the upper limit had been specified the graph would be slightly changed (figure 5.4).

### 5.7.2 Tests for the mean value of a surface—the student t test

Suppose that there are $n$ independent readings with no restrictions on them. This is equivalent to the data set having $n$ degrees of freedom $v$.

If, for example, the values $x_1, x_2, \ldots, xn$ are subject to the constraint

$$a_1 x_1 + a_2 x_2 + \ldots + a_n x_n = 0$$



**Figure 5.3** Two-sided test.



**Figure 5.4** One-sided test.

there would be a constraint of one, so the number of degrees of freedom left is reduced by one, so $v = n-1$. If there are $m$ restrictions then $v = n-m$.

Hence, suppose the standard deviation of $n$ readings is wanted:

$$\sigma = \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - \bar{x})^2 \right)^{1/2} \tag{5.23}$$

has $n$ degrees of freedom if $\bar{x}$ is known previously, but if $x$ has to be worked out from the values of the $x$ this adds one constraint. Thus the formula (5.23) becomes

$$s = \left( \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2 \right)^{1/2} \tag{5.24}$$

which is more correct and is in fact unbiased (tends to the true value as $n \to \infty$). $s$ *is* the sample standard deviation.

Reverting back to the confidence limits, it so happens that a 'statistical parameter' containing means is given by

$$\frac{\mu - \bar{x}}{s/\sqrt{n}}$$

follows a $t$ distribution where $\mu$ is the true mean, $\bar{x}$ is the mean from the sample of $n$ readings, and $s$ is the sample standard deviation given by equation (5.24).

So it is possible to plot this function as if it were the 'statistical parameter' in equation (5.22):

$$\text{Prob}\left[ t_1 \leq \frac{\mu - \bar{x}}{s/\sqrt{n}} \leq t_u \right] = P = 1 - \alpha. \tag{5.25}$$

Rearranging

$$\left[ \bar{x} - \frac{t_1 s}{\sqrt{n}} \leq \mu \leq t_u \frac{s}{\sqrt{n}} \right] = P = 1 - \alpha. \tag{5.26}$$

Thus the true value of the mean ( lies within $x - t_1 s/\sqrt{n}$ and $x + t_u s/\sqrt{n}$ to within a probability of $P$ or with a significance of $\alpha$; $x - t_1 s/\sqrt{n}$ *is* the lower confidence limit and $x + t_u s/\sqrt{n}$ is the upper limit. It just so happens that the distribution of $t$ is symmetrical, so $t_1 = t_u$ except for the sign.

As $n$ becomes large so the value of the $t$ distribution becomes closer to the Gaussian. Therefore, for $P=0.95$ and $n$ large the value of $t$ becomes near to 2 which is the value for 95% confidence from the mean of a Gaussian distribution. When the value of $t$ is needed as in (5.25) it is looked up in tables with arguments $\alpha/2$ and $v$ degrees of freedom.

This test is often used to decide whether or not a sample or standard of roughness or roundness as received is the same as that sent off. This is not simply a question of bad postage. In many 'round robins' of standards committees a lot of time has been spent measuring the wrong specimen.

Under these circumstances the null hypothesis is used, which says in this case that the mean value of the received specimen $\bar{x}_1$ minus that of the sent specimen $\bar{x}_2$ should be zero.

So if $n_1$ measurements have been made on one specimen and its standard deviation evaluated as $s_1$ and $\bar{x}_2$, $s_2$ and $n_2$ are the corresponding values for the questionable sample. The $t$ value will be given by

$$t = \frac{\bar{x}_1 - \bar{x}_2}{s/[n_1 n_2/(n_1 + n_2)]^{1/2}} \tag{5.27}$$

where

$$s = \frac{\sum_{i=1}^{n_1}(x_{1i} - \bar{x}_1)^2 + \sum_{i=1}^{n_2}(x_{2j} - \bar{x}_2)^2}{n_1 + n_2 - 2}. \tag{5.28}$$

Equation (5.27) is really testing the difference between the means $x_1$ - $x_2$ as a ratio of the standard error of the mean. This latter term is the standard deviation of the distribution which could be expected with samples having joint degrees of freedom of $n_1 + n_2 - 2$ (i.e. $n_1 + n_2$ readings).

To work this out, if the measured $t$ value obtained by inserting all the numerical data into the right-hand sides of equations (5.27) and (5.28) is larger than the value of $t$ looked up in the tables corresponding to a value of $P$ and $v$, then the null hypothesis would be rejected—the specimens are different. If it is less then they are not different—to within the significance $\alpha$ (i.e. $1 - P$).

This can also be used to test if an instrument is giving the same answers from a specimen under different conditions. An example might be to test whether an instrument designed to measure a specimen off-line gives the same answers (to a known significance) when it is being used in process.

Note that the usual criterion for small $n$ numbers is about 25. Anything less than this requires the application of the $t$ and other tests. Gaussian or binomial approximations should only be used for large values of $n$, which rarely happens in surface metrology.

### 5.7.3 Tests for the standard deviation — the $\chi^2$ test

The quantity $ns^2/\sigma^2$ is distributed as a chi-squared ($\chi^2$) distribution with $v = n - 1$ degrees of freedom, where $\sigma$ is the true standard deviation, $s$ is the sample standard deviation and $n$ is the number of readings. The confidence limits of $\sigma$ in terms of $s$ can be derived from this $\chi^2$ distribution:

$$\text{Prob}\left[\chi^2_{1-\alpha/2,n-1} \leqslant \frac{(n-1)s^2}{\sigma^2} < \chi^2_{\alpha/2,n-1}\right] = P = 1 - \alpha \tag{5.29}$$

from which

$$\frac{(n-1)s^2}{\chi^2_{\alpha/2,n-1}} \leqslant \sigma^2 < \frac{(n-1)s^2}{\chi^2_{1-\alpha/2,n-1}} \tag{5.30}$$

to within $P$ probability.

So the spread of values to be expected on a surface $\sigma$ can be estimated from the measured value of $s$. The value of the degrees of freedom $v = n - 1$ is because $s$ is worked out using the evaluated mean.

This test is particularly useful for determining the quality of roughness standards where the spread rather than the actual value can be the deciding factor.

As for the $t$ test for $n$ large, the $\chi^2$ value becomes compatible with the standard error of the standard deviation obtained using a Gaussian distribution.

### 5.7.4 Goodness of fit

Another use of the $\chi^2$ test is to check out the type of distribution found in an experiment. This is usually called the 'goodness of fit' test.

Suppose that it is required to check whether a set of readings falls into a distribution pattern. It could be, for example, that someone needs to determine whether a surface has a Gaussian height distribution so that the behaviour of light or friction or whatever can be checked against a theoretical model. The $\chi^2$ distribution can be used to test this.

The method is to compute a statistic which in fact takes a $\chi^2$ value. The statistic $\chi^2$ is given by

$$\chi^2 = \sum_{i=1}^{m} \frac{(O_i - E_i)^2}{E_i} \qquad (5.31)$$

where $O_i$ is the observed frequency of observations falling into the $i$th class (or interval) and $E_i$ is the expected value assuming that the distribution is correct. From (5.31), if the distribution has been correctly chosen, the sum of squares $O_i - E_i$ will be close to zero. If incorrect, the value of $\chi^2$ will be large, so this very simple test can be used.

The degrees of freedom which have to be applied with this test are $m - 1$ because it is important that the total number in the observed classes equal that in the expected classes that is $\sum O_i = \sum E_i = m$.

### 5.7.5 Tests for variance—the F test

This is an alternative to the $t$ test for checking the validity of data and is based more on a comparison of variances rather than a comparison of mean values.

Thus if two samples are of size $n_1$ and $n_2$ and have variances $s_1^2$ and $s_2^2$ respectively, the variance ratio test $F$ *is* given by

$$F \sim \frac{\sigma_1^2}{\sigma_2^2} \sim \frac{s_1^2}{s_2^2} \qquad (5.32)$$

where values of $F$ are for chosen levels of significance $\alpha$. The degrees of freedom to be looked up are $v_1 = n_1 - 1$ and $v_2 = n_2 - 1$.

The two samples are considered to differ considerably in variance if they yield a value of $F$ greater than that given by the table at the chosen level of significance.

### 5.7.6 Measurement of relevance—factorial design

In function and in manufacture the influence of surfaces is often not well understood. It is therefore very relevant before investigating the importance of surface metrology to get an idea of how to test the importance. In practice this is not a very easy task because of the interdependence of parameters.

Perhaps the most difficult problem is to get the first idea of importance: what is the dominant factor? One easy way to tell is to use a factorial design. This gives an indication of which variable is worth concentrating on. An example might be the frictional force in a shaft in which the variables are oil, speed and texture amongst others. It is cost effective to spend effort, money and time on the most significant of the parameters. Factorial design is meant to give this vital clue.

Many experimental situations require an examination of the effects of varying two or more factors such as oil, texture and speed. A complete exploration of such a situation is not revealed by varying each factor one at a time. All combinations of the different factor levels must be examined in order to extract the effect of each factor and the possible ways in which each factor may be modified by the variation of the others. Two things need to be done:

1. Decide on the set of factors which are most likely to influence the outcome of the experiment.
2. Decide the number of levels each one will use [4].

Obviously it is best to restrict such pilot investigations to the minimum number of experiments possible, so it is often best to take two values of each factor, one high and one low. These are usually chosen such that

the influence on the outcome in between the two levels is as linear as possible. Non-linearity can spoil the sensitivity of the experiment. Two-level experiments are used most often and result in what is called factorial 2 design of experiments. There are of course factorial 3 or higher levels that can be used. The idea is that once the dominant factor has been identified more refined experiments can then be carried out using many more than just two levels for this dominant factor.

Take as an example a three-factor two-level experiment. The outcome might be the surface roughness of the finished component and the factors might be cutting speed, depth of cut and feed. There are $2^3$ possible experiments, equal to eight in all. The following shows how the dominant factor is evaluated. The experiments are listed in table 5.1.

**Table 5.1**

| Experiment | Algebraic assignment | Factor $A$ | Factor $B$ | Factor $C$ | Result |
|---|---|---|---|---|---|
| 1 | (1) | Low | Low | Low | $x_1$ |
| 2 | $a$ | High | Low | Low | $x_2$ |
| 3 | $b$ | Low | High | Low | $x_3$ |
| 4 | $ab$ | High | High | Low | $x_4$ |
| 5 | $c$ | Low | Low | High | $x_5$ |
| 6 | $ac$ | High | Low | High | $x_6$ |
| 7 | $bc$ | Low | High | High | $x_7$ |
| 8 | $abc$ | High | High | High | $x_8$ |

In going from experiments 1–2, 3–4, 5–6, 7–8, only factor $A$ is altered, so the total effect of $A$ can be found by comparisons of $(x_1 + x_3 + x_5 + x_7)$ with $(x_2 + x_4 + x_6 + x_8)$. All the former have factor $A$ low and the others high. Similarly the effect of factor $B$ is found by comparisons.

$$(x_1 + x_2 + x_5 + x_6) \quad \text{with} \quad (x_3 + x_4 + x_7 + x_8)$$
$$\downarrow \qquad\qquad\qquad\qquad \downarrow$$
$$\text{low} \qquad\qquad\qquad\qquad \text{high}$$

and $C$

$$(x_1 + x_2 + x_3 + x_4) \quad \text{with} \quad (x_5 + x_6 + x_7 + x_8)$$
$$\downarrow \qquad\qquad\qquad\qquad \downarrow$$
$$\text{low} \qquad\qquad\qquad\qquad \text{high}$$

while comparing $x_3 + x_7$ with $x_4 + x_8$ gives the effect of $A$ when $B$ is high etc. It is convenient to use Yates's symbolism to express these experiments more simply.

Consider the effect of $A$. The effects of all experiments with $a$ ($A$ high) and not $a$ ($A$ low) are $a$, $ab$, $ac$, $abc$ with (1), $b$, $bc$, $c$. (Here (1) does not mean unity—it means all factors are low.) The main effect of $A$ is

$$A = \tfrac{1}{4}(a + ab + ac + abc) - \tfrac{1}{4}[(1) + b + c + bc] \tag{5.33}$$

which can be written purely algebraically as

$$A = \tfrac{1}{4}(a - 1)(b + 1)(c + 1) \tag{5.34}$$

where each bracket is not meant to stand alone. For example, $a-1$ does not mean the result of experiment 2 with the value of unity taken from it; equation (5.34) only means something when multiplied out.

Similarly, the effects of $B$ and $C$ are

$$B = \tfrac{1}{4}(a+1)(b-1)(c+1)$$
$$C = \tfrac{1}{4}(a+1)(b+1)(c-1). \tag{5.35}$$

Note that the factor '?' say which is being investigated has the bracket $(? - 1)$ associated with it.

### 5.7.6.1 The interactions

The interaction between $A$ and $B$ is usually designated by the symbol $AB$ and is defined as the average of the difference between the effect of $A$ when $B$ is high and the effect of $A$ when $B$ is low.

The effect of $A$ with $B$ high is

$$\tfrac{1}{2}(abc + ab) - \tfrac{1}{2}(bc + b) = \tfrac{1}{2}b(a-1)(c-1) \tag{5.36}$$

and the effect of $A$ with $B$ low is

$$\tfrac{1}{2}(ac + c) - \tfrac{1}{2}(c + (1)) = \tfrac{1}{2}(a-1)(c+1). \tag{5.37}$$

The difference between (5.36) and (5.37) is defined as twice the value of the interaction, so

$$AB = \tfrac{1}{4}(a-1)(b-1)(c+1)$$
$$AC = \tfrac{1}{4}(a-1)(b+1)(c-1) \tag{5.38}$$
$$BC = \tfrac{1}{4}(a+1)(b-1)(c-1)$$

The interaction $ABC$ (half the difference between the interaction $AB$ when $C$ is high and when $C$ is low) is given by symmetry as $ABC$ where

$$ABC = \tfrac{1}{4}(a-1)(b-1)(c-1). \tag{5.39}$$

The reason for bringing the half into the interactions is to make all seven expressions similar. Again notice that there should be $(2^N - 1)$ equations for a $2^N$ factorial set of experiments.

The general $2^N$ experiments of factors $ABC \ldots Q$ are

$$A = (\tfrac{1}{2})^{n-1}(a-1)(b+1)\ldots(q+1)$$
$$AB = (\tfrac{1}{2})^{n-1}(a-1)(b-1)\ldots(q+1) \tag{5.40}$$
$$AB\ldots Q = (\tfrac{1}{2})^{n-1}(a-1)(b-1)\ldots(q-1)$$

Because this is so important consider an example of surface roughness and grinding. Take two factors as grain size $g$ and depth of cut $d$. The four factorial experiments might produce the results shown in table 5.2.

**Table 5.2**

| Experiments | Roughness |
| --- | --- |
| (1) | 20.6 |
| $g$ | 26.5 |
| $g$ | 23.6 |
| $dg$ | 32.5 |

The effect of $g$ and $d$ is

$$g = \tfrac{1}{2}(g-1)(d+1) = \tfrac{1}{2}(gd - d + g - (1)) = 7.4$$
$$g = \tfrac{1}{2}(g+1)(d-1) = \tfrac{1}{2}(gd + d - g - (1)) = 4.5$$

where 7.4 and 4.5 are just numbers of merit. Hence the grain size is about twice as important as the depth of cut in determining surface texture value.

The factorial design of experiments gives an indication of the important or dominant factors. Obviously if the relevant factor has been missed out then it is no help at all. Under these circumstances — and in any case — it is useful to do more than one set of experiments. If this is done then it is possible, using the $t$ test, to test the significance of the results obtained. This process is called replication. The difference between each set of experiments is carried out to see if the difference is significant. If it is, either the experiments have been badly performed or there is another factor which is present but has not been identified [3].

An alternative to the factorial design is the Taguchi method described in quality control tests. This arranges inputs and outputs to experiments in an orthogonal matrix method. The technique identifies the critical parameters but does not allow correlations between parameters to be established as in the factorial method described above.

Another test that is often carried out and has been referred to earlier concerns correlation. This is not the autocorrelation function but the simple correlation between two variables (figure 5.5).



**Figure 5.5** Correlation coefficient between x and y.

The correlation coefficient between $x$ and $y$ is given by $\rho$ where

$$\rho = \frac{\sum_{i=1}^{N}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \tag{5.41}$$

### 5.7.7 Lines of regression

Earlier in the book the idea and use of best-fit curves and lines has been employed extensively, but for those unfamiliar with the technique a few comments are given below.

The best-fit line between $y$ and $x$ and *vice versa* depends on which direction the differences or residuals between the actual line and the points are measured. It has been indicated when looking at least-square lines for surface analysis that, strictly speaking, the residual distance should be measured *normal to the line itself.* This is rarely needed, but as can be seen from figures 5.6 and 5.7 there are a number of possibilities.

**Figure 5.6** Lines of regression.



**Figure 5.7** Deviations from best-fit line.

For $y$ on $x$ the slope $m$ is given by

$$m = \frac{N\sum xy - \sum x \sum y}{N\sum x^2 - \left(\sum x\right)^2}.$$ 

(5.42)

For $x$ on $y$ the slope $m'$ is given by

$$m' = \frac{N\sum xy - \sum x \sum y}{N\sum y^2 - \left(\sum y\right)^2}.$$ 

(5.43)

The values of $x$ and $y$ here have not been altered to extract the mean values. It is straightforward to show that the two lines are not the same and that they cross at $(\bar{x}, \bar{y})$, the mean values of $x$ and $y$.

It is conventional to plot the regression with the dependent variable on the independent variable ($y$ on $x$).

The value of such an analysis is important, for instance if the actual dependence or relationship of one parameter with another is being tested.

### 5.7.8 *Methods of discrimination*

It is often required to determine which parameter best describes the functional performance of a workpiece surface. Quite often theoretical predictions are unsatisfactory and lacking in some way. In these circumstances it is possible to use a number of statistical tests to ascertain from practical results which parameter is best. One such method is called 'discriminant analysis'. Suppose there are a number of surfaces which have been tested for performance and have been separated out into 'good' and 'bad' surfaces depending on whether or not they

fulfilled the function. Discriminant analysis is a method whereby as many surface parameters as is deemed sensible are measured on all the surfaces. The idea then is to find which parameter best separates the good surfaces from the bad. If there are $n$ parameters from $m$ surfaces they can be considered to be $m$ points in $n$-dimensional space. Weights are then assigned to each parameter in such a way that the two groups 'good' and 'bad' form different clusters. The weights are manipulated, in much the same way as in neural networks, until the ratio of between-cluster separation to within-cluster separation is maximized. If this is successful the weights form the coefficients of a 'discriminant function,' which separates the function in one-dimensional space (see [5]). Discriminant analysis is another embodiment of the analysis of variance (ANOVA) technique. Dividing the between-cluster to within-cluster sum of squares is a measure of the discrimination. Thomas [5] uses the symbol $\Lambda$ for this. The larger $\Lambda$ is the better the discrimination. The method described above applies to just one parameter; it can be extended by means of linear discriminant functions. In this approach, the various parameters are combined as a weighted sum and the weighted sum and the weights are adjusted to maximize $\Lambda$. The techniques in effect define $\Lambda$ in terms of the coefficients $a_i$ of a linear expression. The ratio is then differentiated with respect to $a_i$ and solved for the values of $a_i$ which maximize the ratio. This is exactly equivalent to solving the normal equations in the least-squares routine.

The discrimination used in the example of beta functions in section 2.1.7.3 uses the discriminant function approach of combinations of surface parameters. In this case, however, the weightings are kept at unity.

The linear discriminant function in effect reduces to finding the eigenvectors of a non-symmetric matrix. The eigenvector is composed of the number of parameters in the linear expansions and the associated eigenvalue is a measure of the effectiveness of the particular discriminant function.

As in the use of the factorial experiments of section 5.7.6, discriminant analysis or the ANOVA methods are only meant to be an aid to the investigation and not an end in themselves.

## 5.8 Uncertainty in instruments—calibration in general

The practical assessment of surface roughness or any of the other surface metrology parameters involves directly or by implication:

(1) recognition of the nature of the signal;
(2) acceptance of measurements based on various approximations to the true surface;
(3) acceptance of a sampling procedure;
(4) an appreciation of the nature of the parameters employed and the methods and instruments for carrying them into effect;
(5) the establishment of calibrating procedures and tests for accuracy.

In this section the sources of error in instruments will be considered. The obvious place to start is at the input. Because most instrumentation in surface metrology at present revolves around the stylus method this will be taken as the basis from which to compare the others. The stylus, then the transducer and, finally, the amplifier will be dealt with in order. Then some procedures of sampling, including some idea of national variations, will be given. Obviously some information will be outdated because the improvements in materials and devices will continue to determine the best value of signal-to-noise ratio that can be achieved.

The best way to develop the test procedure is to establish the uncertainty of the specimen rather than the instrument.

A possible procedure might be as follows:

1. Visual inspection of the workpieces to determine where it is obvious that inspection by more precise methods is unnecessary, for example because the roughness is clearly better or obviously worse than that specified or because a surface defect which substantially influences the function of the surface is present.

2. If the visual test does not allow a decision, tactile and visual comparison with roughness specimens should be carried out.
3. If a comparison does not allow a decision to be made, measurement should be performed as follows. The measurement should be carried out on that part of the surface on which the critical values can be expected, according to visual inspection.

(a) Where the indicated parameter symbol on the drawing does not have the index 'max' attached to it. The surface will be accepted and the test procedure stopped if one of the following criteria is satisfied. Which one of the following to use depends on prior practice or agreement:
   • the first measurement value does not exceed 70% of the specified value (indicated on the drawing);
   • the first three measured values do not exceed the specified value;
   • not more than one of the first six measured values exceeds the specified value;
   • not more than two of the first 12 measured values exceed the specified value.
Sometimes, for example before rejecting high-value workpieces, more than 12 measurements may be taken (e.g. 25 measurements), up to four of which may exceed the specified value. (Corresponding to the transition from small numbers to large numbers of samples.)

(b) Where the indicated parameter symbol does contain the index 'max' initially at least three measurements are usually taken from that part of the surface from which the highest values are expected (e.g. where a particularly deep groove is visible) or equally spaced if the surface looks homogeneous.

The most reliable results of surface roughness inspection (as in roundness and other surface metrology parameters) are best achieved with the help of measuring instruments. Therefore rules and procedures for inspection of the most important details can be followed with the use of measuring instruments from the very beginning, although, of course, the procedure may be more expensive. But if the latter is to be considered then the instrument procedure has to be agreed.

Before embarking on the actual calibration of the instrument and its associated software, it is useful to remember that many errors are straightforward and in fact can be avoided. These are instrument-related errors which in some cases are just misleading. A typical set of things which should be checked are as follows:

1. Motor drive: check worn, loose or dirty parts, connections, switches, misaligned motor mounts, excessive bearing wear.
2. Link between drive and tracer: make sure locking joints are solid, rotating joints free and that the linkage is connected to the drive correctly.
3. Tracer itself: make sure the stylus is not loose or not badly damaged. Check any electrical connections—the most likely source of error.
4. Amplifier: batteries if used should be charged. Any meter should not stick. No worn or defective power cable or switches should be allowed.
5. Workpiece: this should be firmly mounted preferably on some solid mechanical base. The workpiece should be as free from dirt as possible and the area chosen for examination should be as representative of the surface as possible. In the case of stylus instruments some preliminary check should be made to ensure that the stylus pressure (or that of the skid) will not damage the workpiece.
6. Readings: if the output is on a non-digital scale remember that the eye always tends to read towards a division mark. If the output is multiscale select the relevant scale and if necessary, practice reading it.

These and other common-sense points can save a lot of unnecessary trouble later on.

The next section will deal with conventional stylus instruments but it will be necessary to consider scanning probe microscopes from time to time as these have some features in common with ordinary stylus instruments.

## 5.9 The calibration of stylus instruments

Most stylus instruments involving electronic amplification are provided with at least one adjustable control for the amplification. Often there are several. While these may be set by the maker before despatch, they must generally be checked and, if necessary, readjusted by the user after setting up the instrument, and as often thereafter as may be required to maintain the standard of performance. Such controls often have the form of screwdriver-operated potentiometers.

Calibration procedures must be suited to the features of the instrument being calibrated. The following features in various combinations will be considered here:

1. Pick-ups with displacement-sensitive transducers able to portray steps (e.g. modulated carrier systems) cooperating with an auxiliary datum device (commonly known as skidless pick-ups).
2. Pick-ups with displacement-sensitive transducers having a skid cooperating with the surface under test (commonly known as skid-type pick-ups).
3. Skid-type pick-ups with motion-sensitive transducers not able to portray steps (e.g. moving coil or piezo).
4. Graphic recording means.
5. Digital recording means.
6. Meters giving $R_a$ (and sometimes other parameters).

Difficulties may sometimes arise, for example when there is too much floor vibration to use an instrument with a skidless pick-up, or when the skid normally fitted has too small a radius to give good results on a calibration specimen (i.e. it does not integrate properly). Some degree of special adaptation may then be required.

Most instruments offer a range of magnification values for a recorder, and/or of full-scale values for a meter, these being selected by a switch. There may be an individual gain control for each switch position, or one for a group of positions, and often there will be separate controls for the recorder and the meter. When there is one gain control for a group of switch positions, switching errors may arise, and the adjustment may then be made for one preferred position, or averaged over the group. In all cases, the gain at each switch position must hold throughout a range of wavelengths to the extent required by the standard.

Apart from the adjustment of gain, about which this comment is mainly concerned, and the residual inter-range switching errors, errors may result from the non-linearity of circuits, recorders and meters. Such errors are dependent on the deflection of the pointer, if relevant, and will generally be least in the region between 65% and 80% of full scale, where calibration of gain is to be preferred. They are inherent in the apparatus and, apart from verifying that they are exceptionally small, generally there is nothing that can be done to correct them. Internal vibration from actuating motors and electrical fluctuations (both classed as 'noise'), may give rise to further errors, especially at high magnifications. In addition, charts and meter scales are subject to reading errors.

It will be evident that a complete calibration throughout the whole range of amplitudes and wavelengths for which the instrument is designed can involve a somewhat lengthy procedure, so much so that, at least in the workshop, it is usual to accept calibration on the basis of a limited number of spot checks, and for the rest to assume that when the equipment is right for these, it is right for its whole range of operation. While this approach has given generally acceptable control, the trend towards greater accuracy has been urging manufacturers and standardizing bodies to evolve a more comprehensive system of calibration.

Two forms of calibration procedure are recognized. One is basic and complete but costly: it involves the direct evaluation of the magnification of recorded profiles, either analogue or digital, followed by accurate assessment of their parameters from the records. With the aid of digital recording and computer techniques this form has now been raised to a very high level. The other procedure, suitable for the workshop, is based on the use of instrument calibration specimens which, in the context of surface metrology, can be likened to

gauge blocks in the field of dimensional metrology. Unfortunately these specimens, in the range of values ideally required, are still in a state of evolution, and it is for this reason that the range of workshop testing has so far been somewhat restricted.

### 5.9.1 Stylus calibration

Calibration should start by making sure that the stylus tip is in good order. Pending the evolution of instrument calibration specimens suitable for checking the sharpest tips, sufficient inspection to ensure that the tip is not damaged can generally be carried out with a microscope, but as will be seen this is often not sufficient in side view and plan view. The former requires an objective having an adequate working distance generally of at least 3 mm, which will be provided by most 16 mm and many 89 mm objectives. A magnification of $100 \times$ to $300 \times$ will generally suffice. Some fixturing may be required to mount the pick-up on the stage of a microscope. If the tip is seen to require cleaning, great care will have to be taken to avoid damaging the tip or pick-up. Brushing with a sable hair brush may suffice, but if not, the tip can be pressed lightly (with perhaps 10 mg force) into a soft material such as a stick of elder pith or even the end of a matchstick, taking care to avoid destructive shear stresses such as could result from scraping the stick across the tip under such a load.

This point about cleaning is valid for all methods of assessment. It may seem very crude but it is extremely effective.

Referring to chapter 4 in the section on tactile sensors, it is clear that the observed signal is very dependent on the geometry of the stylus. Changing from a 2 $\mu$m tip to a 10 $\mu$m tip can influence the result by almost 10% on a fine surface. The same is true for the resolution spot of an optical probe or the tip of a capacitor electrode. So methods of examining the stylus are absolutely crucial from the point of view of calibration and verification of instruments.

Consider the verification of the stylus. The obvious way is to examine the tip under a microscope. Methods of doing this are not quite as easy as may be supposed. The reason is that the typical tip has dimensions that are very close to the wavelength of light. This poses two problems. The first is that it needs a reasonably good microscope to view the dimension and the second is that the method is hardly practical for instruments that do not have a detachable stylus or pick-up unit.

Methods are available for the investigation of the wear of a stylus during an experiment where the pick-up from a conventional instrument is separately mounted and can be swung around out of its normal position for measurement, but in general this technique is not available except in research work.

The SEM has been used to measure styli but it is not easy. In principle it should give an accurate, highly resolved image of the tip profile. However, for small-radii tips the styluses should be coated with a conducting film to stop the surface charging up. Carbon has been used for this, as also has palladium-gold alloy which is coated to a depth of 100 angstroms. Also the magnification of the images must be calibrated since it can vary up to 20% from one pumpdown to the next. The calibration is usually achieved by means of a calibrated line standard. The stylus and the calibrated line standard have to be mounted in the same rotary sample holder. During operation [5] the line scale is first rotated under the SEM beam to enable the calibration to be determined.

The TEM can also be used as described in the section 4.8.4 on TEMs but this requires a double replica method which is a messy procedure but can be used as an effective check.

Optical micrographs of the stylus are an alternative but the resolution is poor and any structure less than about 1 $\mu$m is likely to be undefinable. Using a numerical aperture of about 0.9 is desirable but the short field depth needed to tell where the shank starts and the diamond finishes precludes it for spherical tips. Flat tips are not so much of a problem as spherical tips because it is the flat dimension only that has to be determined and not the shape. For spherical tips the shape has to be determined in order to evaluate the radius. This involves the use of algorithms from which to get the radius from the profile of the stylus. The use of optical

microscopes for the determination of the dimensions of spherical tips is therefore restricted to 0.3 NA to get a suitably large depth of field. This gives a resolution of worse than 2 $\mu$m and so limits the technique to 10 $\mu$m tips. Also, if a high-resolution microscope is used there is a serious possibility of pushing the stylus into the objective lens.

For optical methods, as in the SEM, the stylus has to be removed from the instrument, although not necessarily the pick-up, so it is more convenient but less accurate than the SEM and TEM.

One practical method is to run the stylus, whilst in the instrument, over a very sharp object such as a razor blade, and to record the profile on the chart of the instrument. Razors have a tip of about 0.1 $\mu$m and an angle of a few degrees and so are suited to measure styluses of about 1 $\mu$m. It has been reported [7] that there is very little deformation of the blade. In the same reference it is highlighted that traversing speeds have to be very low so that the frequency response of the chart recorder is not exceeded. Typical tracking speeds are 0.001-0.01 mm s$^{-1}$. Also, it is advisable to keep the vertical and horizontal magnifications of the chart equal so that there is no ambiguity as to the true shape of the stylus. For accurate assessment using this method, the blade has to be supported very close to its edge (~ 0.5 mm) to avoid deflection.

Once the profile of the stylus has been obtained there remains the problem of evaluating it. If the instrument has a digital output then problems of algorithms and frequency responses have to be agreed. If the chart record is the only output it is convenient to try to fit a template to the graph inscribed with different radii [6].

Ideally, in industry, such complicated procedures should be avoided. There are a number of indirect methods which have been used in the past. The first one used a roughness standard for calibrating the meter to give an indication of stylus wear. This is perhaps the most important factor to consider from the point of view of a user. The form of the profile on these specimens is triangular and called the Caliblock, originally devised by General Motors. The standard has two sections: one is used for a rough surface amplitude calibration (see later) and the other is of about 0.5 $\mu$m $R_a$ for fine surfaces. It was the finer of the two parts which was and still is used to estimate wear of the diamond. The idea is that if the stylus is blunt it will not penetrate into the valleys of the triangle. This results in a reduction in the signal and hence the measured $R_a$ value. It is straightforward to calculate to a first approximation what loss in $R_a$ corresponds to the bluntness of the stylus. However, it is difficult to get quantitative information about the true shape of the tip. It could be argued that this is not needed and that all that is required is a go–no-go gauge. For many cases this is correct but there is a definite need for a more comprehensive, yet still simple, test. Furthermore this test should be performed on the instrument without any modification whatsoever.

One such standard is shown in figure 5.8. This standard comprises a number of slits etched or ruled into a block of glass or quartz.

Each groove has a different width and usually corresponds to one of the typical stylus dimensions (i.e. 10, 5, 2, 1 $\mu$m). Use of these standards enables the tip size and the shape and angle to be determined.

When the stylus is tracked across the standard a picture like figure 5.9 is produced in which the penetrations $P_1$, $P_2$, etc, are plotted against the dimension of the slot (figure 5.10). A straight line indicates a pyramidal stylus, of slope $m$. The stylus slope measurement

$$\frac{d/2}{P} = \tan \theta/2 \tag{5.44}$$

from which

$$P \sim d \cot \theta. \tag{5.45}$$

The penetration of a spherical stylus is approximately given by the spherometer formula (figure 5.11)

$$P = d^2/8R.$$

**Figure 5.8** Stylus calibration block



**Figure 5.9** Profile of stylus locus seen on chart.



**Figure 5.10** Graphical interpretation of chart readings.



**Figure 5.11** Penetration of styluses into grooves.

So, for example, on the graph the spherical stylus has a quadratic penetration and, incidentally, will always drop to some extent into a slot no matter what its width. This is not true for the flat-type stylus.

Problems with this type of standard include the fact that the grooves can fill up with debris and so cause the stylus to bottom prematurely and give a false impression of a blunt stylus. The opposite or negative of the standard has also been attempted giving a set of ridges corresponding to a set of razor blades. Although stiffer than a razor blade, the evaluation of the stylus shape is more difficult because the profile obtained is a sum of the stylus curvature and ridge curvature at every point. Extracting the stylus tip information from the profile is difficult, although debris is less of a problem than with grooves.

One or two extra problems emerge with the use of scanning probe microscopes. The stylus in ordinary instruments follows a given shape such as a cone or pyramid (or tetrahedron as in the case of the Berkovich stylus). These styli have a defined tip dimension and an angle (i.e. $r$ and $\alpha$) as seen in figure 5.12.



**Figure 5.12**

However, the fact that these styli have a finite and usually substantial slope (i.e. 60° or 90°) means that accurately measuring the position of an edge or side of a protein is virtually impossible. In figure 5.12(c) for example instead of B, the edge of the protein registering the contact on the shank C together with the large size of the stylus, indicates that the edge of the protein is at A and not B.

In biological applications error produced in this way can be misleading. This has resulted in attempts to minimise $r$ and eliminate $\alpha$ by attaching a carbon nanotube (by means of an adhesive) to the usual tungsten or silicon stylus used in AFM or STM (Fig. 5.12(b)) [8].

These nanotubes, members of the self-assembling carbon fullerine family, have typical dimensions of about a nanometre diameter and micrometre length. They are tubes and not rods and can be single or multi-walled. It would appear that the multiwalled, say, three layer, are stiffer than the single wall type but no information is available on the internal damping characteristics–damping is the correct word to describe energy loss at this scale. There are other types of linear molecule that could conceivably be used. For example, simple fibres rather than tubes. Silicon or molybdenum selenide fibres having similar dimensions are possibilities, with titanium or tungsten based probes also worthy of consideration [9].

To be of any practical use the tubes have to have the requisite lateral resolution and must be stable over a realistic lifespan. Nguyen et al [10] carried out some exhaustive tests on ultra-thin films of conductive and insulator surfaces in the 5-2 nm region. They also tested resolution by using metal films (e.g. gold and iridium) having grain boundary separation between 10 and 2 nm. Measuring resolution is very difficult at this small scale; it is hardly feasible to find a test artefact which has a range of grooves or spacings similar to that used for the microscale shown in figure 5.8.

Very hard surfaces such as silicon nitride were repeatedly scanned for long time periods to test for wear and degradation. Fortunately they showed that the carbon nanotube (CNT) probe on an AFM offers great potential for measurement down to the nanometre and subnanometre. Most important is the finding that these tips maintain lateral resolution, hence fidelity, over long periods of time. Even after fifteen hours of continuous scanning wear is minimal. It should be pointed out here however, that continuous scanning is not necessarily a valid operation. It is discontinuous operation which shows up weaknesses in bonding and other mechanical factors, not to mention operator error.

Also needed for a profilometer is the force on the surface. This is important from the point of view of potential damage. There are a number of ways of achieving the force. It suffices to measure the static force. Multiplying this by 2 gives the traceability criterion. The classic way of doing this is by using a small balance and simply adding weights until a null position is reached. For a laboratory the balance method is acceptable but as typical forces are 0.7 mN (70 mg) and lower it is much too sensitive for workshop use. One simple practical level [11] actually uses the profilometer for its own force gauge. This method is shown in figure 5.13(a).



**Figure 5.13**

Where the deflection $y(x)$ is a function of $x$ and cantilever constants $b$ for width and $d$ for thickness, and $l$ is the 2nd moment of area.

$$y(x) = Fx^3/3El \tag{5.47}$$

$y(x)$ can be found from the chart (figure 5.13(b)) or from digital output as the probe is tracked along the cantilever.

In practice the free cantilever can be replaced by a beam fixed at both ends. The deflections near the centre of the beam are small but stable. Using the equation (5.47) the force $F$ can be found and from it the maximum force imposed on the surface (i.e. $2F$ in the valleys). If possible it would be useful to determine the damping by using the log decrement method (i.e. letting the stylus dangle in a vertical position). Watching the decrease in oscillations from the chart enables the damped natural frequency to be found (the natural frequency is usually known).

### 5.9.2 Calibration of amplification

Magnification is defined as the ratio of the vertical displacement of the needle on a graph (or its digital equivalent) to that of the movement of the stylus (or equivalent tracer).

The basic requirement is to displace the stylus by accurately known amounts and observe the corresponding indication of the recorder or other output. The stylus can be displaced by means of gauge blocks wrung together. For best answers the blocks should be wrung onto an optical flat (figure 5.14). In some cases it is best to fix them permanently so that the steps can be calibrated by an interferometer rather than just relying on the specified values. According to Reason [12] the technique can be used directly for the calibration of

steps over 50 $\mu$m. There are problems, however, if small displacements are to be calibrated because of the calibration error of the gauge block itself, amounting to one or two microinches (25–50 nm), which is a large proportion of the step height. A convenient way to minimize the errors is to use an accurate mechanical lever. This can produce a reduction of 10 or even 20 times.



**Figure 5.14** Calibration of amplifier with gauge blocks.

Reason's lever [12, 13] was designed so that the position of the pivot could be accurately determined (figure 5.15).



**Figure 5.15** Calibration of amplifier with Reason lever arm.

The reduction ratio is determined by adjusting the position of the pivot P relative to the stylus (which is not traversed) by means of the screw S and gauge block G until equal deflections are obtained on both sides.

An alternative to the use of a lever arm and gauges is a precision wedge as shown in figure 5.16 [14].

A thick glass optical flat with parallel faces is supported in a cradle on three sapphire balls of closely equal diameter (two at one end and one at the other). The balls rest on lapped surfaces and a precise geometry is obtained by using a heavy base plate with a flat lapped surface onto which lapped blocks of gauge block quality are secured and on which the sapphire balls sit. The nominal wedge angle is 0.0068 rad.

A shielded probe with a spherical electrode is rigidly mounted to the base place and forms a three-terminal capacitor with a conducting layer on the undersurface of the glass flat. The micrometer is arranged to move the cradle over 12mm.

**Figure 5.16** Calibration of amplifier using sine bar method.

The displacement of the stylus can be worked out in two ways, one by knowing the horizontal shift and the wedge angle, and the other, better, method using the capacitance probe.

The relationship between electrode spacing and capacitance is

$$y = a + b \exp(\lambda c) \tag{5.48}$$

where $y$ is the gap, $a$, $b$ and $\lambda$ are constants.

A change in the gap is then related to a change in the capacitance by the differential of equation (5.48). Thus

$$\frac{\mathrm{d}y}{\mathrm{d}c} = b\lambda \exp(\lambda c) = K \exp(\lambda c) \tag{5.48a}$$

To relate the change in $c$ to $y$ requires $b\lambda$ to be known, so the probe has to be calibrated. This might sound like a snag but it turns out not to be so because it is possible to calibrate capacitance probes over relatively large displacements and then to interpolate down to small displacements with the same percentage accuracy. So the wedge is moved over its whole range yielding a height change of about $90 \pm 0.25$ $\mu$m.

The wedge angle is calibrated by using interferometry. The ultimate resolution of this technique has been claimed to be 2 nm.

Obviously there are other ways in which a step or vertical movement can be calibrated. The lever and wedge are just two of them. Of these the lever is much the cheaper and simpler but not quite as accurate. How does the lever work?

If $L$ is the length of the long arm engaging the gauge blocks and $2d$ the displacement of the unit from one stylus operating position to the other, the lever reduction is $d/L$. As an example, $L$ is $100 \pm 0.1$ mm and $2d$ is measured accurately with a 25mm or 10mm gauge block G.

The reported accuracy of such a technique is much better than 1%. But again the problem is that much smaller surface height and displacement steps are now being required. The lever method is limited to about $0.25$ $\mu$m steps for normal use but 25 nm has been measured with 1 nm accuracy. Unfortunately, this is barely sufficient for today's needs.

Obviously the accuracy of any gauge block technique depends in the first place on how accurately the gauges have been measured. The same is true for the measurement of the step itself by interference means. In this case there are two basic errors: one is that of uncertainty in measuring the fringe displacement, and the other is the variation of step height measurement when the geometry of the edge is imperfect.

Therefore proper measurement of a step height must take into account:

(1) the geometry on both sides of the step;
(2) the effect of fringe dispersion by imperfect geometry;
(3) the nature of the 'method divergence' between optical interference and stylus methods;
(4) the realization that the accuracy of the height value is ultimately limited by the surface roughness on both sides of the step.

This last point limits the overall accuracy of transferring a step height calibrated by interferometry to a stylus measurement (figure 5.9).

Thus examination of the system of measurement for a stylus and an interferometer can be made using figure 5.17.

Comparison between stylus instruments and optical instruments has been carried out by the PTB in Braunschweig [15]. Two useful optical methods were included in the investigation: the phase shift Mirau-type interferometer and the Linnik interference microscope. Measurements were carried out on very smooth rectangular-shaped groove standards with depths from 3 $\mu$m down to 60 nm, just inside the nanometre range. The steps were obtained by evaporating aluminium onto a zerodur plate using a rectangular aperture. Successively opening the aperture led to six steps forming a staircase of steps.

Due to the limited lateral resolution of the plane wave interferometers, it was found to be impossible to use the same lengths of profile as the stylus and interference microscope. Also high waviness of the upper step levels resulted in differences between the interferometer and the interference microscope. It is not clear whether relocation was used in this exercise. Table 5.3 shows the results taken with the Nanostep stylus instrument, the Zeiss IM/Abbé, the Zeiss IM/UBM and the Wyko IM (interference microscope).

**Table 5.3** Measurement with different instruments.

| Groove Number | Nanostep, nm | Zeiss IM/Abbé, nm | Zeiss IM/UBM, nm | Wyko IM, nm | Mean value, nm | s, nm |
|---|---|---|---|---|---|---|
| 8 | 2803.5 | 2806.9 | 2805.6 | 2799.1 | 2803.8 | 3.8 |
| 7 | 1838.8 | 1636.3 | 1838.4 | 1843.2 | 1837.8 | 1.4 |
| 6 | 1407.8 | 1407.4 | 1407.1 | 1402.9 | 1406.3 | 2.3 |
| 5 | 1001.9 | 1004.1 | 1003.9 | 1003.3 | 1003.3 | 1.0 |
| 4 | 485.9 | 484.8 | 488.5 | 487.5 | 486.7 | 1.6 |
| 3 | 240.5 | 236.9 | 242.1 | 245.2 | 241.2 | 3.4 |
| 2 | 95.4 | 95.0 | 96.7 | 93.2 | 95.1 | 1.4 |
| 1 | 61.2 | 59.5 | 61.7 | 58.8 | 60.3 | 1.4 |

The deviations expressed as a linear function of groove depth $d$ came within the estimated instrument uncertainties. At the $2\sigma$ level these were: Nanostep $1.5 \times 10^{-3} d + 0.7$ nm, Zeiss IM/Abbé $1.5 \times 10^{-3} d + 4$ nm, Zeiss IM/UBM $1.5 \times 10^{-3} d + 1.5$ nm and the Wyko IM $1.5 \times 10^{-3} d + 4$ nm. The constraint part is added to take into account the higher deviations at smaller grooves. Among these, the stylus method has smaller uncertainty probably because the stylus tends to integrate random effects. The experiments showed that the methods had gratifyingly close agreement as can be seen in the table. However, problems of surface contamination and phase jumps caused by the influence of varying extraneous film thicknesses, especially in the bottom of the grooves caused problems with all the optical methods as discussed in chapter 4. Polarization problems were not addressed.

**Figure 5.17** Calibration relative to fringe pattern.

Here the height $h$ of the step measured interferometrically as a function of the fringe spacing $t$ and dispersion $\Delta$ is

$$h = \left( n + \frac{\Delta}{t} \right) \frac{\lambda}{2}. \tag{5.49}$$

The value of the fringe $n$ can be obtained by following the fringe or by use of the method of 'exact fractions' [16] using different wavelengths of light.

The variability of $h$ as a function of $\Delta$ and $t$ can be obtained using the propagation formulae given in section 5.6 as

$$\mathrm{var}(h) = \left( \frac{\partial h}{\partial t} \right)^2 \sigma_t^2 + \left( \frac{\partial h}{\partial \Delta} \right)^2 \sigma_\Delta^2. \tag{5.50}$$

If the mean $t$ and mean $\Delta$ are $\bar{t}$ and $\bar{\Delta}$, the mean percentage error [14] is

$$\frac{\delta h}{h} = \frac{\delta h}{h} = \left( \frac{1}{1 + n\bar{t}/\bar{\Delta}} \right) \left[ \left( \frac{\sigma}{\bar{\Delta}} \right)^2 \left( \frac{\sigma}{\bar{t}} \right)^2 \right]^{1/2}. \tag{5.51}$$

This is again obtained from the error formulae making the assumption that $\sigma_\Delta = \sigma_t =$ standard deviation $\sigma$ of the measurement.

Errors due to geometry rather than simple fringe measurement errors shown in equation (5.49) can be incorporated by expanding the expression $s = t/\Delta$ into two terms of Taylor's theorem, yielding

$$\frac{\delta h}{h} = \frac{\sigma_h}{h} \left[ 1 + \left( \frac{\mathrm{d}(\delta s)}{\mathrm{d}(s)} \right)^2 \right]^{1/2} \tag{5.52}$$

where the first term represents equation (5.51) and the addition is the term arising from geometry and lateral averaging problems. Errors in step height of $0.1\lambda/2$ are not unknown. Generally the variation in height is strongly dependent on the particular geometry and its slope relative to that of the reference surface. Geometrical errors in the lower part of the step and in the reference surface would also contribute to the total uncertainty of the height measurement.

The effects of lateral averaging are particularly important if the interferometer magnification is at $10 \times$ or less. Practical verification of the typical step error for stylus instruments has indicated [14] that the error is about equal to the $R_q$ value of the surface texture. Remember that the need for a step to be calibrated is so that the magnification of the pick-up can be determined.

Notwithstanding the problems of diffraction, phase shift and polarization, there are difficulties in putting a number to what is, in effect, a simple step height. The reason is that the step itself, usually of nanometre height, can easily be buried in extraneous signal. Deviations from linearity of less than one per cent can result in long wavelength bows. Tilt of the specimen is another problem. Figure 5.18 shows an example:



**Figure 5.18** Actual step height in presence of noise.

In this case the step height value can be only a small component of the recorded signal (figure 5.18(a)). Removal of trends and 'bowing' leaves a very indistinct step as shown in figure 5.18(b). This is another example of noisy data as discussed in chapter 3 on processing. One attempt at resolving this problem is to fit a 'best fit step' to the data [17].

Thus if the data taken across the probe signal is $z(x)$ it can be represented by a power series given in equation (5.53) where

$$z(x) = Z(x) + C_0 + C_1 x + C_2 x^2 \ldots \ldots \tag{5.53}$$

$Z(x)$ is the desired step and $z(x)$ is the recorded data.
Evaluating the coefficients $C_0$, $C_1$ etc.

$$S = \sum (z_i - Z_i)^2 = \sum_{\text{top}} (z_i - Z_i)^2 + \sum_{\text{bottom}} (z_i - Z_i)^2 \tag{5.54}$$

The partitioning of the data in $x$ (according the shape required—the $Z_i$ values) need not be a step; it could be a staircase. The desired shape is stored in a look-up table and the value of $Z_i$ corresponding in $x$ to the $z_i$ measured value is brought out within the summation sign (5.54) in order to obtain the residuals. This method is an example of the partitioning developed in roundness evaluation for partial arcs [18]. One example of this is removing the effects of keyways in shafts. In this case there are two steps in opposition. The reason why this type of calculation is needed in SPM is because the step heights are very small when compared with gauge block steps used for calibration in general engineering. It is possible to scale down the relatively large steps attainable with gauge blocks by using a mechanical reduction lever system devised by Reason [12] as shown in figure 5.15.

It is fruitful to consider other ways of establishing height steps for calibration which involve absolute physical quantities such as the wavelength of light. The obvious alternative is to use the inviolate properties of crystal lattice spacings as a way of getting a known small displacement. This affects SPM calibration.

### 5.9.2.2 *Calibration of SPM*

In the previous section amplitude calibration in terms of step height has been discussed. In engineering surface metrology height is considered to be the most important feature. However, in scanning probe instruments position is considered to be most critical; presumably because of the requirement for high accuracy masks in lithography.

Positional accuracy in the $x$ and $y$ directions is usually guaranteed by means of laser interferometers, but the $z$ direction is more difficult. In all three axes hysteresis, non-linearities, drift motion errors, squareness and straightness of axes have to be eliminated or at least minimized. The beauty of SPM instruments is that they have been able to overcome the poor lateral resolution that is inherent in conventional optical scanners due to the laws of physical optics. Instead the resolution is limited by the probe dimension which can be of nanometre scale.

The advent of SPMs such as the scanning tunnelling microscope (STM) and the atomic force microscope (AFM) has greatly increased the need for high resolution of position along the $x$ and $y$ axes e.g. [20] which makes use of closed loop rather than open loop control.

Sometimes the $z$ axis uses a laser interferometer to calibrate height, or a capacitative system which is calibrated off-line with an interferometer.

Any instrument which uses laser interferometers to measure and control the three axes has the name 'metrological' attached to it. This means that all three axes' movement and position are traceable to the wavelength of light. This does not mean that the measurand is traceable, e.g., if the measurand is current it is not necessarily traceable to the standard AMP. Figure 5.19 shows the arrangement used in [21].



**Figure 5.19** Plane mirror interferometer 'metrological' system.

Movements of many micrometres in $x$ and $y$ are accompanied by subnanometre positional accuracy.

The $z$ direction is usually restricted to a few micrometres. $M_1$, $M_2$, $M_3$ are mirrors $S_1$, $S_2$, $S_3$ sources usually obtained from one laser suitably split. The drives have been ignored. In practice the interferometers are heterodyne especially in the $z$ of a 'Zygo' which in the system indicated uses parts of a 'Zygo' heterodyne interferometer.

The classical case of a metrology AFM is the $M^3$ system devised by Teague at NIST [22].

Although not specifically measuring surface roughness, this instrument does measure within the nanometric range where the distinction between roughness, form and dimensional characteristics becomes blurred.

The specification requires a spatial resolution of 0.1 nm of the distance between any two points within a $50 \text{ mm} \times 100 \, \mu\text{m}$ volume with a net uncertainty for point-to-point measurement of 1.0 nm.

The maximum specimen size is $50 \times 50 \times 12 \text{ mm}^3$. The system has been designed to operate at three levels. Large scale mapping ~ 100 mm$^2$, small scale 0.01 ~ 1 $\mu$m$^2$ and point-to-point measurement. Each mode of operation works with a different probe: with AFM for the small scale and confocal microscope for the large scale (figure 5.20).

**Figure 5.20** Schematic drawing of the MMM (M³) showing the metrology reference system and the X and Y interferometer arrangement. The interferometer beams for X and Y displacement straddle the probe tip to minimize Abbé offset errors.

**Table 5.4** M³ Uncertainties

| Calibration Type | Dimension Range | Uncertainty |
|---|---|---|
| Surface roughness | 1 nm to 3$\mu$m | 0.3 to 100 nm |
| Step height | 10 nm to 10 $\mu$m | 1.3 to 200 nm |
| Line spacing | 1 to 50 $\mu$m | 25 to 250 nm |
| Line spacing | 5 to 50 mm | 10 to 30 nm |
| Line width | 0.5 to 10 $\mu$m | 50 nm or less |
| M³ point to point | 100 nm to 50 nm | 1 nm |
| M³ height measure | 1 nm to 10 $\mu$m | 0.01 to 1nm |

Uncertainties have been determined as follows in table 5.4 taken from [22].

A schematic view of the metrology system is shown in figure 5.20. Detailed metrology and constructional issues are discussed elsewhere [51]. The main point is that the design philosophy is that of *repeatability*. The required accuracy is achieved by systematic error compensation of position and motion. For this reason great attention is paid to thermal and vibrational stability. The basic metrology frame of position between probe and specimen is provided by using ultra high resolution heterodyne interferometers having a four pass configuration to get λ/8 displacement resolution which, when interpolated, gives the required repeatability. The differential mode of operation of this interferometer system reduces the effect of angular errors considerably. All design methodologies, such as those described in chapter 4, have been taken into account in the design of this extraordinary instrument. It has squeezed the last 1% out of existing technology. However, it may well be that advances in physics rather than technology will provide the next step in ultra high resolution instruments. Quantum effects are already being noticed in instruments.

A general point always made when discussing these instruments is that they are basically planar in concept: they address areal rather than volumetric (3D) problems. This refers back to the semiconductor industry. For surface metrologists this dependence reflects the importance of surface rather than volumetric attributes. In surface instruments with such small $z$ range the probe can be regarded as normal to the measurements in virtually all applications. What will happen when the $z$ axis becomes an equal partner in range is that probe orientation will take on a far more important role than it does at present.

Another point is that the $M^3$ instrumentation is a position/length measuring instrument: the probe measures geometry, making the instrument traceable to international length standards which is good. For instruments measuring any surface property other than geometry the advantages are not so obvious. This aspect will be discussed later on in the conclusions of the chapter.

It is fortunate that another possibility does exist because the wavelength of light is very large when compared with that of the small surface textures and displacements now being requested. The establishment of height standards in the nanometre range has long been needed yet the interpolation of light fringes, although good, is not adequate for the future. What is needed is an atomic unit.

The first attempt to use atomic-state characteristics of materials was made by Whitehouse in the late 1960s and early 1970s [23]. This was an attempt to use the lattice spacing of a crystal as a mechanical step height.

Two materials have been tried, mica and topaz. The idea is that if a crystal is cleaved it will not split indefinitely along one plane. It may jump planes from time to time, thereby leaving steps. Because the lattice spacing is fixed the step would be known. The lattice spacing itself could easily be found by x-ray diffraction methods.

To make such a method practical a number of criteria have to be met:

1. The crystal material has to have a plane which cleaves easily.
2. No other lattice direction should cleave easily with a spacing close to the preferred one and should preferably have no other low-energy lattice spacing.
3. The crystal should be very hard, smooth and flat.

Mica has the property of splitting easily in a laminar way, which is very practical. Unfortunately it tends to flake and is somewhat spongy because more than one layer tends to separate out. Also it is not very hard. Topaz on the other hand cleaves well, is hard and does not give a wavy or soft surface. Unfortunately there are practical problems which is outside the accuracy of 5%. Furthermore it is rare that the crystal fracture jumps over just one lattice spacing. Often it is quite a large number such as 10. The problem here is that the actual number is not known and, if the value of $n$ is large, the possible ambiguities between the two possible accentuations becomes greater. It is possible to unravel the steps, again using the method of exact fractions, but so far this has not been attempted. It may be that more suitable materials can be found, such as rocksalt. Certainly it is a simple, elegant way of achieving a very small step of the same scale of size as that required.

### 5.9.3  X-ray methods

It makes metrological sense to use a unit which is derived from atomic structure to measure down to nanometres. An alternative method, which has great potential, uses the crystal lattice as a moiré grating on the atomic scale. This is a crude approximation to the true physics of the situation, but it suffices [24]. The equivalent of light rays on this scale are x-rays. Basically the stylus is rested on one blade of an interferometer. This is moved and the movement monitored by two independent methods, one the x-ray moiré and the other the transducer of the stylus instrument. In this way the same displacement is monitored by two independent methods. Obviously the method utilizing the x-rays and crystal lattice dimensions is the reference.

The theory is dependent on the availability of very high-quality semiconductor-grade silicon [25]. Crystals of up to about 100 mm diameter, which are free from long-range imperfections such as dislocations, are read-

ily available. They may contain small dislocation loops and local inhomogeneity, but commercial-grade material of 1 part in $10^7$ is acceptable.

The method of x-ray interferometry allows a crystal to be used as what is in effect a spatial encoder. A translation gives a sine wave with a spatial frequency of one lattice plane spacing. The x-ray configuration is shown in figure 5.21.



**Figure 5.21** X-ray path through silicon interferometer.

Three crystal blades A, B and C are employed with their crystal axes very accurately aligned. Two blades A and B are held stationary and the third blade C is translated in a direction perpendicular to the lattice planes used for diffraction. The x-ray intensity of $\alpha$ and $\beta$ will oscillate with the periodicity of the lattice planes. Furthermore, given that the blades are thin enough, waves $\alpha$ and $\beta$ will be in phase and quadrature which means that the measurement is also sensitive to direction.

The point to note is that the wavelength of the x-rays (40 kV or thereabouts) does not influence the periodicity.

The complete mathematical description of the wave intensities is complex despite the simple periodic results. This is because of absorption in the crystal which modifies the phase/quadrature relationship of $\alpha$ and $\beta$.

Simply, the incident monochromatic x-ray beam of wavelength $\lambda$ is diffracted from the crystal planes $d$ at $\theta$ which satisfies the well-known Bragg equation

$$n\lambda = 2d \sin \theta. \tag{5.55}$$

When diffraction is considered by the dynamic theory it is found that three beams emerge from the crystal: the attenuated normal beam T, the diffracted beam R and the forward diffracted beam (figure 5.21).

The incident beam is of no importance except as a source of noise but it is clear that both the other two beams 'know' the spatial phase of the crystal lattice that they have traversed. The second blade B redirects these beams towards each other, and also produces two more beams that have to be washed out so that they overlap at the third crystal blade C. The wavefield created in the third crystal is, therefore, a linear combination of the wavefronts produced by these two coherent beams. This wavefront will have maximum amplitude when the spatial coherence of all three blades is identical. When the third blade is translated with respect to the other two, then the output beams oscillate in intensity with the periodicity of the crystal planes.

Note that if the centre blade is translated and the two other ones held fixed then the oscillation has twice the spatial frequency.

To promote temperature stability and rigidity the whole interferometer is usually made out of a monolithic crystal of silicon. The signal-to-noise ratio is determined by absorption, which is minimized by the following precautions:

1. The crystal blades are kept thin (0.5mm). They have to be as thin as possible yet maintaining maximum strength.
2. (111) crystal planes are used for the diffraction because they give the strongest diffraction and weakest absorption.
3. A silver-target x-ray tube is used for penetrating radiation.

The periodicity of the 111 plane is 0.31 nm.

The problem of optimizing the design of such a system is complicated but possible. There are a number of features of importance:

1. the test bed, to isolate the interferometer from the environment and vibration;
2. the design of the interferometer itself;
3. the mounting;
4. the drives;
5. electronics and control.

A general interferometer is shown in figure 5.22.



**Figure 5.22**

A typical output is shown in figure 5.27. Figure 5.28 shows a comparison of the x-ray interferometer with a laser interferometer. The difference in performance is striking.

The three blade $X$ interferometer is not necessarily useful for displacement measurement. If the third moveable blade is rotated as a roll relative to the x-ray beam moiré fringes can be produced, and if the rotation is a yaw Vernier fringes can result. In fact the presence of moiré fringes can be used in the setting up procedure of the translatory interferometer. Alternatively, the monolith can be designed to use moiré fringes directly as will be seen.

One modification to the basic design is shown in figure 5.23 [26]. This design is to extend the translation range. It has larger blades to keep the working area free of edge effects. More significantly the three blades are placed above the traversing mechanism. Also the blade and sensing point move together more closely. The addition of a reflectance monochromator ahead of the first blade allows only the x-rays matched to the blade set to actually go through. By this means (Bragg reflection) the signal-to-noise ratio is improved.

Ranges of up to 10 $\mu$m are feasible using this technique even under modest metrology room standards but it has to be concluded that for ranges larger than this, the monolith concept has to be dropped in favour of the independent blade configurations, which allow an arbitrary range at the expense of difficult setting up procedures.



**Figure 5.23** The novel long-range monolith proposed by Chetwynd *et al*.

Another modification to the basic scheme is shown in figure 5.24. This enables translation in the *x* and *y* directions rather than just one dimension. The passage of the x-ray beam is shown in Figure 5.25.

In this method x-rays scattered from the (111) plane of the silicon are used for the vertical *y* direction and the (220) rays for the horizontal *x* direction. It is possible to get *x* and *y* movements simultaneously. The biggest problem with such a system is the presence of parasitic variation coupling the two directions when the analysis blade (number 3) is being moved.

The reason why the preferred (111) plane is reserved for the vertical axis is that many instruments for measuring surface texture and form rely to some extent on gravity so it makes sense to use this plane [27].



**Figure 5.24** Isometric sketch of *x*–*y* interferometer monolith. The base dimensions are $40 \times 30 \times 25$ mm.

**Figure 5.25** Schematic diagrams of the $x$–$y$ monolith showing x-ray paths.

One of the major practical problems of the translational x-ray interferometer discussed above is the control of parasitic rotations of the moving blade due to coupling of the drive motion with the movement of the blade. This interaction causes Moiré fringes across the field of the detector which causes loss of contrast.

Recently the availability of low cost x-ray cameras [28] based on CCD arrays opens up the possibility of using the Moiré fringes for submicroradian angle measurement. Although not directly involved as surface metrology, very small angle measurement is an important aspect of nanometrology, which overlaps both surface and dimensional features.

One example is shown in figure 5.26. Part (a) of the picture shows the path of the x-ray through the system and a schematic isometric view. Part (b) shows the angle interferometer. This is a silicon monolith as before but slightly larger; the blades are approximately $30 \times 35$ mm in size.

While angle metrology is the most obvious application of Moiré fringes produced by x-ray devices it is pointed out [28] that it might have further applications.

If the blades of pitch a have a rotation of $\alpha$ the Moiré fringe spacing is $\lambda_m$

$$\text{Where} \quad \lambda_m = \frac{d}{\sin \alpha} \cos(\alpha/2) \sim \frac{d}{\alpha} \tag{5.56}$$

If yaw or strain causes a small change in $d$ (i.e. the Vernier fringes) $\lambda_v$ is given by

$$\lambda_v = \left( \frac{1+\varepsilon}{\varepsilon} \right) . d \sim \frac{d}{\varepsilon} \tag{5.57}$$

The general situation involves relative rotation $\alpha$ and expansion or reduction $\varepsilon$ of pitch. The combination produces fringes at an angle $\gamma$ from the normal from the unstrained grating of $\lambda_c$.

**Figure 5.26** A monolithic x-ray interferometer for angular measurement: (a) schematic illustrations showing typical ray paths, illustrations (b) photograph—the 'blades' are approximately 35 x 30 mm.

$$\lambda_c = \frac{d}{\sin \alpha} \cos(\gamma - \alpha) \text{with} \tan \gamma \simeq \frac{1}{\alpha}\left(t + \frac{\alpha^2}{2}\right) \qquad (5.58)$$

If the fringe pitch is measured along the original direction of the gratings, the projection of $\lambda_c$ into the line of the grating gives

$$\lambda_H = \frac{\lambda_c}{\cos \gamma} \sim \frac{d(1 + \varepsilon)}{\sin \alpha} \qquad (5.59)$$

This gives the relationship between $\lambda_H$ and $\alpha$.

In other words, measurement of the Moiré pitch can give the rotation $\alpha$ directly even in the presence of a small Vernier effect $\varepsilon$.

Use of a monolith as shown in figure 5.28(b) is therefore the basis of an angle standard. A system such as this can readily be used for x-ray.

In the various configurations using x-ray interferometry the apparently simple geometry is misleading. To minimize error of motion parasitic effects small weights have to be added especially to reduce unwanted rotations. Where to put these weights and the values of them is found out to some extent empirically. There are other special manufacturing procedures that allow the system to be optimized and matched to a specific application. At this stage of development building x-ray interferometers is as much art as it is science.

Phase imaging [28] using phase shift rather than absorption of intensity as the physical mechanism causing variation in the interference plane is exactly the same for x-rays as it is for normal optics. It represents the difference between envelope detection and phase detection. Remember that these are the two options in, for example, white light interference described in chapter 4. In the phase mode the system becomes in effect an

x-ray phase microscope. The x-ray Moiré fringe pattern can be used in the same way as an optical interferometer uses tilt fringes to highlight form error on a workpiece.

### 5.9.4   Practical standards

For convenience in workshops, there are portable standards made in much the same way as the stylus calibration standard. There are variants naturally, but one such standard is shown in figure 5.29.

The standard consists of a ruled or etched groove which is calibrated. The material used is often quartz but can be made from glass and chrome plated to make them suitable for use as workshop standards.



**Figure 5.27** Experimental results of x-ray interferometer: (*a*) experimental record of three-phase interferometer; (*b*) deduced displacement plotted versus demand ramp.



**Figure 5.28** Calibration result with overinterpolated HP interferometer.

There may be three 'lands' but one preferred method is to use only the centre height and the adjacent two valley bottoms. To minimize the possible noise values only a small region with each of the lands is used, such as $R_1, R_2, R_3$ in figure 5.29.

The difference between standards used for optics and those used for stylus methods is that the former usually involve the metallizing of the surface to ensure high and uniform reflection from both the upper and lower levels of the step or groove.

Flat-bottomed grooves have the advantage that the calibration of magnification is not affected by the quality of the tip. Moreover, they tend to ensure that the surface quality of the two levels is identical, which is as desirable for the optical method as it is for the stylus.

The important thing to remember about such standards is that the instrument is being calibrated as a system; from front end to recorder. This is vital to maintain the high fidelity needed and is so easily lost when dealing with workshop quantities.



**Figure 5.29** Practical instrument.

Automatic calibration can be achieved by having a fixed movement programmed to cover the regions of $R_1, R_2$, and $R_3$. These height values would have a best-fit line fitted through them to get the calibrated height value.

There are other possibilities for measuring and calibrating height displacement as well as other dimensional features, such as roundness, at one and the same time. One method is to use beads but the technique is at present limited to optical microscope calibration. However, it may be that in the future harder materials will become available and the method extended. At present the arrangement is to manufacture small precise beads made out of styrene monomer and which swell around hydrocarbon seeds in water [29]. Spheres of $9.89 \pm 0.04 \mu$m have been made in space and represent another possible way of calibrating small sizes. The basic problem is that the spheres are soft and only suited to viewing. Fortunately for surface metrologists it is easier to make small balls than large ones, so 1 $\mu$m calibrated spheres may be available soon. Large quartz and glass balls of about 20 mm are now being used to calibrate the amplitude of roughness and form.

The use of space frames such as a tetrahedron ball frame or a rack of balls with cylinder ends [30] to calibrate the volumetric accuracy of coordinate-measuring machines will not be considered here. Neither will

the use of interferometers directly [31,32]. These are well known and not especially relevant to surface metrology calibration.

### 5.9.5 Calibration of transmission characteristics

This involves the transmission characteristics from the input, that is the stylus or probe, through to where the output is registered. Basically what is required is how the instrument sees waves of a given length and communicates them, whole or in part, to the storage device or to be acted on as part of an adaptive control loop in a bigger system.

There are a number of ways of approaching this, but some involve the direct use of time and others the indirect use. Consider first the direct use. This was originally needed because the filters, transducers, etc, are analogue, that is temporal in response. This situation is changing rapidly because now the filtering and metering are carried out digitally and so have to be dealt with differently.

However, the basic idea is to evaluate the transfer function of the system. This concept assumes a linear system. Unfortunately there are occasions when the system is anything but linear. One such instance is when the stylus bounces off the workpiece because the traversing speed is too high. Another is because the stylus itself acts as a non-linear mechanical filter if it does not penetrate all the detail. Having said that it is still possible and necessary to have a clear idea of the metrology system response, not only for calibration purposes but also to enable it to be incorporated into the larger manufacturing process loop.

There are two basic ways of checking the transmission characteristic. Both involve oscillating the stylus with a constant amplitude and waveform through a sufficient range of frequencies to check all those wavelengths which might be important on the surface. One way is to use a vibrating platform, and the other is to use a periodic calibration surface. The former will be considered first.

The basic set-up is shown in figures 5.30 and 5.31.



**Figure 5.30** Vibrating-table method of instrument calibration.

A vibrating platform energized from a low-frequency oscillator has been used. The basic problem is not that of monitoring frequency or even waveform, but that of keeping the amplitude of vibration constant over the range of frequencies. Notice that this type of calibration has to take place with the skid removed or not touching the vibrating table, otherwise no signal is seen in the amplifier.

A number of ways have been devised to monitor the amplitude of vibration. Originally a small slit was put on the table and this was illuminated by a light source. The slit or flag impedes the amount of light falling onto a photodetector. Depending on the output the voltage into the drive circuit of the vibrator is adjusted to keep the movement amplitude constant. This turns out to be less than satisfactory.

**Figure 5.31** VNIIM calibrator using interferometer.



**Figure 5.32** Step and groove measurement.

The preferred way to monitor the displacement is to use an interferometer [31,32] similar to figures 5.31 and 5.32.

In this the armature of a moving coil is used for two purposes. The upper face is used for the table and the lower face for the mirror of the interferometer. There are endless variations on this theme but one is that the reference mirrors are slightly tilted so that fringes move across the face of the detector. Alternatively the

detector could be a simple counter. Various ways have been used to enhance the sensitivity of such a device when the movement has been required to be less than $\lambda/2$, for example the use of Bessel function zeros [32] when more complex optics are used. It seems that this technique will ultimately be used for the overall calibration of instruments for both amplitude and frequency. Laser interferometers have been used to measure both step and groove heights in this system [33].

Despite being a method which is easily controllable, the vibrating-table method has two serious disadvantages. The first is that it is difficult to transport and therefore it cannot be used as a workshop method. The second is that it does not truly measure the system characteristics because the pick-up is not traversing horizontally when the frequency response is being checked. The problem of dynamic calibration of SPM is the same as for surface texture instruments only on a smaller scale. Very little has been reported. The following technique addresses these two problems.

### 5.9.6 Filter calibration standards

The gradual emergence of substantially sinusoidal instrument calibration standards has been leading towards a simpler method. The idea of a portable standard is attractive and has been used to calibrate the meter for years. The question is whether a similar type of standard can be used to check the transmission characteristics.

Basically a specimen of suitable spacing of say, 250 $\mu$m and 2.5 $\mu$m $R_a$ is traversed preferably with a skidless pick-up at various angles giving effective spacings from 250 to 2500 $\mu$m (i.e. about 10:1 range). The output, which has only to be expressed as a percentage of the maximum indications, should fall for each transmission value within the specified range of the standard.

The method can be used with skid-type pick-ups providing that the skid has a long enough radius both along and across the direction of traverse to make the skid error negligible (<1%). The requirement for the specimen above would be a spherical skid of 10 mm radius.

Concerning the shape of the waveform of the specimen, it needs to be said here that there are differences in the value of $R_a$ if the specimen is non-sinusoidal. Various other shapes of cross-section have been tried out as possibilities (figure 5.33) for the simple reason that sine wave standards are not easy to make. Triangular and cusp-like shapes have been made [33, 34] and the various values of $R_a$ evaluated for different angles of traverse. In general the curves for all shapes are very similar except at places where the attenuation is meant to be high, for example at long wavelengths. Care has to be exercised in this region. Also, note how difficult it is to make a suitable standard: it has to have an integral number of cycles *and* have practical non-wearing ends (figure 5.33).



**Firgure 5.33** Different standard waveforms for filter characterization: (*a*) sine wave standard; (*b*) sawtooth standard; (*c*) cusp standard.

Great care has to be taken when using oblique tracking to make sure that the actual assessment length includes an integral number of cycles of the waveform, otherwise the answer can vary from place to place. If there is any doubt in this respect it is sometimes wise to use a peak parameter as a measurement of the output rather than the $R_a$ because peaks are not affected by end effects.

Specimens have been made by etching, photographic methods [34] and nowadays using diamond turning [35]. This last method, whilst being more or less satisfactory (in the sense that the harmonic content over and above the sine wave is less than 1 %), has the problem that the roughness mark is slightly arcuate. This is caused by the method of generation and so obliquity has to be considered carefully. The problem of using filter calibration specimens is how to mark the angles adequately for tracking on the specimen so that no error occurs when it is being used. When the tracking occurs, there is still the problem of actually calibrating the whole instrument through to the meter or display. Ideally one standard should be used to calibrate both the filter and the amplifier. This latter job is conventionally achieved by means of meter calibration standards. The word meter is still used although it really refers to the output value in whatever form. These standards take many forms and have changed over the years.

In its earliest embodiment meter calibration meant just that the procedure was to record the profile of a suitable surface with accurately calibrated recording means and to note the meter reading; then to evaluate that portion of the profile to which the meter reading referred and to adjust the meter to it. At one time graphic records were evaluated by dividing the graph into successive sections of equal length—in fact the sample length described in BS 1134—and drawing a centre line through each section and determining the $R_a$ value from an area measurement with a planimeter or by counting squares (see section 3.9.1 on the use of planimeters).

Since the advent of computers the output is often the processed signal and not a meter reading.

The types of meter calibration specimen differ from each other primarily in the degree of uniformity of texture which is secured within an indicated area of the surface and hence in the possible accuracy of calibration. They can be considered together because the basic instrument approach is the same.

Strictly speaking, today it could reasonably be argued that if it were possible to traverse the stylus or optical probe system infinitely slowly across the surface it would not be necessary to use meter calibration standards at all; the parameter fidelity would depend only on software. However, dynamics cannot be ruled out and it is often beneficial to present to the instrument a known and sometimes typical surface. Whereas this might be relatively straightforward for surfaces simulating single-point cutting it is not so for random surfaces. Two methods have been used for this. The simplest is to use a random surface itself as the calibration specimen. This would automatically have the necessary mix of different wavelengths and amplitudes to test the instrument. The problem is that it is very difficult to arrange a random signal and at the same time to ensure that the variability across the standard is very small. The only satisfactory method of achieving this so far has been by using either plunge grinding or creep feed grinding. This has to be done with extreme care. Ideally the width of the wheel should be equal to the assessment length of the instrument being checked. It must be so because the uniformity should be within 2% so that instruments can be calibrated without paying too much attention to the initial position of the stylus on the specimen [34]. Hillman at the PTB has been a leader in this field.

The nanometre range to traditional range—calibration standards have been the object of much recent investigation. One big problem is the range of surface roughness, another is the range of instruments. A concentrated project to investigate possibilities has been carried out between German, Danish and UK investigators [36]. Key to the project has been the use of injection moulding of plastic negatives with different step height, sinusoidal, triangular and arcuate profiles covering $R_t$ values from 50 nm to 100 $\mu$m and $S_m$ values from 0.8 $\mu$m to 800 $\mu$m. Primary standards were produced by ion beam and plasma etching (step height standards), by holographic generation of sinusoidal structures with two beam interference exposure and finally by diamond turning for triangular and arcuate standards. From these primary standards plastic negatives and nickel negatives were made from which positive nickel standards were produced. The idea behind

the project was to completely control the manufacture including cleanliness sealing and packaging aspects of the industrial class standards based on the aforementioned primary standards.

The parameters $R_t$, $R_a$, $R_q$, $R_{sk}$, $R_{ku}$, $R_{\Delta q}$ and $R_{sm}$ were those selected for comparison. The report uses the designation $P$ rather than $R$ indicating no filtering and the materials used were steel, silicon, glass, nickel and PVC with other plastics.

So consistent were the replicas from the primaries that the need for separate calibration of replicas could be questioned in the future. Also a mini disc cage with a moveable window was developed to ensure that only one stylus track could be made on any single part of the standard. The uniformity of the replica and its potential cheapness points the way to throw away standards! This intriguing philosophy assumes that the variability across the standard is much smaller than the errors that might result if the standard was used often in one place. This means that the lifetime of the secondary standard need no longer be an issue.

This fine piece of work needs now to be backed up with measurements using different instruments. It is feared that the diamond turning option and the use of plastics perhaps uncoated, may cause optical problems.

An alternative method devised by Thwaites [37], as with the wedge method versus the lever method in figures 5.15 and 5.16, shows the possibility of a more accurate, yet repeatable, method. In this he simulates the pseudo-random signal on a generator, which in turn drives the vibrating table to give the impression of a random surface. In this case the randomness is exactly repeatable, unlike that on an actual specimen, but it does have the considerable disadvantage that it is not as portable or convenient as a tactile standard—yet it is more repeatable. Also the effects of horizontal movement of the stylus cannot be simulated.

It may seem that it is unnecessary to go to the trouble to generate such random surfaces for parameter evaluation using digital methods. The argument is that providing that the sampling rate is within the Nyquist criterion, all parameters could be adequately measured. This is true, yet the validity of algorithms often depends on the differences between measured ordinates and the rate of change of differences. These can be very small for any periodic surfaces, especially sine waves and the like, but not for random surfaces. Algorithms which are well behaved for smoothly varying surfaces can break down for jerky surfaces. For this reason a random input check is very useful, especially as there is no universal agreement between manufacturers on methods of testing filters or algorithms.

## 5.10 Calibration of form instruments

This section refers to the calibration of roundness instruments and instruments measuring form. The problem here is at the same time more difficult in some respects and simpler in others. First there is the determination of the magnification and second the generation in space (or by an object) of a perfect circle or form as a reference from which to measure.

### 5.10.1 Magnitude

Taking first the problem of magnification, it has to be some method that can be used on the instrument if at all possible and not require the pick-up to be removed. The step equivalent in roundness is not readily available because the tracer may be rotating and not traversing a linear track.

Conventionally, the method for establishing the magnification has been to use a flat plane which has been lapped on a cylinder. This produces a blip on the recorder which can or should be able to be used as a reference (figure 5.34) [38]. This certainly produces the equivalent of a height change but unfortunately it is not the same as a step used in surface roughness calibration—for the reason that the change in height is short-lived and abrupt. This means that very high frequencies are introduced through to the recorder.

**Figure 5.34** Calibration of roundness instrument: (*a*) cylinder with flat; (*b*) magnified view; (*c*) chart representation.

Another alternative is to use a tilted cylinder (figure 5.35) which does not have the same problem as the milled flat because it is well behaved. Assuming that the stylus tip is infinitely sharp, the height variation in traversing a tilted cylinder is

$$2\rho - 2r = 2r(\sec \alpha - 1). \tag{5.60}$$



**Figure 5.35** Use of tilted cylinder to calibrate amplitude of roundness.

This technique is quite sound but suffers from two fundamental problems. The most important of these is that the cylinder has to be tilted by an exact amount. Unfortunately, although it is possible to mount a cylinder at a very precise angle in its mount, it is not possible to guarantee that the mount sits properly on the table of the instrument; small dirt particles or burrs on the table or the mount can throw out the angle by a small amount.

The finite-size stylus accentuates the problem of calibration. It will normally tend to increase the apparent height variation between the major and the minor axis, so for calibration using this method it is advisable to use a fine diamond stylus and also to mount the cylinder on a wide base—an optical flat base— so that dust or debris does not change the tilt of the cylinder from the calibrated amount. Reason, in his book on roundness [38], sets limits for the tilt and the size of the stylus that can be allowed under practical conditions.

### 5.10.1.1  *Magnitude of diametral change*

The change in the apparent diameter of the part in the direction of tilt will depend on the nature and radius of the part, the angle of tilt and the radius of the stylus in the axial plane (also on the extent of the radius, but this will be neglected here).

(*a*) *Shafts and holes engaged by a sharp stylus*
The change in the major diameter caused through a tilt of $\theta$ is to a first approximation $\Delta d$, where

$$\Delta d = r\theta^2 \tag{5.61}$$

where $r$ is the radius of the workpiece.

(*b*) *Shaft engaged by a hatchet*
When the shaft is tilted by an angle $\theta$ and the radius of the hatchet or ball is $s$

$$\Delta d = (r + s)\theta^2 \tag{5.62}$$

to the same degree of approximation.

(*c*) *Hole engaged by a hatchet*
This is surprisingly different because the hatchet (or stylus) radius works against that of the size of the part:

$$\Delta d = (r - s)\theta^2. \tag{5.63}$$

In the special case when $s=r$ the measurement will be insensitive to tilt of the part. Obviously this could not be used as a basis for a calibration standard for roundness. However, it does illustrate that careful instrumentation can reduce variability.

From the formula (5.63) the sensitivity of the change in diameter for constant shaft or cylinder size is

$$\delta d = \frac{\partial d}{\partial s}\delta s + \frac{\partial d}{\partial \theta}\delta \theta \tag{5.64}$$
$$\delta d = \theta^2 \delta s + 2\theta\, \delta\theta$$

so that the change in angle $\theta$ caused by debris on the mount can be more important than the radius of the stylus.

### 5.10.2 *Separation of errors — calibration of roundness and form*

A classical but nevertheless empirical rule of metrology is that the measuring instrument should be an order of magnitude (i.e. ten times) more accurate than the deviations expected in the test piece. Failure to achieve this degree of separation is commonplace. So much, in fact, that metrologists have devised a large number of ways of achieving the desired accuracy indirectly. Such methods include self calibration, reversal methods, redundancy and error separation techniques. One prerequisite is that the measuring instrument should be reasonably stable. A review of techniques has been compiled by Evans [39]. Most of these techniques are covered in this book under the heading of error separation. The theoretical basis is usually included.

There are two issues to be resolved. One is the calibration of the reference of the instrument, the spindle bearing (i.e. the circle generated in space by the instrument), and the other is the calibration of the specimen used to check the spindle.

Fortunately these two problems can be solved at the same time using a method described in the earlier section 2.3.6.13 on roundness and form (figure 5.36). This method presupposes that the spindle is repeatable and that random components have a bandwidth that lies outside the shape deformation of the specimen and the spindle.

**Figure 5.36** Calibration of specimen and spindle errors.

A simple way is as follows [40]. Let the spindle error be $e$ and the signal from the specimen $s$. A probe will detect a voltage $V_1(\theta)$

$$V_1(\theta) = s(\theta) + e(\theta). \tag{5.65}$$

Turning the part through half a revolution, moving the probe on its carriage through the axis of rotation without touching the spindle and then changing the direction of sensitivity of the probe gives a voltage $V_2(\theta)$ where

$$V_2(\theta) = s(\theta) - e(\theta). \tag{5.66}$$

Simple arithmetic gives

$$e(\theta) = (V_1 - V_2)/2 \quad s(\theta) = (V_1 + V_2)/2 \tag{5.67}$$

The idea is that with the two measurements the specimen errors and spindle errors have to be arranged to add in the one case and subtract in the other.

This method, although simple, does have its problems. Perhaps the most important of these is that the pick-up needs to be handled when it is passed through the centre and reversed in sensitivity. Also it may be necessary to stop the spindle rotation during the evaluation. These considerations could be important in highly critical cases. For example, the spindle bearing could be hydrodynamic, in which case the oil film separating the bearings would be disturbed when stationary and the true errors would be distorted or become effectively non-stationary.

To get the separate errors it is necessary to plot the functions represented by equation (5.60) on graph paper or to store them in a computer.

Once $e(\theta)$ has been determined for all values of $\theta$ from 0° to 360°, then the error in the spindle has been determined. Similarly for $s(\theta)$, the specimen systematic error, but the reference of angle has to be preserved, that is the position corresponding to $\theta = 0$ has to be marked on the specimen. It is important to note that this $s(\theta)$ represents the systematic error, and the hope is that this is true over the very long term; recalibration should have a timescale of months or years and not days.

Normally this requirement is met with precision instruments which tend to have stable datums. For instance, the spindle on a roundness instrument would probably have been aged for months before installation to reduce the possibility of distortion due to stress relief when the instrument is in use.

If this is true then it becomes possible to redefine the meaning of uncertainty of form measurement. Instead of relying on the absolute accuracy as being the limit for the measurement, it becomes a precision or repeatability of the system that is the limiting factor. This means that the effective uncertainty of measurement of the instrument can be reduced by a factor of about 5:1. This is because the repeatability is about this factor less than the systematic error.

This is a simplistic approach to errors because the nature of the systematic error could well depend on the type of instrument. As an example in roundness, if the instrument is of the rotating-table type referred to earlier in section 2.3.6.3 the systematic error is a questionable quantity because it can depend on the shape and weight of the part being measured. If the centre of gravity of the part is not coaxial to the bearing axis this can distort the bearing thereby reducing accuracy. Thus, for this type of instrument the accuracy can change with the workpiece. For this reason the rotating-table type is not used for the highest accuracy and also because the systematic error may be somewhat changeable. Guarantees should be obtained on tilt stiffness of the instrument before this method is used. Rotating-spindle methods are free from this problem.

At first sight there would appear to be two alternative approaches to this method. Take the possibility of using two probes simultaneously at $180°$ to each other (figure 5.37(a)) or of using two orientations of the workpiece relative to the datum (figure 5.37(b)).



**Figure 5.37** Ways of separating specimen and spindle error: (a) two probes; (b) two orientations.

In figure 5.37(a), if probe 1 gives $V_1$ and probe 2 gives $V_2$ then

$$V_1(\theta) = s(\theta) + e(\theta)$$
$$V_2(\theta) = s(\theta - \pi) - e(\theta).$$

(5.68)

What effect this has can be seen by using Fourier analysis techniques. Fourier methods are used extensively to describe the behaviour of roundness instruments because the signal is naturally of a periodic nature.

Thus if the signal is analysed at discrete harmonics $n = 0, 1, 2$, where $n = 1$ corresponds to one undulation per circumference of the part, and the results of the analysis for $V_1, V_2, s$ and $e$ are $F_1, F_2, F_s$, and $F_e$ respectively, then

$$F_1(n) = F_s(n) + F_e(n)$$
$$F_2(n) = F_s(n)\exp(jn\pi) - F_e(n).$$

(5.69)

From this $F_s(n)$, the transform of $c = V_1 + V_2$, the combination which eliminates $e$, is found by simple arithmetic to be

$$F_c(n) = F_s(n)[1 + \exp(j\pi n)].$$

(5.70)

The term $1 + \exp(jn\pi)$ is a harmonic distortion factor and represents the error in the harmonic value in amplitude and phase that would occur if $F_s(n)$, the true harmonic content of the surface, were approximated by $F_c(n)$.

Examination of equation (5.70) shows that $1 + \exp(jn\pi)$ has a magnitude of $2\cos(n\pi/2)$. From this it is obvious that if $n$ is odd no harmonic value is measurable, but if $n$ is even then it is. The well-known geometrical implication of this statement is that only diametric information is retrievable using two probes in this way (as is the case with a micrometer).

The second possibility uses two orientations of the specimen shown in figure 5.37(b). Here one probe is used instead of two. The orientations are shown at 180°.

The two outputs from the probe will be $V_1$ and $V_2$ where

$$V_1(\theta) = s(\theta) + e(\theta)$$
$$V_2(\theta) = s(\theta - \pi) + e(\theta)$$

(5.71)

and their associated transforms are

$$F_1(n) = F_s(n) + F_e(n)$$
$$F_2(n) = F_s(n)\exp(jn\pi) + F_e(n).$$

(5.72)

The combination signal $c = V_1 - V_2$ required to eliminate $e(\theta)$ gives

$$F_c(n) = F_s(n)[1 - \exp(jn\pi)].$$

(5.73)

This time the magnitude of $1 - \exp(jn\pi)$, the harmonic distortion, is given by $2\sin(n\pi/2)$ which is zero for all values where $n$ is even—exactly the opposite of the two-probe system. Obviously combining the two gives the required result.

Despite the fact that both of these alternatives to Donaldson's method [40] are mechanically simpler they give less harmonic information. In fact combining equations (5.70) and (5.73) gives the Donaldson result because the resultant is independent of $n$. This is equivalent to saying that the Donaldson result could be achieved by using two probes and two orientations.

Note that, in the formulation given in equations (5.70) and (5.73), no spindle error information is obtained. By reversing the procedure of developing $c$ the specimen error can be removed and the error signal of the bearing remains.

Thus the equivalents of (5.70) and (5.73) become

$$F_c(n) = F_e(n)[1 + \exp(jn\pi)]$$
$$F_c(n) = F_e(n)[1 - \exp(jn\pi)]$$

(5.74)

which are completely symmetrical with the other equations.

Remember that, in the two alternatives to Donaldson's method, the combination signal $c$ does not contain spindle (or table) error. It does contain somewhat scrambled information of the component out-of-roundness. For those occasions where either may be used in preference to Donaldson's method it is necessary to show how the useful component signal can be extracted. To do this it is essential to understand the nature of the harmonic losses. Consider the two-orientation method. Let the orientations be separated by an angle $\alpha$. The Fourier transform of the combination signal becomes $F_c(n)$, where

$$F_c(n) = F_s(n)[1 - \exp(jn\alpha)]. \tag{5.75}$$

This will have zeros at $\alpha n = 2\pi N$ where $N$ is an integer. The range of frequencies that can be obtained without having zeros will extend from the first harmonic to $2\pi/\alpha$. Hence, measuring the Fourier coefficients of the combination signal up to $2\pi/\alpha$, modifying them according to the weighting function contained in equation (5.75) and then synthesizing them should produce the out-of-roundness of the component, subject only to a high-frequency cut at $2\pi/\alpha$. An instrument with an in-built computer can do this. If there is any danger of higher-frequency components causing aliasing the combination signal can be attenuated by using a high-cut analogue filter prior to analysis. Thus $F_s(n)$, the true coefficient of the component out-of-roundness signal, can be obtained from equation (5.75) using an equalization technique:

$$F_s(n) = F_c(n)/[1 - \exp(-jn\alpha)]. \tag{5.76}$$

The amplitude modification necessary, $A(n)$, is given by

$$A(n) = [(1 - \cos n\alpha)^2 + (\sin n\alpha)^2]^{1/2} \tag{5.77}$$

and the phase by

$$\Phi(n) = \tan^{-1}[\sin(n\alpha)/(1 - \cos n\alpha)]. \tag{5.78}$$

These equations can be conveniently simplified to $A(n) = 2\sin(n\alpha/2)$ and $\Phi(n) = (\pi - n\alpha)/2$. To get the correct values of $F_s(n)$ the amplitudes of $F_c(n)$ are divided by (5.77) and the phases shifted in the opposite direction to the phase angles given in equation (5.78).

The number of zeros in the weighting function can be reduced by increasing the number of orientations. Take the case of three orientations, for example, at angles $\alpha$ and $\beta$ to the original one. Then

$$\begin{aligned} V_1(\theta) &= s(\theta) + e(\theta) \\ V_2(\theta) &= s(\theta + \alpha) + e(\theta) \\ V_3(\theta) &= s(\theta - \beta) + e(\theta) \end{aligned} \tag{5.79}$$

The combination signal to remove $e(\theta)$ is $c = 2V_1 - V_2 - V_3$ which has a transform

$$F_c(n) = F_s(n)[2 - \exp(-jn\alpha) - \exp(-jn\beta)]. \tag{5.80}$$

Equation (5.80) will only have zeros when both $n\alpha$ and $n\beta$ are together a multiple of $2\pi$. There is a further degree of freedom that can be used to extend the coverage of harmonics and this is the use of different probe sensitivities for different orientations. For instance, taking the first orientation as unit gain, the second as $a$ and the third as $b$, the transfer function corresponding to equation (5.80) is

$$F_c(n) = F_s(n)[1 - a\exp(jn\alpha) - b\exp(-jn\beta)] \tag{5.81}$$

which has amplitude and phase characteristics of $A(n)$ and $\Phi(n)$:

$$A(n) = [(1 - a \cos n\alpha - b \cos n\beta)^2 + (b \sin n\beta - a \sin n\alpha)^2]^{1/2}$$
$$\Phi(n) = \tan^{-1}[(b \sin n\beta - a \sin n\alpha)/(1 - a \cos n\alpha - v \cos n\beta)].$$

(5.82)

The real significance of the use of variable sensitivity in the method will become clear in the case of variable error suppression.

So far, using the multiple orientation technique, the out-of-roundness has been obtained by a synthesis of modified Fourier components. There are other ways. One such simple but novel method is to solve a set of linear simultaneous equations. In effect what needs to be done in the two-orientation method, for example, is to detect that part of the signal which has moved by the angle $\alpha$. The part which moves is identified as component out-of-roundness. The part which remains stationary is attributed to instrument error. So, if the discrete voltages from the probe on the first orientation are $V_{11}, V_{12}$, etc, and those on the second are $V_{21}, V_{22}$, etc, the following matrix equation has to be solved (as it stands this is not straightforward; unless suitable partitioning is used it can give unstable numerical results):

Letting the output at spindle position $\theta$ and specimen position $\alpha$ be $_\alpha V_\theta$ the measurement array is given by equation (5.83).

Readings $\alpha\,\gamma\,\theta$        $\theta$ is position (spindle)  
$\alpha$ is position (specimen)

|  | $\theta_1$ | $\theta_2$ | $\theta_3$ | $\theta_4$ | $\theta_5$ |  |
|---|---|---|---|---|---|---|
| $\alpha_1$ | $_1\gamma_1$ | $_1\gamma_2$ | $_1\gamma_3$ | $_1\gamma_4$ | $_1\gamma_5$ | $\alpha_1\ \bar\gamma$ |
| $\alpha_2$ | $_2\gamma_1$ | $_2\gamma_2$ | $_2\gamma_3$ | $_2\gamma_4$ | $_2\gamma_5$ | $\alpha_2\ \bar\gamma$ |
| $\alpha_3$ | $_3\gamma_1$ | $_3\gamma_2$ | $_3\gamma_3$ | $_3\gamma_4$ | $_3\gamma_5$ | $\alpha_3\ \bar\gamma$ |
| $\alpha_4$ | $_4\gamma_1$ | $_4\gamma_2$ | $_4\gamma_3$ | $_4\gamma_4$ | $_4\gamma_5$ | $\alpha_4\ \bar\gamma$ |
| $\alpha_5$ | $_5\gamma_1$ | $_5\gamma_2$ | $_5\gamma_3$ | $_5\gamma_4$ | $_5\gamma_5$ | $\alpha_5\ \bar\gamma$ |
|  | $\bar\gamma_{\theta_1}$ | $\bar\gamma_{\theta_2}$ | $\bar\gamma_{\theta_3}$ | $\bar\gamma_{\theta_4}$ | $\bar\gamma_{\theta_5}$ |  |

Steps specimen positions (left label)

Specimen errors (right brace label)     (5.83)

Spindle errors (bottom brace label)

Multipoint methods have been described in chapters 4 and 2 of this book for roundness and cylindricity measurement based upon the constraint table (2.14). There have also been other variants e.g. for the measurement of spindle error in an in-process possibility at the same time measuring roundness.

For simplicity only five points per orientation have been shown. Solving for the spindle and component values (here called S) in terms of the matrix M and the input voltages V

$$\mathsf{S} = \mathsf{M}^{-1}\mathsf{V}.$$

(5.84)

This method still suffers from exactly the same frequency suppressions as the synthesis technique. As before, the effect can be reduced by making $\alpha$ small but other problems then arise: differences between

measurements become small—the readings become correlated and the matrix inversion becomes susceptible to numerical noise. For any given $\alpha$, however, it is possible to remove the need for a matrix inversion and at the same time improve the signal-to-noise ratio. This is accomplished by repeating the shift of the specimen until a full 360° has been completed, that is having $m$ separate but equi-angled orientations [41, 42]. The reduction of noise will be about $m^{-1/2}$ in RMS terms. Once this exercise has been carried out it is possible to isolate the component error from the instrument error simply by sorting the information. For example, to find the component signal it is necessary to pick one angle in the instrument reference plane and then to identify the changes in probe voltage at this angle for all the different orientations in sequence. To get instrument errors a fixed angle on the specimen has to be chosen instead. Before this sifting is carried out the data sets from each orientation have to be normalized. This means that the data has to be adjusted so that the eccentricity and radius are always the same. These are the two Fourier coefficients which cannot be guaranteed to be the same from one orientation to the next because they correspond to setting-up errors and do not relate to instrument datum or component errors. Figure 5.38 shows a typical result in which a magnification of one million has been obtained using this method. The plot illustrates the systematic error in a typical spindle. Providing that these errors do not change in time they can be stored and offset against any subsequent data runs, therefore enabling very high magnifications to be obtained. Preliminary results [41] suggest that out-of-roundness of workpieces down to approximately $10^{-6}$ mm may now be measurable. The use of multiple probes has also been considered [43].



**Figure 5.38** Specimen error found by multiple-step method.

Similar techniques can be used for measuring the systematic errors of straightness of a slideway. The equivalent to Donaldson's method is well known. The part is measured normally, it is then turned through 180° and the transducer adjusted to keep contact with the same track on the part. This mechanical trick changes the phase of $e(x)$ with respect to $s(x)$ as in roundness [42, 41]. Alternatively, to get the equivalent of a rotational shift of $\alpha$ in roundness, the part has to be moved linearly parallel to the datum axis [41].

   Obviously there are differences between the roundness and the straightness cases because in the latter case the end signal does not join up with the origin. This means that end information is lost. A more realistic method would be to move the probe to another position on its carriage rather than to shift the workpiece bodily. This has the effect of presenting a different $s(x)$ to the same $e(x)$ (figure 5.39). It does not matter that some end information is lost. All this does is to preclude any absolute distances of the part to the datum from being determined. With the two-position method, changes in the distance between the datum and part can be obtained, from which the out-of-straightness can be found [41]. This method is quite general [44].

   Other calibration requirements are sometimes found in surface metrology instruments. One of these is that of radius calibration. With the recent surface roughness instruments the pick-up has a wide dynamic

**Figure 5.39** Systematic error removal in straightness.

range. This allows the roughness to be measured on curved parts. An added benefit is that the radius of the curve can also be measured at the same time providing that the instrument has been calibrated for it. This can be readily achieved by tracking across the top (pole) of a smooth hemisphere of known radius. Curvature values relative to this known one can then be computed.

### 5.10.3  General errors due to motion

The methods described above are often sufficient to determine the significant errors of a form instrument in the sensitive direction, using Zhang's terminology (reference [31] in chapter 4), whether in roundness, straightness or whatever. However, there are important components of the error which have not been addressed so far, but which are very important in machine tools and coordinate-measuring machines and will to a larger extent be important in the form instruments of the future. These refer to the error motion of spindles in the general case (see e.g. [45, 46]). These errors can appear to be geometric deviations of the workpiece.

It is convenient to deal with the relative motion of the gauge and the workpiece in terms of two line segments. Of these the axis of rotation is embedded in the specimen and moves with it. The other is fixed with respect to the probe at the average position of the axis of rotation and is referred to in the USA as the axis average line (figure 5.40). Equations in this section will as a rule use time as the argument to indicate the dependence on movement rather than $\theta$, for example.



**Figure 5.40** AB, average axis; CD, axis of rotation.

In general, the workpiece has six degrees of freedom consisting of the three angular rotations and three linear translations, as shown in figure 5.41, which appear in the *xz* projected plane.

However, this orientation depends on the fact that the probe is fixed in this position relative to the workpiece. Every configuration of workpiece to probe has to be considered on its merits. For example, the situation is more difficult to visualize when the probe is moving in both rotation and translation, as it may well be in the rotating-probe-spindle method described earlier for cylindricity than it does in the case shown in figure 5.42. Here this case corresponds to the rotating-table method. Again, obviously movement in the *z* direction would not be too serious if only roundness in one plane were being measured, but it would if sphericity were being measured. However, general terms have been applied to those errors usually considered sensitive. These are:

1. Pure radial motion. This is $x(t)$ in figure 5.42 in which the axis of rotation remains parallel to the 'axis average line' and moves perpendicular to it in a sensitive direction and which would obviously show up as a roundness error independent of the *z* position.
2. Axial motion $z(t)$ in figure 5.42 in which the axis of rotation remains coaxial with the average axis line and moves axially relative to it.
3. Angular motion $\alpha(t)$ in which the axis of rotation moves angularly with respect to the axis average line and in the plane of the axial and pure radial motions.



**Figure 5.41** Relative motion and the six degrees of freedom.

*5.10.3.1  Radial motion*

In general, angular motion and pure radial motion occur at the same time, and the sum at any particular axial position is referred to as radial motion. A knowledge of radial motion $r_o(t)$ at one axial position and angular motion $\alpha(t)$ allows the radial motion $r(t)$ at another axial position to be predicted as shown in figure 5.42 (*a*):

**Figure 5.42** Geometry of radial and face motion: (*a*) radial motion variation with axial distance; (*b*) face motion variation with radius.

$$r(t) = r_0(t) L\alpha(t) \tag{5.85}$$

where $L$ is the distance between the two axial locations. Since radial motion varies with axial position, it is necessary to specify the axial location of a radial motion measurement.

### 5.10.3.2 Face motion

Another special term is face motion, which denotes error motion in the axial direction at a specified distance $R$ from the axis average line, as shown in figure 5.42(*b*). Face motion $f(t)$ is related to axial and angular motion by

$$f(t) = z(t) - R\alpha(t). \tag{5.86}$$

Since face motion varies with radial position, it is necessary to specify the radius of a face motion measurement.

### 5.10.3.3 Error motion—general case

The most general case of error motion involves an arbitrary angle $\varphi$ of the sensitive direction with respect to the axis average line, as for the spherical surface shown in figure 5.43. The error motion depends on both the axial and radial locations, which must be specified together with $\varphi$. The equation for error motion $e(t)$ in terms of axial, radial and angular motion is

$$
\begin{aligned}
e(t) &= r(t)\sin\varphi + f(t)\cos\varphi \\
&= r_0(t)\sin\varphi + z(t)\cos\varphi + \alpha(t)(L\sin\varphi - R\cos\varphi)
\end{aligned}
\tag{5.87}
$$

It can be seen from equations (5.85), (5.86) and (5.87) that error motion in general or any of the special cases can be obtained from a knowledge of axial motion $z(t)$, angular motion $\alpha(t)$ and radial motion $r_0(t)$ at a

**Figure 5.43** General case of error motion of accurate spindle.

known axial position. Figures 5.44 and 5.45 show schematic diagrams of two test arrangements which can be used to measure the necessary motions. In both cases, the radial and axial motions are measured directly. In figure 5.44 angular motion is derived from face motion by use of equation (5.88):

$$\alpha(t) = \frac{1}{R}[f(t) + z(t)]. \tag{5.88}$$



**Figure 5.44** Basic motions using radial, face and axial motion measurements.



**Figure 5.45** Basic motions using two radial and one axial motion measurements.

In figure 5.45 a second radial motion measurement is used to obtain angular motion from equation (5.89):

$$\alpha(t) = \frac{1}{L}[r_2(t) - r_1(t)]. \tag{5.89}$$

It should be noted that pure radial motion does not appear in any of the above equations. It is useful as a concept in understanding error motion geometry, but it is not a factor which must be measured in determining the behaviour of an axis of rotation. See also [47].

### 5.10.3.4  *Fundamental and residual error motion*

The term which will be used to refer to the once-per-revolution sinusoidal component of an error motion polar plot is fundamental error motion. Since a test ball is perfectly centred when this component vanishes, it follows that fundamental radial motion of an axis of rotation does not exist.

Similarly, fundamental angular motion does not exist. This can be understood by visualizing a perfect cylinder mounted on an imperfect axis of rotation. If the mounting is adjusted so that the cylinder has no centring error at either end, then there can be no once-per-revolution angular motion. Since familiar terms such as 'coning', 'wobble' and 'swash' suggest a once-per-revolution component, they are inappropriate names for angular motion.

In contrast, fundamental axial motion does exist, and is not caused by a master ball mounting error, as with centring error. It consists of a once-per-revolution axial sliding motion of the axis of rotation along the axis average line, and can arise, for example, from out-of-square on thrust bearing components.

Regarding face motion, $f(t)$ in equation (5.86) shows that fundamental face motion does exist and is equal to fundamental axial motion. This can be understood by visualizing a perfectly flat disc mounted on a perfect axis of rotation. Mounting error can result in a once-per-revolution sinusoidal face motion (increasing in direct proportion to radius), but this will vanish if the disc is perfectly square to the axis of rotation. Assuming perfect squareness and then changing from a perfect axis to an axis having fundamental axial motion, it follows that the same fundamental motion will occur at all radii. Thus a perfectly flat disc is square to an imperfect axis of rotation if the fundamental face motion is the same at all radii. It is possible to cancel the fundamental face motion by mounting the disc out-of-square to the axis of rotation, but this cancellation can only occur at one radius. The out-of-squareness angle necessary for this cancellation becomes larger as the radius becomes smaller, reaching an impossible situation at zero radius.

The existence of fundamental face motion has an interesting consequence in machining and measuring flat faces. If a flat disc is faced on an axis, which is perfect except for the presence of fundamental axial motion, then the part can be viewed as made up of many flat-faced thin rings, each of which is out-of-square with the axis of rotation by an amount that increases with decreasing radius. Such a part is not flat over its full area. However, if the part is mounted in a roundness-measuring machine with the transducer sensing axially, then the part can be tilted so that no flatness error is sensed during a trace around a circular path concentric with the part centre. Such a part is said to have circular flatness. Since it does not have area flatness, it follows that 'circular' flatness measurements can be misleading if they are not properly understood.

Residual error motion is the general term applied to the difference between average and fundamental error motion. The consequences of residual error motion are analogous to those of average radial motion. For example, residual face motion during machining leads to errors in circular flatness in the same way that average radial motion leads to errors in roundness.

In the general case of error motion with an arbitrary sensitive direction angle φ from the axis average line, the fundamental error motion is proportional to cos φ times the fundamental axial motion (see equation (5.80)). Thus a 45° taper involves 70.7% as much fundamental error motion as a flat face.

### 5.10.3.5 Error motion versus run-out (or TIR)

It should be noted that error motion measurements differ from measurements of run-out or TIR (total indicator reading) in several respects. It is important to understand these differences, since run-out tests have been used extensively in the past in assessing the accuracy of rotational axes. Run-out is defined as 'the total displacement measured by an instrument sensing against a moving surface or moved with respect to a fixed surface'. Under this definition, a radial run-out measurement includes both the roundness error and the centring error of the surface that the gaugehead senses against, and hence radial run-out will be identical to radial motion only if both of these errors are zero. As noted previously, neither of these conditions is easily accomplished. While centring error unavoidably makes the run-out larger than the error motion, it is possible for roundness errors to make the run-out either larger or smaller than the error motion. The latter situation can arise if the surface against which the displacement transducer is sensing was machined in place on the axis bearings, as discussed previously. Similar comments apply to face motion versus face run-out; the latter measurement includes non-squareness and circular flatness errors.

### 5.10.3.6 Fixed sensitive direction measurements

Use of an oscilloscope (figure 5.46) [46] for radial motion measurement with a fixed sensitive direction requires a separate means for generating the base circle. Figure 5.47 shows a method described by Bryan *et al* [35]. Two circular cams, eccentric by 25 mm in perpendicular directions, are sensed by comparatively low-magnification gaugeheads to generate sine and cosine signals for the base circle; single cam with the gaugeheads 90° apart could also be used, as could a vector separator used on the drive [31] in another alternative which produces sine and cosine signals. Radial motion is detected by a third high-magnification gaugehead sensing against a master test ball which is centred (as closely as possible) on the axis average line. The sine and cosine signals are each multiplied by the radial motion signal using Hall effect multipliers and are then fed into the two axes of the oscilloscope. The modulation of the base circle by the signal from the fixed radial motion gaugehead yields a polar plot of radial motion versus the angular position of the axis of rotation. Vanherk has tested a modification in which a small (50 g) commercial synchro unit is physically attached to the axis of rotation to replace the eccentric cams and low-magnification gaugeheads as the sine-cosine signal generator. The advantages are lower cost, less difficulty in obtaining an accurately round base circle, and simplification of the test set-up, with negligible influence on the axis from the synchro attachment except in exacting situations.



**Figure 5.46** Test method of Tlusty for radial motion with a rotating sensitive direction.

**Figure 5.47** Test method for Bryan for radial motion with a fixed sensitive direction.

### 5.10.3.7 *Considerations on the use of the two-gaugehead system for a fixed sensitive direction*

Since the oscilloscope test method of figure 5.47 requires electronic equipment that may not be readily available, it is natural to consider substituting the two-gaugehead system for measuring radial motion with a fixed sensitive direction. If this substitution is made, the resulting radial motion polar plot will not be representative of the potential part out-of-roundness. If $\theta = 0°$ is the fixed sensitive direction, then the polar plot reflects radial motion in this direction only in the vicinity of $\theta = 0°$ and $\theta = 180°$. Moreover, if a given localized movement of the axis of rotation occurring at $\theta = 0°$ appears as a peak on the polar plot, the same movement occurring at $\theta = 180°$ will have an undesired sign reversal and will appear as a valley. At $\theta = 90°$ and $\theta = 270°$, the same movement will not register on the polar plot.

Despite the above observation, it still appears intuitively plausible that the radial motion value should be roughly the same for both fixed and rotating sensitive directions, even if the details of the polar plot are different. This view appears reasonable if the factor of concern is random radial motion. However, for average radial motion, Donaldson [40] has noted a case giving precisely the opposite result, in which an axis that exhibits an elliptical pattern when tested in a fixed sensitive direction is free of radial motion when tested in a rotating sensitive direction. Thus, in the first case it gives the results of equations (5.90) and (5.91):

$$\Delta x(\theta) = A\cos 2\theta \tag{5.90}$$

$$\Delta y(\theta) = A\sin 2\theta. \tag{5.91}$$

However, with a fixed sensitive direction along the $x$ axis, the radial motion polar plot has the equation

$$r(\theta) = r_\text{o} + A\cos 2\theta \tag{5.92}$$

where $r_\text{o}$ is the base circle radius. Equation (5.92) represents an elliptical shape, having a value $r_\text{o} + A$ at $\theta=0°$ and $\theta=180°$ and a value of $r_\text{o} - A$ at $\theta = 90°$ and $\theta = 270°$. The radial motion value based on any of

the polar profile centres is 2$A$. If the sensitive direction rotates with angle $\theta$, the radial motion is given by the equation

$$r(\theta) = r_0 + \Delta x(\theta)\cos\theta + \Delta y(\theta)\sin\theta. \tag{5.93}$$

Figure 5.47 shows the resolution of $\Delta x(\theta)$ and $\Delta y(\theta)$ into components along the rotating sensitive direction that leads to equation (5.93). Combining (5.90) and (5.91) with (5.93) and using the trigonometric identities

$$\cos\alpha\cos\beta = \tfrac{1}{2}[\cos(\alpha - \beta) + \cos(\alpha + \beta)] \tag{5.94}$$

$$\sin\alpha\sin\beta = \tfrac{1}{2}[\sin(\alpha - \beta) - \sin(\alpha + \beta)] \tag{5.95}$$

results in

$$
\begin{aligned}
r(\theta) &= r_0 + \frac{A}{2}(\cos\theta + \cos 3\theta) + \frac{A}{2}(\cos\theta - \cos 3\theta) \\
&= r_0 + A\cos\theta
\end{aligned}
\tag{5.96}
$$

which is the equation of a distorted circle (the limaçon) that is offset from the origin by a distance $A$, and hence the axis would appear to be perfect if tested by the two-gaugehead system.

Two additional comments can be made on the above finding. First, it can be argued that if the offset circle is assessed by concentric circles from the polar chart (PC) centre, then a value of 2$A$ is obtained, as with the fixed sensitive direction. However, there is no way to carry out the initial electronic zeroing to locate the PC centre, since the base circle cannot be generated independently of the polar profile using the test method of figure 5.38. Second, the view might be taken that the above example is a mathematical oddity which is unlikely to occur in practice. In this regard it can be noted that radial motion polar plots commonly exhibit an elliptical pattern, and that to the extent that the overall patterns in the $x$ and $y$ directions contain components as given in equations (5.90) and (5.91), these components will not contribute to the measured radial motion value.

### 5.10.3.8  Other radial error methods

The relative merits of the different configurations have been examined in some papers (e.g. [49, 50]) and various alternatives have been suggested (e.g. [34]). As an example the use of two balls on a common spindle has been suggested (figure 5.48). This is an elegant way of measuring angular error when compared with the two probes used in figure 5.46.



**Figure 5.48** Double-sphere method for axis of rotation error.

The two spherical masters separated by a distance $L$ have to be mounted as concentric as possible with respect to the axis of rotation. The angular motion can be deduced from the difference in the output signals of the probes. Thus $\alpha(t)$ is given by

$$\alpha(t) = \frac{x_1(t) - x_2(t)}{L}. \tag{5.96a}$$

This can be displayed and compensated in the same way as for radial motion.

If the two spheres are not concentric with the axis it is possible to remove the relative eccentricities in a computer.

From what has been shown in this section it is clear that the errors that might be present in a surface metrology instrument or a more general instrument can be very complicated to ascertain, but they can usually be found using some of the methods suggested. However, this is the first step in removing them. Providing that the basic instrument is itself well designed and well made, these errors should not vary with time. If this is so they can usually be removed, or at least considerably reduced, by computer methods as described.

Instrument errors are one problem; another is associated with the signal being measured. It has already been shown in chapters 2 and 3 that the measurement of peaks and other surface parameters digitally can in itself be difficult. A number of errors have been identified. These are mainly a result of sampling, quantization and the numerical model. There are other problems, which are concerned with how the parameter values depend on the type of surface. What is the likely variation? This obviously depends on the parameter in question and the type of surface. In section 5.11 some idea of how to find out such variability will be given.

## 5.11 Variability of surface parameters

Picking the most suitable parameter for surface evaluation has always been a difficult problem. The choice has invariably been between parameters which could be readily measured and were reliable, and those parameters perceived as being the most significant functionally. This has usually meant choosing between average values, such as $R_a$, and peak values such as $R_t$. This does not mean that $R_a$ has no functional significance; it just means that it is usually taken to be one step further away from direct functional usefulness than peaks. Peak measurement and extrema of any sort are unfortunately very likely to have a large spread even within the confines of likely statistical variation and, even more important, extreme value measurement can encompass typical flaws and scratches. This is the reason why early attempts at specifying peak parameters such as $R_t$ invariably relaxed into measuring the mean of a number of estimates giving rise to hybrid parameters such as $R_{tm}$ and $R_{3z}$. In what follows the general issues will be considered. Fundamentally the reliability is determined by the number of independent degrees of freedom in the estimate. What these degrees of freedom are for each parameter and how the number of them is determined in a few examples will be given next. In addition there is a further requirement, which is to know to what extent the instrumental procedure can influence the value and spread of values of a parameter.

Assume that a surface $z(x)$ is random. Any roughness parameter $P$ (using Lukyanov and Lisenko's terminology [51]) can be imagined to be formed from the profile by a transformation $T$. This transformation need not be linear. In fact it rarely is. Hence the mean parameter $\bar{P}$ is given by

$$\bar{P} = E[P] = \frac{1}{L}\int_0^L T(z(x))\mathrm{d}x \tag{5.97}$$

where $L$ is the assessment length.

The variance of $P$ is

$$\text{var}(P) = E[P^2] - [\overline{P}^2]$$

$$= \frac{1}{L^2}\int_0^L\int_0^L A_T(x_2 - x_1)\mathrm{d}x_1\mathrm{d}x_2 \tag{5.98}$$

$$\text{var}(P) = \frac{1}{L^2}\int_0^L\int_0^L A_T(x_2 - x_1)\mathrm{d}x_1\mathrm{d}x_2 \tag{5.99}$$

because $A^T(x^2 - x^1) = E[T(z(x_1))T(z(x_2))] - E[P]^2$ assuming stationary statistics. Notice the central role of second-order statistics in calculations of variability. The two-dimensional probability function can be mapped in one dimension using suitable reorganization of the variables $x_1 x_2$ in terms of $\tau$. Thus

$$\text{var}(P) = \frac{2}{L}\int_0^L (1 - \tau/L)A_T(\tau)\mathrm{d}\tau \tag{5.100}$$

as $L \gg T$ for practical cases. Here, to avoid mathematical problems $AT(\tau)$ must be zero for $\tau => L$. Therefore

$$\text{var}(P) \sim \frac{2}{L}\int_0^L A_T(\tau)\mathrm{d}\tau \sim 2\sigma_T^2 \frac{\tau_\mathrm{c}}{L}. \tag{5.101}$$

This says what has already been said in words. The variance of a parameter for measuring roughness is inversely proportional to the number of degrees of freedom in the assessment length, that is $L/TC$ where $TC$ is the correlation length of the autocorrelation function $A_T(\tau)$ and $\sigma_T^2$ is the variance ($=A_T(0)$).

The point of importance here is that only $\tau_c$ is needed, not the shape of the autocorrelation of $T(z(x))$ [52].

As an example consider the $R_\mathrm{a}$ value of a roughness profile

$$T(z(x)) = |z(x)|. \tag{5.102}$$

If the surface is Gaussian the value of $\sigma_T$ is found from $A_{R\mathrm{a}}(\tau)$. Now

$$A_{R_\mathrm{a}}(\tau) = \int_0^\infty\int_0^\infty z_1 z_2 p(z_1 z_2)\mathrm{d}z_1\mathrm{d}z_2 - \int_0^\infty\int_{-\infty}^0 - \int_{-\infty}^0\int_0^\infty + \int_{-\infty}^0\int_\infty^0 \ldots \tag{5.103}$$

$$A_{R_\mathrm{a}}(\tau) = \frac{2\sigma^2}{\pi}\{[1 - A(\tau)]^{1/2} + A(\tau)\sin^{-1}(A(\tau)) - 1\} \tag{5.104}$$

giving

$$\sigma_T^2 = R_\mathrm{a}^2\left(\frac{\pi}{2} - 1\right) = \sigma^2\left(1 - \frac{2}{\pi}\right). \tag{5.104a}$$

Hence

$$\text{var}(R_\mathrm{a}) = (\pi - 2)R_\mathrm{a}^2\tau_\mathrm{c}/L. \tag{5.105}$$

The coefficient of variation, $\text{var}(R_\mathrm{a})/E[R_\mathrm{a}]^2$, is found from (5.105). The apparent reduction of variance to $\sigma_\mathrm{t}^2$ from $\sigma^2$ in equation (5.104) is necessary because the mean line of the rectified signal for $R_\mathrm{a}$ is not at $z = 0$ but at $\sqrt{\frac{2}{\pi}}\cdot\sigma$. Hence the coefficient of variation is

$$(\pi - 2)\frac{\tau_c}{L} \sim \frac{(\pi - 2)}{N} \tag{5.106}$$

where $N$ is the number of degrees of freedom in the chart or assessment length.

Expression (5.106) represents the coefficient of variation for analogue signals. This can be extended to digital signals, in which case sampling and quantization also play a part.

Following the same procedure as for the analogue only, making the discretization approximations $\Delta z$ for the quantization interval and $\delta x$ for the sampling, the discrete autocorrelation function $A_T(i\Delta x)$ can be obtained from the analogue version.

$$A_T(\tau) = \int_{-\infty}^{\infty} \int z_1 z_2 p(z_1, z_2) \mathrm{d}z_1 \mathrm{d}z_2 \tag{5.106a}$$

by making the assumption, as before, that $p(z_1, z_2)$ is a multinormal distribution and therefore can be expanded in terms of Hermite polynomials

$$p(z_1, z_2) = \frac{1}{\sigma^2} \sum \varphi^{(i+1)}(z_1) \varphi^{(i+1)}(z_2) A^i \left( \frac{x_2 - x_1}{i!} \right) \tag{5.107}$$

and

$$\varphi^{(i+1)}(z) = (-1)^i \varphi^{(1)}(z) H_i(z) \tag{5.108}$$

where $H_i(z)$ is the Hermite polynomial equal to

$$\exp(-z^2)\frac{(\mathrm{d}^i \exp(z^2)}{\mathrm{d}x^i}. \tag{5.109}$$

The equivalent of equation (5.109) is for zero quantization:

$$\mathrm{var}(A_T(q\Delta x)) = \frac{2}{n^2} \sum_{q=0}^{n-1} [n - |q| A_T(q\Delta x)] \tag{5.110}$$

where $n$ is the number of samples, $q$ is the discrete lag number

$$\mathrm{var}(A_T(q\Delta x)) = \frac{\Delta x}{L}(\pi - 2)R_{a_\mathrm{D}}^2 \left[ 1 + 2\sum_{q=1}^{n-1} \left( 1 - \frac{q}{n} \right) A_T'(q\Delta x) \right] \tag{5.111}$$

where $A'$ is normalized with respect to $R_q^2$ of the profile and $R_{a_\mathrm{D}}$ is the discrete $R_a$.

For the case when $q \ll n$ and the ordinates are uncorrelated, $A_T = 0$, $q\Delta x = \tau_c$ and

$$\mathrm{var}(R_{a_\mathrm{D}}) = \frac{\Delta x}{L}(\pi - 2)R_{a_\mathrm{D}}^2. \tag{5.111a}$$

The coefficient of variation is given by

$$(\pi - 2)\frac{\tau_c}{L} \sim \left( \frac{\pi - 2}{n} \right) \tag{5.111b}$$

as before in equation (5.106). If the ordinates are correlated and the correlation can be approximated by an

exponential function then some idea of the effect can be seen, that is

$$\mathrm{var}(R_{a_D}) \sim \left( \frac{\pi - 2}{n(1 - \rho^2)} \right) \qquad (5.112)$$

where $p$ is the correlation at $(q\Delta x)$. The term $n(1 - p^2)$ is the effective number of independent ordinates.

The difference between the discrete and the analogue values for the variability in the $R_a$ value of a surface profile is that, in the latter case, the correlation length only is needed, whereas in the former the correlation function is needed because the data points of a digitized profile should be correlated, otherwise curvature and slope information can be jeopardized. Actually if the form of the correlation function is required quickly, it can usually be estimated from a profile chart by evaluating $k$, the ratio of the zero-crossing density to that of peaks shown in the typology estimations in section 2.1.3. This gives some idea of the type of surface.

Such considerations of the variability in surface parameters are important because many features not directly associated with roughness can be affected. One case, as an example, is the position of the centre of a best-fit circle for a set of roundness data discussed earlier in chapter 3 on processing.

Another important parameter that needs to be assessed for variability is the $R_q$ value. The variance in the estimate of the $R_q$ value is similar to that of the $R_a$ value, only it does not have the mean line shifted by $\sqrt{2/\pi}$ but by $2/\pi$.

Traditionally, the accepted values of $n\sqrt{1 - \rho^2}$ and $N$ should be 20 or more (within a sampling length preferably). Sometimes this cannot be achieved, especially if the surface is small. Here $n$ is the digital sample number and $N$ is the number of correlation lengths.

For a significance level in $R_a$ of 5% such a rule is adequate and is allowed for in most standards.

Other parameters can be treated [53] but may be more difficult. This is especially true for peak parameters where amplitude density functions have to be assumed. Some possibilities are the log normal distribution or gamma function. An important parameter case is that of $S_m$, the zero-crossing parameter which is equivalent in length characterization to that of $R_a$ in height. This is related to the average wavelength $\lambda_a$ [54]. Crossing estimates for variability are very difficult to get even though, as in the case of $\lambda_a$ (or $S_m$), they may be effectively averages. The variability of $\lambda_a$ is larger than for $R_a$ because of the nature of its definition: $\lambda_a = 2\pi R_a/\Delta_a$. The slope is badly affected by high-frequency noise so that $\lambda_a$ could be estimated realistically as being four or five times more variable than $R_a$ for the same bandwidth and assessment length. A similar argument cannot be applied so easily to the estimate of the correlation length. $S_m$ is usually estimated from the crossings of the mean line. This is relatively reliable provided that the high frequencies are ignored (which means ignoring the effect of the local peaks and their associated count $S$).

Many of the parameters of roughness are related (as in the case above for $\lambda_a$). Also $\sigma^*\eta R$ is a constant, where $\sigma^*$ is the peak standard deviation, $\eta$ is $S$ and $R$ is the average radius of the peaks. Another check can be obtained between the correlation length and $S_m$, for a random wave. Thus if $\tau_c$ is the correlation length then there will be a zero crossing every $N$ correlation lengths where

$$N = \sum_{i=2}^{\infty} i \left( \tfrac{1}{2} \right)^i \Big/ \sum_{i=2}^{\infty} \left( \tfrac{1}{2} \right)^i = 3 \qquad (5.113)$$

so that $\lambda_a \sim 6\tau_c = S_m$.

If the correlation function is known then $\tau_c$ can be estimated from the mean line itself. Thus if the mean line is the mid-point locus line (chapter 3), it has its own correlation function given by $A_m(\tau)$ where

$$A_m(\tau) = \frac{2}{L} \int_0^L \left( 1 - \frac{\tau}{L} \right) A(\tau) \mathrm{d}\tau. \qquad (5.114)$$

Remembering that the correlation length $\tau_c$ is $\int_0^\infty A(\tau)\mathrm{d}\tau$ and letting $T_m = \sqrt{-A_m(0)/A_m''(0)}$ gives another approximate relationship between $\tau_c$ and crossings, namely

$$\tau_c \sim \frac{T_m^2}{L\pi^2} \tag{5.115}$$

where $T_m$ can be found by the crossings of the mean line with a straight line drawn in the general direction of the surface profile.

Note:

It can therefore be seen that variation in any parameter $P$ of surface roughness can be due as much to how it is evaluated as to statistical variations across the surface.

The situation described above is even more complex if areal properties are being evaluated.

Also, cases can be considered that take into account the sum of more than one type of surface signal. The usual case is that of a long-wavelength signal together with a shorter-wavelength surface signal. Such a combination is often encountered in surface metrology. One instance is the presence of waviness on a ground specimen, in which case the extraneous signal would be the low one. Another case is the presence of micro-roughness in turning. In this case the turning marks would be the required signal and the variability would be the random high frequency. In most cases the situation is simplified if it is considered to be the sum of two random signals, the longer-wavelength one being of a narrow-band spectral type. A number of people have investigated this type of signal, particularly with respect to crossing densities and extreme behaviour [55]. Lukyanov and Lisenko have examined some cases [51] and report on them for broaching, grinding and simulated grinding. The results confirm the variability for $R_a$ but indicate that it is not improbable to get variations in $S_m$ values of 30% or more. Surprisingly this variation is tolerated in the ISO TC57 standard, More work needs to be done in this direction, but it is reassuring to note that the rule of thumb for the variability of parameters using the correlation length as the basic yardstick is acceptable.

## 5.12 National and international standards

### 5.12 General

With the advent of global economies the need for countries to have common standards has increased dramatically. This means that the role of ISO (The International Standards Organization) has become very important. All member countries look to ISO for setting and supporting standards. A block diagram Figure 5.49 shows the relationship

Within a country there are a number of ways in which the central agency, in the UK the NPL, oversees traceability from industry to the approved international standards. There is set up a facility NAMAS which is a facilitator of quality problems which is intermediate between industry and the NPL. NAMAS is a set of



**Figure 5.49** International standard structure.

**Figure 5.50** UK structure for traceability.

centres of excellence which is spread around the country and is accessible to all. Typical laboratory equipment, operators and procedures are all in the chain of standards.

### 5.12.1    *Geometrical Product Specification (GPS)*

[For full master plan see ISO/TR 14638 1995(E) pages 228–257 in ISO Standard 'Limits Fits and Surface Properties' (1999)]

The concept of GPS covers several kinds of standards: some deal with the fundamental rules of specifications (fundamental GPS standards); some deal with global principles and definitions (global GPS standards) and some deal directly with the geometrical characteristics (general and complementary GPS standards). These are shown in Figure 5.51.

| Fundamental GPS standards | Global GPS standards | | | | | |
|---|---|---|---|---|---|---|
| | General GPS matrix | | | | | |
| | Chain link number | 1 | 2 | 3 | 4 | 5 | 6 |
| | Size | | | | | | |
| | Distance | | | | | | |
| | Radius | | | | | | |
| | Angle | | | | | | |
| | Form of line independent of datum | | | | | | |
| | Form of line dependent on datum | | | | | | |
| | Form of surface independent of datum | | | | | | |
| | Form of surface dependent on datum | | | | | | |
| | Orientation | | | | | | |
| | Location | | | | | | |
| | Circular run-out | | | | | | |
| | Total run-out | | | | | | |
| | Datums | | | | | | |
| | Roughness profile | | | | | | |
| | Waviness profile | | | | | | |
| | Primary profile | | | | | | |
| | Surface imperfections | | | | | | |
| | Edges | | | | | | |

**Figure 5.51** The GPS matrix.

Also Figure 5.52 includes the complementary GPS Matrix.



**The Global GPS Standards**

GPS standards or related standards which deal with or influence several or all General GPS chains of standards

**General GPS Matrix**

General GPS chains of standards

1. The **Size** chain of standards
2. The **Distance** chain of standards
3. The **Radius** chain of standards
4. The **Angle** chain of standards
5. The **Form of a line** (independent of a datum) chain of standards
6. The **Form of a line** (dependent of a datum) chain of standards
7. The **Form of a Surface** (independent of a datum) chain of standards
8. The **Form of a Surface** (dependent of a datum) chain of standards
9. The **Orientation** chain of standards
10. The **Location** chain of standards
11. The **Circular run-out** chain of standards
12. The **Total run-out** chain of standards
13. The **Datums** chain of standards
14. The **Roughness profile** chain of standards
15. The **Waviness profile** chain of standards
16. The **Primary profile** chain of standards
17. The **Surface defects** chain of standards
18. The **Edges** chain of standards

**Complementary GPS Matrix**

Complementary GPS chains of standards

**A. Process specific tolerance standards**

A1. The **Machining** chain of standards
A2. The **Casting** chain of standards
A3. The **Welding** chain of standards
A4. The **Thermal cutting** chain of standards
A5. The **Plastic moulding** chain of standards
A6. The **Metallic and inorganic coating** chain of standards
A7. The **Painting** chain of standards

**B. Machine element geometry standards**

B1. The **Screw thread** chain of standards
B2. The **Gears** chain of standards
B3. The **Splines** chain of standards

**The Fundamental GPS Standards**

**Figure 5.52** The GPS matrix model—GPS masterplan — overview.

There are some definitions which are useful. Although this handbook is concerned primarily with surface properties they cannot be divorced from more general geometrical characterization.

(a) Fundamental GPS standards: These establish the fundamental rules for the GPS dimensioning and tolerancing of workpieces and products. An example is ISO 8015 which is a standard (1985), Technical Drawings—Fundamental Tolerancing Principle.

(b) Global GPS standards : An example is ISO 370 (1975), Tolerance dimensions — conversion from inches into mm and vice versa.

(c) General GPS standards: An example is ISO 1302 (1992), Technical drawings—methods of indicating surface texture.

(d) Complementary GPS standards: An example is process-specific tolerance standards 1, machining 2, casting etc.

### 5.12.2 Chain of standards

This refers to all related standards concerning the same geometrical characteristic(e.g. primary profile).

(a) A chain of standards is such that each single standard (i.e. each link) affects the other standards in the chain so that a full understanding requires a knowledge of all standards in the chain.

(b) The task of all the chains is to link symbol unambiguously to the SI unit of length in such a way that the tolerance limits are defined in every case possible—irrespective of the deviations from ideal geometry of the toleranced features and other deviations from theoretical correct conditions.

(c) Each of the six links in the chain has a specific task. These are:-

1. Chain link 1—Product documentation indication codification. These deal with the drawing indication for the characteristic of the workpiece—it can be a symbol.

2. Chain link 2—Definition of tolerances—theoretical definition and values. These define the rules of translating from the 'code' of link 1 to 'human understandable' and 'computer understandable' values into SI units.

3. Chain link 3—Definitions for actual feature—characteristic or parameter. These extend the meaning of the theoretically exact feature so that the non-ideal real world (actual feature) is always unambiguously defined in relation to the tolerance indication (code symbol) on the drawing. The definition of the actual feature characteristics in this link would be based on data points.

4. Chain link 4—Assessment of the deviations of the workpiece—comparison with tolerance limits. This defines the detailed requirement for the assessment of the deviations of the workpiece from the one on the drawing, taking into account the definitions in chain links 2 and 3.

5. Chain link 5—Measurement equipment requirements. These describe specific measuring equipment or types of measuring instrument. The standard may include values for the limits of maximum permissible error for the defined characteristic of the measuring equipment.

6. Chain link 6—Calibration requirements—measurement standards. These describe the calibration standards and the calibration procedures to be used. Also verifying functional requirement of the specific measuring equipment (limits of permissible error) in chain links with traceability to the definition of the SI unit concerned.

Figure 5.53 shows how the chain link of standards applies to a specific geometrical feature—in this case the surface texture. There are others included in the same family (with different ISO references) e.g. straightness, roundness etc.

### 5.12.3 Surface standardization

Although the need for surface metrology standards has been acknowledged since the 1920s and despite some countries having their own standards the ISO version has taken a long time to consolidate. This is because the various camps have found it difficult to compromise. For many years the UK developed the 'M' system based on filter 'mean' line in which all the profile was used in order to establish a reference line from which to measure. This was used primarily throughout the English speaking world and was championed by Dr R E

| Chain link number | | 1 | 2 | 3 |
|---|---|---|---|---|
| Geometrical characteristic of feature | Parameters | Codification on a drawing | Definition of tolerance | Definitions for actual feature |
| Surface roughness | M-system $R_a$ | 1302 | 4287–1,–2, 468 | 4288 |
| | M-system other | | 4287–1, 468 | |
| | Motif method R | | | 12086 |

| Chain link number | | 4 | 5 | 6 |
|---|---|---|---|---|
| Geometrical characteristic of feature | Parameters | Comparison with tolerance limits | Measurements equipment requirements | Calibration requirements |
| Surface roughness | M-system $R_a$ | 4288, 2632–1,–2 | 3274, 1878, 1879, 1880, 2632, 11562 | 5436, 2632 |
| | M-system other | | 3274, 1880, 11562 | |
| | Motif method R | | | 12086 |

**Figure 5.53** The surface roughness chains of standards.

Reason. On the other hand the 'E' system pioneered by Professor van Weingraber in Hanover, Germany was put forward as an alternative. Later still the French under the leadership of Professor Biel came up with the 'Motif' system. These two later systems relied on applying rules to the peaks and valleys of the profile.

The ISO standard of the 1950s referred mainly to the 'M' system because of the ease by which the reference could be established instrumentally. The 'E' system and then the Motif system were graphical until recently when the use of computers enabled these reference lines to be calculated readily. Both the E and Motif systems react clearly to changes on the surface geometry such as patches or high peaks. This sensitivity makes them more function-orientated, which is a good point, but more variable, which is bad for process control. The clear distinction between the reference line procedures and their suitability is usually forgotten.

There is another factor which has influenced the development of ISO standards for surface metrology. This is the influence of the country which holds the secretariat.

After World War II the secretariat for texture and roundness was offered to the UK. This offer was promptly rejected. The same outcome occurred in a number of other countries. Eventually after a period of over ten years, the USSR accepted the task and as a result the responsibility for the ISO standard on surface metrology was that of Gosstandart under Director V. V. Boitsov. (Gosstandart being the name given to the state committee for standards of the USSR council of ministers). Surface metrology came under the heading of measuring instruments for product quality improvement.

Immediately after the USSR took control the procedures became formalized in some detail; the somewhat nebulous documentation of the past became part of the powerful Gosstandards. In the USSR it was mandatory to follow the standards set down by the state. The various industries and institutions had no option but to obey the directives laid down. On the other hand, following the standards in the US, Europe and Japan was voluntary. There was some reluctance to undertake all the detail set out in the standards.

Despite the fact that the job of the secretariat was to collate and unify surface metrology worldwide and to formulate proposals for international approval, the actual documents inevitably took on the character of the country. In the case of the USSR the characteristics were that of being highly theoretical and very detailed.

In some respects they were too complicated for industrial users in the UK and US in particular. The USSR approach was scientific rather than 'user friendly' In principle the USSR was aware of the need for user bias because surface metrology came under VNIIMS — the state metrology service based in Moscow rather than VNIIM — the state metrology institute based in what was then Leningrad. There were some problems with foreign instrument makers wishing to sell in the USSR. All instruments had to be vetted by the state factory 'KALIBR' before any sale could take place. This caused delay and some frustration because most makers were reluctant to divulge commercial secrets.

Fortunately, the director of surface metrology management Dr VS Lukyanov was both practical and theoretical and initiated good work on instruments error analysis and practical artefacts for calibration. Many of the ideas on referred standards (e.g. for cylindricity and other forms) originated with him. Also parameters such as $S$ and $S_m$ in present day form were formulated properly in Moscow. During the late eighties, for some reason, roundness and form were separated from surface texture. More recently the secretariat has been taken over by Germany and is now more user based.

### 5.12.4    *Role of technical specification documents*

Many of the procedures and topics discussed in chapters 2 and 3 involve various aspects of sampling and filtering. The theory behind these considerations is well understood. This does not mean that they are suitable for insertion into national and international standards documents. It could be that the technique is sensitive to noise, as in differentiation of the profile to get slopes and/or curvatures. For this reason any method or parameter should be thoroughly investigated before adoption to avoid unnecessary variability of the measured results.

A great many of the problems associated with surface properties are due to filtering in order to separate different components of the profile and areal signal. Another main source of variation is in sampling. In what follows, some of the issues currently being investigated are outlined. Nothing in this should be interpreted as an ISO standard. Also not all the complexity is necessary: mathematical analysis is far easier to carry out than practical relevance. Ideally each potential technique or parameter should be tested for usefulness before it is incorporated into a standard. Unfortunately this is difficult to enforce, so many marginally important techniques and parameters end up in the standard, hopefully to be filtered out in time.

One of the roles of the technical committees supporting each branch of the subject is to propose, examine, clarify and submit subjects that could be incorporated into future standards. The main task of a technical committee is to prepare international standards but it can propose the publication of a technical report of one of the following types:

Type 1: a publication with the complete support of ISO
Type 2: a technical development which has the potential of being included in a future standard.
Type 3: a technical committee which has collected data that may be generally useful such as a 'state of the art' contribution.

An example of such technical committee work is ISO TC 213/Ag 9 which has resulted in ISO/TS 16610 work on geometrical product specifications (GPS)—data extraction techniques by sampling and filtration—(secretariat UK). The letters TS mean 'technical specification.'

The work is split up into a number of parts:

1. Basic terminology
2. Basic concepts of linear filters

3. Spline filters
4. Spline wavelets
5. Basic concepts of morphological operations and filters
6. Morphological operations and filters
7. Morphological scale space techniques
8. Robust spline filters
9. Basic concepts of linear surface filters.

The breakdown of the work in TS 16610 shows the intention of the committee.

Whilst it is to be applauded that a continuous stream of analytical methods are submitted for scrutiny two factors should be borne in mind:

(a) Many of the proponents of current suggestions are not aware of past examination of the same or similar function and as a result have not benefited from previous experience. These omissions are usually obvious from sparse bibliographies. In the case above phase-corrected filters and splines are good examples.

(b) Not many of the TS documents provide evidence of practical results, with the result that there is no weighting of the suggestions made within the specific TS. The benefits of most suggested functions do not warrant inclusion into what is essentially a practical working standard.

All the topics come under the general heading of geometrical product specification (GPS) — data extraction techniques by sampling.

### 5.12.6 *Selected list of international standards applicable to surface roughness measurement: methods; parameters; instruments; comparison specimens*

| | |
|---|---|
| ISO 1878 | Classification of instruments and devices for measurement and evaluation of the geometrical parameters of surface finish. (1983). |
| ISO 1879 | Instruments for the measurement of surface roughness by the profile method—vocabulary. (1981). |
| ISO 3274 GPS | Instruments for the measurement of surface roughness by the profile method—contact (stylus) instruments of consecutive profile transformation—(1996). |
| ISO 4287 GPS | Surface texture—profile method—terms, definitions and surface texture parameters (1997). |
| ISO 4288 GPS | Rules and procedures for the measurement of surface texture profile method (1996). |
| ISO 4291 | Methods for the assessment of departures from roundness—measurement of variations in radius. (1985). |
| ISO 4292 | Methods for the assessment of departure from roundness—measurement by two and three point methods. (1985). |
| ISO 5436 | Calibration specimens—stylus instruments—types, calibration and use of specimens. |
| ISO 6318 | Measurement of roundness—terms, definitions and parameters of roundness. |
| ISO 8503/1 | Preparation of steel substrates before application of point and related products—surface roughness characteristics of blast cleaned steel substrates—Part 1: Specifications and definitions of ISO surface profile comparators for the assessment of abrasive blast cleaned surfaces. (1988). |
| ISO 8785 GPS | Surface imperfections—terms, definitions and parameters. (1998). |
| ISO 11562 GPS | Metrological characterization of phase corrected filters—Surface texture. Profile method (1996). |
| ISO 12085 GPS | Surface texture. Profile method—motif parameter. (1996). |

| ISO 13565/1 GPS | Texture profile method—surface having stratified functional properties—Part 1: Filtering and overall measuring conditions (1996). |
| ISO 13565/2 GPS | Characterization of surface texture profile method surfaces having stratified functional properties—Part 2: Height characterization using the linear material ratio curve. |
| ISO 13565/3 GPS | Surface texture profile method—surfaces having stratified functional properties—Part 3: Height characterization using the linear material ratio curve. (2000). |

*5.12.7   International (equivalents, identicals and similars)*

| ISO 468 | Equivalent to AS 2536 (Australia). Similar to BS 1134 Part 1. (1988), BS 1134 Part 2. (1972). |
| ISO 1880 | Similar to DIN 4772. |
| ISO 1302 | Implemented in BS 308 Part 2. (1985). |
| ISO 2632/1 | Identical to BS 2634 Part 1. (1987). Equivalent to DIN 4769 Parts 1, 2 and 3. Equivalent to NFE 05-051. |
| ISO 2632/2 | Identical to BS 2634 Part 2. (1987). Equivalent to NFE 05-051. |
| ISO 2632/3 | Identical to BS 2634 Part 3. (1980). Equivalent to NFE 05-051. |
| ISO 3274 | Similar to DIN 4768. Similar to DIN 4772. Equivalent to NFE 05-052. |
| ISO 4287/1 | Similar to BS 1134 Part 1. (1988). Identical to BS 6741 Part 1. (1987). Equivalent to DIN 4762. Similar to NFE 05-015. |
| ISO 4287/2 | Identical to BS 6741 Part 2. (1987). |
| ISO 4288 | Similar to DIN 4768. Equivalent to DIN 4768 Part 1. Equivalent to DIN 4775. Similar to NFE 05-054. |
| ISO 4291 | Similar to BS 3730 Part 2. (1982). Identical to BS 6740. (1986). Similar to NFE 10-103. |
| ISO 4292 | Similar to BS 3730 Part 3. (1987). |
| ISO 5436 | Identical to BS 6393. (1987). |
| ISO 6318 | Identical to BS 3730 Part 1. (1987). Similar to NFE 10-103. |
| ISO 8503 | Identical to BS 7079. (1989). |

Note:
Because of the rapidly changing national standards scene it is difficult to compile a current accurate list. The following is a representative list for guidance only. For updated national standard information consult the following: International Organization for Standardization, Central Secretariat, PO Box 56, CH 1211 Geneva 20, Switzerland. The following are not necessarily the latest versions. The preferred practice is to use the ISO standards directly.

### *Belgium*
| NBN 863 | Basic concepts and standardized data for surface roughness. (1970). |

### *Canada*
| CSAB 95 | Surface texture: roughness, waviness and lay. (1962). |

### *Czechoslovakia*
| CSN 01 4450 | Surface roughness, parameters and definitions. (1980). |
| CSN 01 4456 | Evaluation of surface roughness of castings. (1983). |

| CSN 01 4451 | Surface roughness. Basic parameters, numerical values. (1980). |
| CSN 25 2302 | Surface roughness. Comparison specimens. Technical requirements. (1980). |
| CSN 73 2582 | Test for abrasive resistance of surface finish of building structures. (1983). |

### *France*

| NF E 05 015 | Surface texture of products: regulations, general terminology, definitions. (1984). |
| NF E 05 016 | Surface texture of products: specification of surface texture on drawings. (1978). |
| NF E 05 018 | Surface texture of products: economic aspects. (1969). |
| NF E 05 050 | Surface texture. Synoptic tables for electronic pick up equipment. (1970). |
| NF E 05 051 | Roughness comparison specimens. (1981). |

### *Germany*

| DIN 4760 | Form deviations; concepts; classification system. (1982). |
| DIN 4761 | Surface character; geometrical characteristics of surface texture terms; definitions, symbols. (1978). |
| E DIN 4762 | Part 1. Surface roughness; Terminology. (1981). |
| DIN 4763 | Progressive ratio of number values of surface roughness parameters. (1981). |
| DIN 4764 | Surfaces of components used in mechanical engineering and light engineering; terminology according to stress conditions. (1982). |
| DIN 4765 | Determination of the bearing area fraction of surfaces, terms. (1974). |
| DIN 4766 | Surface roughness associated with types of Part 1 manufacturing methods; Attainable arithmetical mean value of peak-to-valley height. $R_z$ according to DIN 4768 Part 1. (1981). |
| DIN 4766 | Surface roughness associated with types of Part 2 manufacturing method; Attainable arithmetical mean value $R_a$ according to DIN 4768 Part 1. (1981). |
| DIN 4768 | Determination of surface roughness $R_a$, Part 1 $R_z$, $R_{max}$ with electric stylus instruments; Basic data. (1974). |
| DIN 4768 Supplmt. 1 | Determination of roughness parameters $R_a$, Part 1 $R_z$, $R_{max}$, by means of electrical stylus instruments; conversion of parameter $R_a$, to $R_z$, and vice versa. (1978). |
| DIN 4769 | Roughness comparison specimens; Part 1 Technical conditions of delivery; application. (1972). |
| DIN 4769 | Roughness comparison specimens; Part 2 Surfaces produced by cutting with periodic profile. (1972). |
| DIN 4769 | Roughness comparison specimens; Part 3 Surfaces produced by cutting with a periodic profile. (1972). |
| DIN 4769 | Roughness comparison specimens; metallic Part 4 Surfaces, cleaning with abrasives. (1974). |
| DIN 4771 | Measurement of the profile height P, of surfaces. (1977). |
| DIN 4772 | Electrical contact (stylus) instruments for the measurement of surface roughness by the profile method. (1979). |
| DIN 4774 | Measurement of the depth of waviness by means of electrical contact stylus instruments. (1981). |
| DIN 4775 | Measuring the surface roughness of workpieces; visual and tactile comparison, methods by means of contact stylus instruments. (1982). |

### *Italy*

| UNI 3963 | Surface geometrical errors: Part 1 General definitions. (1978). |
| UNI 3963 | Surface geometrical errors: Part 2 Surface roughness. (1978). |

### Japan

| | |
|---|---|
| JIS B 0031 | Method of indicating surface texture on drawings. (1982). |
| JIS B 0601 | Definitions and designations of surface roughness. (1982). |
| JIS B 0610 | Waviness. (1976). |
| JIS B 0659 | Roughness comparison specimens. (1973). |
| JIS B 0651 | Instruments for the measurement of surface roughness by the stylus method. (1976). |
| JIS B 0652 | Instruments for the measurement of surface roughness by the interferometric method. (1973). |
| JIS B 7451 | Roundness measuring machines. (1991). |

### Netherlands

| | |
|---|---|
| NEN 3631 | Surface roughness: terms and definitions. (1977). + *corrigendum* April, 1982. |
| NEN 3632 | Surface roughness: evaluation. (1974). |
| NPR 3633 | Surface roughness. Various characteristics of surface roughness. (1978). |
| NEN 3634 | Surface roughness: method of indicating surface texture on drawings. (1977). + *corrigendum* February, 1979. |
| NEN 3635 | Surface roughness. Measurement of $R_a$ roughness. (1980). |
| NEN 3636 | Surface roughness. Roughness comparison specimens. (1980). |
| NPR 3637 | Surface roughness: directives for the relation between the function of a workpiece surface and the roughness value $R_a$. (1983). |
| NPR 3638 | Surface roughness: directives for the attainable roughness values $R_a$ for various machining operations. (1983). |

### Spain

| | |
|---|---|
| UNE 1037 | Method of indicating surface texture on drawings. (1983). + *erratum*. |
| UNE 82 803 | Instruments for measurement of surface roughness by the profile method. Vocabulary. (1976). |
| UNE 82 804 | Instruments for measurement of surface roughness by the profile method. Contact (stylus) instruments and progressive profile transformation. (1976). |
| UNE 82 302 | Classification of instruments and devices for the measurement and evaluation of the geometrical parameters of surface finish. (1976). |

### Sweden

| | |
|---|---|
| SMS 671 | Surface roughness: terminology. (1975). |
| SS 572 | Method of indicating surface texture on drawings. (1981). |
| SMS 673 | Surface roughness. Roughness criteria. Standard values. Sampling lengths and cut-off values. (1975). |
| SMS 674 | Surface roughness. Guidance for the choice of surface roughness. (1975). |
| SMS 675 | Surface roughness. Measurement of surface roughness by means of electrical profile recording instruments. (1975). |

### Switzerland

| | |
|---|---|
| VSM 58070 | Surface typology: geometric characteristics (SNV 258070) of surface texture. (1976). |
| VSM 58101 | Surface roughness—stylus instruments for (SNV 258101) the measurement of roughness according to system M. (1976). |
| VSM 58102 | Surface roughness—instructions for the (SNV 258102) use of instruments having a stylus. (1976). |
| VSM 58250 | Surface roughness—analysis of surfaces of (SNV 258250) a part by the measurement of roughness. (1978). |

### UK

| | |
|---|---|
| BS 1134 | Assessment of surface texture methods and instrumentation. Part 1. (1988). |
| BS 1134 | Guidance and general information. Part 2. (1990). |
| BS 2634 | Roughness comparison specimens turned, ground, bored, milled, shaped and planed. Part 1. (1987). |
| BS 2634 | Spark eroded, shot blasted, grit blasted, polished. Part 2. (1987). |
| BS 2634 | Cast specimens. Part 3. (1980). |
| BS 3730 | Assessment of departures from roundness—glossary of terms. Part 1. (1987). |
| BS 3730 | Methods of determining departures from roundness using two and three point measurement. Part 3. (1982). |
| BS 6393 | Calibration of stylus instruments. (1987). |
| BS 6740 | Determining departures from roundness by measuring variations in radius. (1987). |
| P06461 | Vocabulary of metrology—basic and general terms. (1985). |
| BS 6741 | Surface roughness—terminology; surface and its parameters. Part 1. (1987). |
| BS 6741 | Surface roughness—terminology; measurement of surface roughness parameters. Part 2. (1987). |

### USA

| | |
|---|---|
| ANSI B.46.1 | Surface texture: surface roughness, waviness and lay. (1978). |
| SAE AS 291 | Surface finish (RMS). |
| SAE AS 291 D | Surface texture, roughness, waviness and lay. (1964). Inactive for new design after March 31, 1979. |
| SAE J 448 | Surface texture. |
| SAE J 449 | Surface texture control. |
| ASME B.89.3.1 | Measurement of out of roundness. (1988). |
| ASME Y.14.36 | Surface texture symbols—drafting standards. (1978, R 1993). |
| ASTM F.I 048 | Test method for surface roughness (total integrated light scatter). (1987, R 1992). |

### USSR (Russia*)

| | |
|---|---|
| GOST 8.296 | State special standard and all-union verification schedule for measuring surface roughness parameters $R_{max}$ and $R_z$. (1978). |
| GOST 2789 | Surface roughness: parameters and characteristics. (1973). |
| GOST 19299 | Instruments for the measurement of surface roughness by the profile method: types, main parameters. (1973). |
| GOST 8.300 | State system for the uniformity of measurements. Roughness comparison specimens. Calibration methods and means. (1978). |
| GOST 8.241 | State system for ensuring the uniformity of measurements. Contact profile meters system. Methods and means of verification. (1977). |
| GOST 8.242 | State system for ensuring the uniformity of measurements. Profile recording instruments. Methods and means of verification. (1977). |
| GOST 9847 | Optical surface roughness measuring instruments: types, dimensions and standards of accuracy. (1979). |
| GOST 9378 | Roughness comparison specimens. (1975). |

\* Other republics assumed to be the same.

## 5.13 Specification on drawings

### 5.13.1 Surface roughness

The technical drawing is still in one form or another the most usual form of communication between the design engineer and the manufacturing engineer. Any ambiguity in this link is disastrous. For this reason many of the instructions from one to the other have been standardized. Of these the method of specifying surfaces and their geometric features is pertinent here.

The relevant ISO document is ISO 1302 (2001). The basic symbol for roughness is two legs of unequal length as shown in figure 5.54. This is called the tick system. The legs are 60° apart.

If the removal of material by machining is required a bar is added to the basic symbol as shown in figure 5.54(*b*). If the removal of material is not permitted a circle is added to the basic symbol as in figure 5.54(*c*).



**Figure 5.54** Basic symbols for surface.

There are some occasions when messages are implicit in the symbol but these have to be agreed. For example, the symbol used in figure 5.54(*c*) can be understood as meaning that the surface is to be left in the state resulting from a preceding manufacturing process. This applies to whichever method has been used. There are other occasions where more information rather than the somewhat negative kind given above is required. The basic symbol is intended to be able to cater for this eventuality. The longer arm has a flat attached to it as in figure 5.55. For just the surface roughness alone the symbols of figure 5.54 have the actual value of roughness added in the vee. Thus figure 5.56 shows a roughness *a* where the actual value *a* can be obtained by any production method. If it accompanies a completed triangle as in (*b*) it must be obtained by the removal of material by machining and (*c*) it must be obtained without the removal of metal.



**Figure 5.55** General symbol for special requirements.



**Figure 5.56** Basic symbol on drawing.

Note that if only one value of roughness is obtained it refers to the maximum permissible value of surface roughness.

To avoid confusion between the different units of the Imperial system and the metric system, sometimes a grade number system is used. For example, if the roughness parameter is the $R_a$ value it is advisable to use the grade number as shown in table 5.5 rather than the numerical value of $R_a$ required just in case the measurement system is confused.

**Table 5.5**

Roughness value $R_a$

| ($\mu$m) | ($\mu$in) | Roughness grade no |
|---|---|---|
| 50 | 2000 | N12 |
| 25 | 1000 | N11 |
| 12.5 | 500 | N10 |
| 6.3 | 250 | N9 |
| 3.2 | 125 | N8 |
| 1.6 | 63 | N7 |
| 0.8 | 32 | N6 |
| 0.4 | 16 | N5 |
| 0.2 | 8 | N4 |
| 0.1 | 4 | N3 |
| 0.05 | 2 | N2 |
| 0.025 | 1 | N1 |

Should the production process be specific to the part it has to be written in clear language above the bar on the symbol, as shown in figure 5.57.



**Figure 5.57** Symbol with process specified.

Also, any information regarding coatings should go on the space above the bar (figure 5.58).



**Figure 5.58** Coating roughness before and after coating operation.

Should the texture be required before and after plating it should be specified with two symbols, the first should show the roughness value of the surface before coating as $a_1$, the second symbol shows the roughness requirement $a_2$ after plating or coating.

The sampling length is placed under the bar as shown in figure 5.59. The lay direction is that of the predominant surface pattern determined by the process (figure 5.59). The symbols for the lay are varied and are shown in figure 5.60.

**Figure 5.59** Symbols for other surface details.

| Symbol | Interpretation |
|--------|----------------|
| = | Parallel to the plane of projection of the view in which the symbol is used  |
| ⊥ | Perpendicular to the plane of projection of the view in which the symbol is used  |
| X | Crossed in two slant directions relative to the plane of projection of the view in which the symbol is used  |
| M | Multi-directional  |
| C | Approximately circular relative to the centre of the surface to which the symbol is applied  |
| R | Approximately radial relative to the centre of the surface to which the symbol is applied  |

**Figure 5.60** Some lay pattern specifications.

Dimensional tolerance can be included within the symbol, usually to the left of the vee as shown in figure 5.61. The value should be put in the same units as those used in the rest of the drawing.



**Figure 5.61** Machining tolerance.

Taken as a whole figure 5.62 shows the positions around the symbol as just described.



**Figure 5.62** General texture symbolism: *a*, roughness; *b*, production process; *c*, sampling length; *d*, direction of lay; *e*, machine allowance; *f*, other roughness values.

### 5.13.1.1 *Indications generally—multiple symbols*

Figure 5.63 shows how the component might be completely covered for surface texture assessment.

### 5.13.1.2 *Reading the symbols*

Because there could be many surfaces to any component some convention has to be reached in how to specify them. Some examples of this are given in the following figures. Many of them are common sense, but for continuity in the system what follows agrees with ISO/R 129, which says that the symbol as well as the inscriptions should be orientated so that they may be read from the bottom or from the right-hand side of the drawing as seen in figure 5.63. This figure also shows where a leader line is used in those cases where difficult surfaces are being specified. Such leader lines should terminate in an arrow on the surface in question and may also branch. Figure 5.63 shows that symbols can be put anywhere providing they do not carry any indications of special surface characteristics.



**Figure 5.63** Ways of writing multiple position symbols.

Another general rule is that the symbol will point from outside of the material either to a line representing the surface or to an extension of the outline (figure 5.64).

In accordance with the general rules on dimensioning the symbol should only be used once per surface, preferably on the view carrying the dimension determining the size or position of the surface (figure 5.64).



**Figure 5.64** Positioning of symbols relative to dimension lines.

Should the same texture be wanted on all the surfaces of a part it is not necessary to put symbols on all the surfaces. One symbol only is sufficient providing that it is specified by a note near the title block, in the space devoted to general notes or following the part number on the drawing (figure 5.65).



**Figure 5.65** Texture is the same all over the part.

### 5.13.1.3   General points

An indication of roughness is only given when necessary for the functional working of the part, in the same way that tolerances or special processes are. Often the specification of texture is unnecessary if the surface is made by a process which ensures an acceptable finish.

### 5.13.1.4   Other points

Surface texture specification on drawings is one aspect of the typology of surfaces. What was an acceptable typology at one time may not be later on. For example, some modern surfaces are being made to better

finishes than the fixed surface grade numbering system in Table 5.5. In these circumstances it is wise to use the actual value in clear units of the system in use and to forego the grading system until updated grades are in the standards.

Also many surfaces are functionally important in such a way that conventional surface measurement parameters will not be sufficient, for example in the case where there is a specification for flaws or scratches or where areal properties are required. Again, a clear note to this effect should be near to the surface in question.

It can be seen therefore that the specification of surfaces on drawings presupposes that the surfaces can be adequately specified. As has been seen in chapter 1, this is not necessarily so. Until such specification has been firmed up it is wise to add an explanatory note in words or by a picture of what is required in addition to the drawing.

## 5.14 Summary

Surface metrology is, as its name suggests, primarily concerned with measuring surfaces. This means assigning a value, usually roughness, to a surface, making sure that the numerical value obtained can be repeated by someone else and, finally, ensuring that the surface parameter chosen is significant and reliable. This chapter has been concerned with just these issues.

Factors which have had to be considered include the following:

1. The setting down of the best procedure to be used to measure the surface. This aspect is considered in section 5.7.
2. Having an understanding of the nature of the measuring system and controlling the errors that can arise from the use of a particular instrument. This means setting down a careful calibration procedure. A great deal of emphasis has been given to this aspect in this chapter. Another aspect of this consideration is understanding 'method error', which may arise between different techniques for measuring the surface. In this case the word 'error' is wrong because there is no 'correct' value. It just means accepting that different types of instruments use a different physical basis for interrogating the surface with the result that different numerical values for the surface can be produced. These aspects are briefly outlined in sections 5.8 and 5.9.
3. It is important to choose the correct parameter with which to characterize the surface and to know how sensitive the parameter is to the numerical technique. Although not covered specifically in this chapter in great detail (see section 5.11), it has been considered in some detail in chapters 2 and 3.
4. Knowing the nature of the interaction of the surface with the instrument and the parameter estimation is important. In section 5.7 the first part of this issue was addressed and in section 5.11 the second. It should be stressed that, because the latter problem is very dependent on the parameter and the surface, the technique for assessing the variability has been given rather than listing results.
5. Another aspect of surface metrology has been included in this chapter, which is that of estimating errors in composite parameters that arise in cylindricity and other form deviations. This has meant that some space has been devoted to error propagation in the first part of the chapter (sections 5.1 to 5.6).
6. Statistical testing and methods of discriminating are an essential part of assessing the credibility of measured values. For this reason an outline has been given in section 5.7 with some examples.

This part of the chapter is not meant to provide an in-depth course of basic statistical methods. It is intended to be more of an illustration of the type of statistical test often used by the metrologist.

A final part of the chapter has been used to provide some standards information. It is impossible to provide a complete list because of the changing scene. However, some major examples have been given for instruments, calibration standards and drawing specifications. These standards are fundamental to the exchange of information and the traceability of measured values.

Perhaps the biggest challenge in the past few years has been that of calibrating SPM. Even though the $x$ and $y$ movements can be made traceable to international standards down to nanometres, the $z$ direction is more difficult. Unlike traditional engineering metrology instruments the $z$ variable is not necessarily geometric; tunnelling current and atomic forces are well known possibilities. As a result the interaction between the probe and surface is more complex than with conventional instruments. In the case of the AFM, for example, how much deflection of the probe does a calibrated unit of force produce? It is easily possible to contour on constant force by closed loop force control but what exactly is the force? The more exotic SPMs pose worse problems. It is not sufficient to name the instrument as 'metrological' just because movement of the probe in the three directions is traceable; the $z$ is indirect. It is not satisfactory to infer the force on an AFM cantilever by noting its deflection; there may be two or more sources of force or interacting fields, or a different force to that specific force being (hopefully) measured by the AFM. The stylus properties or geometry may be different from that expected. For true calibration a unit of force should be applied *in situ*. The same argument holds for other phenomena such as charge density and magnetic flux.

Another problem is that the probe has now to be specified not only in size and shape but also its intrinsic chemical and physical properties. Of course the probe properties have to be known in terms of hardness and conductivity in conventional instruments but not to the same extent. Properties of the stylus have had to be controlled to avoid distorting the signal, whereas now the probe properties have to be carefully matched to the atomic/molecular phenomenon being investigated.

The SPM instruments are in the same position as the SEM; some very impressive pictures are being produced with lots of detail in the $x$ and $y$ directions but is the $z$ meaningful? In some respects the situation is worse for SPMs rather than SEMS because the SPM is subject to the same requirement for dynamic calibration as the standard surface texture instrument; the relative movement of the probe and surface has mechanical constraints. For some reason these are rarely addressed in SPM performance.

Should there be any problem concerning the content or the availability of any national standard it is usually best to refer instead to the relevant international standard either by itself or via a compilation such as the ISO 1999 handbook.

Chapter 6, which follows, considers how surfaces are generated by the manufacturing process and what characteristics are important. It also deals in some depth with the use of the surface to control the process and to provide some degree of monitoring the machine tool condition.

## References

[1]    Hoffman D 1982 *Measurement Errors*, *Probability and Information Theory* vol 1 (New York: Wiley)
[2]    Dietrich C F 1973 *Uncertainty*, *Calibration and Probability* (London: Hilger)
[3]    Woschini E G 1977 Dynamics of measurement *J. Phys. E: Sci. Instrum.* **10** 1081
[4]    Davies O L (ed) *The Design and Analysis of Industrial Experiments* 2nd edn (London: Longman)
[5]    Thomas T R (ed) 1982 *Rough Surfaces* (Harlow: Longman)
[6]    Teague E C 1976 *NBS Note* 902 p 99
[7]    Witzke F W and Amstutz H N 1976 Sources of error in surface texture measurement *SME Paper* IQ76-585
[8]    Stevens R M D *et al* 2000 Carbon nanotubes as probes for atomic force in microspy *Nanotechnology* **11** 1–5
[9]    Shevyakov V, Lemeshko S and Roshin V 1998 Conductive SPM probes of base $T_i$ or W refractory compounds *Nanotechnology* **9** 352–355
[10]   Nguyen C V *et al* 2001 *Nanotechnology* **12** 363–367
[11]   Chetwynd D G 1998 low cost testing of profilometer stylus forces *6th IMEKO Symposium Tech University Wren Sept.*
[12]   Reason R E 1971 *The Measurement of Surface Texture* ed Wright Baker (London: Longman)
[13]   Whitehouse D J, Bowen D K, Chetwynd D G and Davies S T 1988 Micro calibration of surfaces *J. Phys. E: Sci. Instrum.* **21** 40–51
[14]   Whitehouse D J 1973 *Stylus Methods—Characterization of Solid Surfaces* ed Kane and Larobee (New York:Plenum) ch 3
[15]   Brand U and Hillman W 1995 Calibration of step height standards for nanotechnology using interference microscopy and stylus profilometry *Precision Engineering* **17** 22–32
[16]   Ewart Williams E 1954 *Applications of Interferometry* (London: Methuen Monographs) p 83

[17] Edwards H 1997 New method to estimate step heights in scanning-probe microscope images *Nanotechnology* **8** 6–9
[18] Whitehouse D J (partial ones)
[19] Reason R E 144
[20] Schneir J and McWaid T H 1994 Design of atomic force microscope with interferometric position control *J. Vac. Sci. Tech.* **12** 3561–3566
[21] de Groot P 2001
[22] Teague E C 1989 NIST Molecular measuring machine, project metrology and precision engineering design *J. Vac. Sci. Tech.* **7** 1898
[23] Teague E C and Evans C 1988 Tutorial Notes A.S.P.E. Raleigh, North Carolina p 99
[24] Hart M 1968 An angstrom rule *J. Phys. D: Appl. Phys.* **1** 1405–8
[25] Chetwynd D G, Siddon D P and Bowen D K 1988 X-ray interferometer calibration of micro displacement transducer *J. Phys. E: Sci. Instrum.* **16** 871
[26] Chetwynd D G, Krylova N O and Smith S T 1996 Metrological x-ray interferometry in the micrometre region met
[27] Chetwynd D G, Schwarzenberger and Bowen D K 1990 *Nanotechnology* **1** 14–26
[28] Chetwynd D G, Krylova N O, Bryanston Cross P J and Wong Z 1998 *Nanotechnology* **9** 125–132
[29] Raloff J 1985 *Sri. News 1* **28** 92
[30] Yaro H, Nakamura T and Matsuda J 1985 *Bull. NRL Metrol.* January p 36
[31] Barash V Y and Uspensky Y P 1969 *Measurement Technology* **2** 83 (in Russian)
[32] Barash V Y and Resnikov A L 1983 *Measurement Technology* **3** 44 (in Russian)
[33] Hasing J 1965 *Werkstattstechnic* **55** 380
[34] Hillman W, Kranz O and Eckol T 1984 *Wear* **97** 27
[35] Bryan J, Clouser R and Holland E 1967 *Am. Mach.* **612** 149
[36] Trumpold H and Frenzel C 2000 Cali Surf *Contract SMT* 4-CT97–2176
[37] Thwaites E G 1974 *Int. Symp. on Metrology. INSYMET* **74**, *Bratislava*
[38] Reason R E 1966 *Report on the measurement of roundness* Rank Organisation
[39] Evans C J, Hocken R J and Estler W T 1996 A multipoint method for spindle error motion measurement *Ann. CIRP* **42** 441
[40] Donaldson R R 1973 *Ann. CIRP* **21** 125–6
[41] Whitehouse D J 1976 Error separation techniques in surface metrology *J. Phys. E*: *Sci. Instrum.* **9** 531
[42] Spragg R C 1972 *Proc. Joint Meas. Conf.*, *NBS* p 137
[43] Iizuka K and Goto M 1974 *Proc. Int. Conf. Prod. Eng. (Tokyo)* pt 1, p 451
[44] Thwaite E G 1973 *Messtechnik* **10** 317
[45] ANSI, Standard axes of rotation, B89
[46] Tlusty J 1957 Systems and methods of testing machine tools *Microtechnic* **13** 162
[47] Zhang G X, Zhang Y H, Yang S M and Li Z 1997 A multipoint method for spindle error motion measurement *Ann. CIRP* **42** 441
[48] Kakino K, Yamamoto Y and Ishii N 1977 *Ann. CIRP* **25** 241
[49] Murthy T S R, Mallina C and Visvesqeran M E 1978 *Ann. CIRP* **27** 365
[50] Peters J and Van Herk P 1973 *Proc. 14th MTDR Conf. (Manchester)*
[51] Lukyanov V S and Lisenko V G 1981 Accuracy of assessment of 2D parameters by discrete and analogue analysis *Izmen. Tech.* **9** 5 (in Russian)
[52] Levin A I and Memirovskii G 1970 Establishment of the interval of independence between values of separate functions *Izmer. Noya Tech.* **3** 8 (in Russian)
[53] Skakala L 1976 Die Analyse des Messverfahrens zur Bestimmung des Rauheits Keennwertes Rž bie Betrachtung der Oberflache. als normales stationarer Prozess *Paper 27 IV Oberflachenkolloquium (Feb. 1976, Karl Marx Stadt, DDR)*
[54] Spragg D J and Whitehouse D J 1970 A new unified approach to surface metrology *Proc. IMechE* **184** 397–405
[55] Bendat J S 1958 *The Principles and Applications of Random Process Analysis* (New York: Wiley)
[56] Teague C T 1978 *Metrologia* **15** 39-44
[57] Shannon H B 1966 NEL specimens for surface texture calibration *NEL Rep.* p 231

# Chapter 6
# Surface metrology in manufacture

## 6.1 Introduction

The role of workpiece geometry in manufacture is shown below in figure 6.1 as a subset of the block diagram drawn in the first chapter (figure 1.1).



**Figure 6.1** Role of manufacture and measurement.

The present chapter shows that there is a two-way interaction between the surface and the manufacturing process. The first interaction is concerned with the nature of the geometric characteristics produced on the surface by the manufacturing process. This is considered in some detail in sections 6.2 to 6.6. Physical characteristics are examined only briefly in section 6.7. The other interaction concerns ways in which the surface roughness and form can be used to detect changes in the process and also; in some cases, the machine tool. Some ways of doing this are examined in section 6.8. It is shown that the Fourier kernel and its derivatives, such as the Wigner function, can be used to great effect.

In what follows the way in which most conventional machining processes and some less conventional processes affect the surface will be examined. In order to do this it is often necessary to describe the basic mechanisms of the process and sometimes to attempt to determine how these mechanisms are reflected in the surface characteristics. Consequently, although the process will be described, the investigation will only be carried out as far as is necessary to explore the surface implications. Some recent advances will be discussed. These include dry cutting and the implications of miniaturization.

For many comprehensive process details one of the standard references, such as the Metcut machinability handbooks, should be consulted.

## 6.2 Manufacturing processes

### 6.2.1 General

In order to transform raw material into a workpiece having the desired shape, size and surface quality, it has to be processed by some means. There are many different ways in which this transformation can be

Process — Roughness values (µm $R_a$)

50  25  12.5  6.3  3.2  1.6  0.8  0.4  0.2  0.1  0.05  0.025  0.0125

Flame cutting
Snagging
Sawing
Planing, shaping
Drilling
Chemical milling
Electro-discharge machining
Milling
Broaching
Reaming
Boring, turning
Barrel finishing
Electrolytic grinding
Roller burnishing
Grinding
Honing
Polishing
Lapping
Superfinishing
Sand casting
Hot rolling
Forging
Permanent mould casting
Investment casting
Extruding
Cold rolling, drawing
Die casting

key:

average application     less frequent application

**Figure 6.2** Typical roughness values produced by processes.

achieved. Each has its own particular advantages and disadvantages. Some workpieces are produced by one process and others by many. In this section the shape and roughness connotations will be concentrated on. How the roughness is produced will be one consideration, but as important will be what the roughness can reveal.

Figures 6.2 and 6.3 show the relative cost of different machining processes and the roughness produced according to BS 1134 and similar standards.



**Figure 6.3** Breakdown of typical machining times as a function of roughness.

The processes will be split up into a number of subdivisions:

1. Cutting with single or multiple tool tips—this includes turning, milling, broaching, planing.
2. Abrasive machining—this includes grinding, polishing, honing.
3. Physical and chemical machining—this includes electrochemical machining, electrodischarge machining, etc.
4. Forming, casting, extrusion.
5. Other macroscopic machining to include laser machining, high-power water jet.
6. Ultra-fine machining (nanomachining) including ion beam milling and energy beam machining.

Within these there is a considerable overlap, for example diamond turning is capable of producing very fine surfaces, yet it is included here in group 1. In fact the groups have been assembled according to generic likeness rather than magnitude of roughness, because it is the former that imparts the unique character of the process to the workpiece.

### 6.3 Cutting

#### 6.3.1 Turning

##### 6.3.1.1 General

This process is most common for generating the primary dimension of the part and, as the name implies, involves an axis of rotation somewhere in the generation. Figure 6.4 shows some of the turning modes.



**Figure 6.4** Typical turning configurations: (*a*) radial; (*b*) axial; (*c*) face.

Typical variables are cutting speed—workpiece peripheral speed relative to the tool axial feed — the advancement of the tool per revolution of the workpiece, the shape of the tool and the depth of cut of the tool into the workpiece material.

There are other very important aspects that are not shown on the diagrams but which contribute a considerable difference to the form and roughness. These include the absence or presence of coolant and, if present, its constitution and the method of supply, whether fluid, mist or drip, and so on. In addition to these effects is the effect of the machine tool itself. But in general the actual value of the roughness can be estimated at least crudely in terms of height and form from a knowledge of the basic process parameters, unlike the case for grinding.

As far as turning is concerned, it is rarely (except in the case of diamond turning) used for very fine finishes. It is a very useful method for removing stock (material) in order to produce the basic size and shape. As a general rule, the surface roughness tends to be too rough to be used in very critical applications in which high stresses could be detrimental, but nevertheless there are many applications where turning (single-point machining) is used because of the speed and efficiency of the process. There is a form of turning, "hard turning" where the tool is very hard like CBN and it is being used to replace grinding. This is because of the extra flexibility and control available in turning processes.

The undisputed advantage of single point cutting is the flexibility of the tool position and orientation relative to the workpiece. This flexibility enables a wide variety of shapes to be machined.

This wide selection of shapes is not open to processes using more complicated tools such as in milling and grinding.

### 6.3.1.2 Finish machining

Although single-point machining is often used to remove enough material to get the workpiece down to nominal size, it can be used as a finishing process. It then goes by the name of finish machining. $R_a$ values of roughness are typically in the range of 0.5 to 1 $\mu$m (20-40 $\mu$in).

There are two distinct types of finish produced by a cutting operation involving a single tooth. These are (i) roughness due to the primary cutting edge and (ii) roughness due to the secondary cutting edge. The first usually refers to surface broaching and form turning. The second is conventional turning with a round-nosed tool. Figure 6.5 shows the position of the cutting edges on a typical tool.

### 6.3.1.3 Effect of tool geometry—theoretical surface finish—secondary cutting edge [1]

In its simplest case the tool could be thought of as a triangle (figure 6.6). If $f$ is the feed and $d$ is the depth of cut then

$$R_t = d \qquad R_a = d/4. \tag{6.1}$$



**Figure 6.5** The two cutting edges with respect to the tool.



**Figure 6.6** Surface produced by angular tool.

In this case the surface texture is independent of the tool feed. Equation (6.1) does not involve $f$, the feed. The triangular tip is hardly practical. What is more often used is a curved tip of radius $R$ and, in its simplest form, is as shown in figure 6.7. In this case the roughness, at least in principle, is a function of the feed. Assume that all the cutting takes place on the radiused part of the tool as shown in the figure. The roughness can be given by

$$R_t = R - \sqrt{R^2 - f^2/4} = R(1 - \sqrt{1 - f^2/4R^2} \tag{6.2}$$

for

$$f \ll R \qquad R_{\mathrm{l}} \sim f^2/8R \qquad\qquad (6.3)$$



**Figure 6.7** Surface produced by curved tool.

From equation (6.3) it can be seen that for a given tool tip radius the surface roughness can be improved quite dramatically, simply by reducing the feed. In fact the roughness is very much more sensitive to feed than it is to tip radius. Before the $R_{\mathrm{a}}$ can be found, the mean line has to be established. This is obviously trivial in the case of the triangular tip—it is half way down the depth of cut. In the case of the curved tip it is $R_{\mathrm{t}}/8$ from the bottom of the trough.

Using the mean line the $R_{\mathrm{a}}$ value is approximately

$$0.032 f^2/R \qquad\qquad (6.4)$$

which is about one-quarter of the $R_{\mathrm{t}}$ value, as can be seen from, equation (6.3). This is the same ratio as for the triangular tip and can be assumed for most single-shaped tool tips.

For the lowest ratio of $R_{\mathrm{t}}/R_{\mathrm{a}}$ the surface has to be a square wave surface; then the ratio is 1:1. For a sine wave surface the ratio is $R_{\mathrm{t}}/R_{\mathrm{a}} = \pi$.

Most often a round-nosed tool has straight flanks and, depending on the depth of cut, these flanks enter into the calculation of the theoretical surface roughness. Three conditions apply (figure 6.8):

(1) round nose only producing the surface (case I);
(2) round nose and straight edge in primary cutting position (case II);
(3) round noise and straight edge in primary cutting position and straight edge in secondary cutting position (case III).



**Figure 6.8** Different regions of tool causing roughness.

Thus in terms of the $R_t$ value and the radius $R$

$$\text{case I: } \frac{R_t}{R} = \frac{1}{8}\left(\frac{f}{R}\right)^2 \tag{6.5}$$

which is the same as equation (6.2). As more of the flank is included, the roughness ratio $R_t/R$ [2] becomes more involved, as seen in equation (6.6) below. Thus

$$\text{case II: } \frac{f}{R} = \left[2\frac{R_t}{R} - \left(\frac{R_t}{R}\right)^2\right]^{1/2} + \sin\beta + \left(\frac{R_t}{R} - 1 + \cos\beta\right)\cot\beta \tag{6.6}$$

which can be extended to the more general case of

$$\frac{R_t}{R} = \frac{f/R}{(\tan\alpha + \cot\beta)} - \frac{\cos(\pi/4 - \alpha/2 - \beta/2)}{\sin(\pi/4 + \alpha/2 + \beta/2)} \tag{6.7}$$

which reduces, when $f/R \sim 1$, to

$$\frac{R_t}{R} = \frac{f/R}{(\tan\alpha + \cot\beta)}. \tag{6.8}$$

In figures 6.8 and 6.10 $\alpha$ is the side cutting angle and $\beta$ is the end cutting angle. In practice the only part of the tool used is as shown in figure 6.8 and is suitable for values of $f/R < 3$ (figure 6.9).

In this treatment $R_t$ is a theoretical value of roughness determined by the geometrical aspects of the tool tip. In practice the texture will not be equal to this for a number of reasons — built-up edge included. A definition of efficiency of cutting has been used in the past [3].

The $R_a$ value is also obtainable from the simple geometry in other cases. If the triangular tip is not regular as in figure 6.10 (*a*) then



**Figure 6.9** Theoretical roughness with tool shape and feed variables.

(*a*) triangular

$$R_a = \frac{f}{4(\tan \alpha + \cot \beta)}$$

$$(6.9)$$

(*b*) curved

$$R_a \sim \frac{0.443 f^2}{R} \left( \frac{3\alpha' - \beta}{\beta} \right) (\alpha' + \beta)$$

where $f > 2R \tan \beta$ and $\alpha'$ is given by $\cos \alpha' = H/R$; that is $\alpha = \cos^{-1}(H/R)$ where $H$ is the distance of the centre line from the centre of radius $R$.



**Figure 6.10** Tool shape relative to workpiece.

From these equations:

1. Plane approach angle $\alpha$ has an effect on the surface roughness.
2. Plane trail angle $\beta$ should be kept low because the surface roughness value increases in proportion.
3. In general, the lower the feed and the bigger the radius the better the surface roughness, unless chatter sets in.
4. Rake and clearance angles affect the cutting action.

Checking the formulae against practical tests is very difficult, if not impossible, for conventional turning. For example, in condition 1 the tool tends to break when it engages with the surface. For radiused tools the agreement can be checked more readily because alignment of the tool is not too important.

Secondary cutting edge roughness is present in what is called conventional turning. This has been considered above. The secondary edge is separated from the primary edge by the nose radius as seen in figure 6.8. The use of the secondary edge to generate the roughness introduces several complications. The most important ones are:

1. The geometry of the tool at its nose is replicated in the surface at feed mark intervals.
2. There is some uncertainty of the geometry of the cut at the trailing edge because the chip thickness reduce gradually to a small value.
3. The metal at the trailing edge of the tool is subjected to unusually high normal stress and tends to flow to the side in order to relieve this stress. This can produce a furrow which can contribute to the roughness, especially in the case of a soft ductile metal. In this case the tool profile is not properly replicated in the surface.

Surfaces generated by the secondary cutting edge are in general more complicated than those produced by a primary cutting edge. The surface roughness quoted is always taken as the maximum value. For the secondary edge roughness is in the axial direction.

### 6.3.1.4 Primary cutting edge finish

This is measured in the circumferential direction as opposed to the axial direction in secondary edge roughness. Primary edge roughness is important not only in its own right but also for interpreting the wear on the clearance of turning tools employing a nose radius.

The surface profiles of the primary cutting edge and the secondary cutting edge are compared in figure 6.11.



**Figure 6.11**

The roughness produced in both cases is strongly dependent on the cutting speed. The roughness for low speeds is very poor, particularly that produced by the primary edge, whereas the roughness obtained by use of a secondary cutting edge approaches the theoretical value at high cutting speeds (figure 6.12). This figure shows that there is quick convergence to the theoretical value when the cutting speed increases. Figure 6.13 shows a comparison between theory and practice.

There are a number of factors that can influence the roughness and introduce extra factors into the evaluation; one is fracture roughness and another is built-up edge on the tool.

### 6.3.1.5 Fracture roughness [3]

When cutting steel at very slow speeds a special type of roughness is frequently observed owing to subsurface fracture. When the metal is cut at about room temperature the chip is found to fracture periodically.



**Figure 6.12** Effect of cutting speed on surface finish.

**Figure 6.13** Comparison of theory and practice—function of speed.

This crack tends to follow, not unsurprisingly, the path of maximum shear stress. When the cracks run completely across the chip the surface has the appearance almost of a moiré fringe. The shiny areas correspond to where the cut has taken place, while the dull patches are the uncut regions. As the speed increases the cracks tend not to go across the whole chip. Above a medium to high cutting speed of about 50 m min $^{-1}$ this pattern of roughness disappears.

### 6.3.1.6 Built-up edge (BUE)

As the cutting speed is increased the friction between the chip and the tool increases. At some point this will be large enough to cause a shear fracture in the chip in the vicinity of the tool face—a built-up edge (BUE) will be formed. There is no BUE at very low speeds of about 0.5 m min$^{-1}$ because the temperature on the face of the chip is not sufficient to cause the chip surface to behave in a ductile manner.

As the speed (and hence temperature) increases the chip metal in contact with the chip face becomes ductile and the plastic flow can cause welds between the chip and the tool. The extra plastic flow on the chip face causes hardening and a further increase in force which causes an even stronger adhesion. When the bonding force between the chip and the tool becomes greater than the shear strength of the metal in the chip the BUE forms.

As the BUE grows forward it will usually grow downward, causing in some cases the finished surface to be undercut. For a small BUE the direction of the resultant force puts the BUE into compression, which makes for a stable configuration. If the BUE becomes large the force loads it as a cantilever and eventually the BUE comes off. The gradual growth and rapid decay of the size of the BUE causes a sawtooth roughness on the surface which is very characteristic of a BUE. The way in which the BUE forms can be seen in figure 6.14. The BUE and its effects are therefore symptomatic of the middle range of cutting speeds above 50 m min$^{-1}$.

The BUE has been studied extensively by Nakayama and colleagues [4, 5], who give ways of eliminating it thus:

1. Increase the cutting speed.
2. Make the metals less ductile. Brittle materials find it difficult to get into the plastic mode and hence the BUE does not get the chance to form because of fracture.
3. Increase the rake angle.
4. Have a fluid present. This can eliminate BUE at low speeds.

**Figure 6.14** Influence of speed on chip formation.

The entire range of surface roughness produced by a primary cutting edge is dependent on the temperatures of the face and flank of the cutting tool. Figure 6.15 [2] shows the boundary between BUE as a function of cutting speed and the chip-tool interface temperature.

The BUE is partly responsible for the dependence of the surface roughness on cutting speed. The influence is shown practically in figure 6.16, which also shows the rather less dependence of feed.



**Figure 6.15** Built-up edge as function of speed and chip tool temperature.



**Figure 6.16** Variation of roughness ($R_t + R_a$) with cutting speed and feed (after Shaw).

### 6.3.1.7 Other surface roughness effects in finish machining

#### (a) Effect of minimum undeformed chip thickness

Sokolski [6] suggested that there is a nominal thickness below which a chip will not be formed and the tool begins to rub the surface. This criterion is called the minimum undeformed chip thickness. It depends on the

tip radius, the cutting speed and the stiffness of the system. When this concept is applied to secondary edge roughness it is found that a small triangular portion of material which should have been removed is left. This constitutes an additional roughness called a Spanzipfel [7].

The theoretical surface roughness has to be modified accordingly, to take into account this remnant material. Thus

$$R_t = \frac{f^2}{8R} + \frac{t_m}{2}\left(1 + \frac{Rt_m}{2}\right)$$

(6.10)

where $t_m$ is minimum chip thickness.

The second term represents the Spanzipfel. In some cases, for example for dry cutting, this can amount to a sizeable fraction of about 25% of the actual roughness. Figure 6.17 shows the formation of a Spanzipfel.



**Figure 6.17** Formation of Spanzipfels on a surface.

In most estimates of the Spanzipfel the value of the minimum chip thickness can be taken to be about 1 $\mu$m.

### (b) Physical properties of the workpiece

The hardness of the workpiece plays an important role in finish machining. If the work is too soft, side flow will result which is unacceptably high. This means that it takes too high a speed to get rid of the BUE and to produce a thermally softened layer on the tool face. Harder materials do not suffer in this way. For this reason the finish-machining performance of low-carbon steels can be improved by cold drawing.

The ductility of the workpiece is also important. Additives like sulphur and lead decrease chip ductility and consequently improve surface roughness by providing a smaller BUE and a smaller tendency towards side flow.

### (c) Other factors

The sharpness of the cutting edge is important when the tool is new and it also influences the way in which the tool wears. Carbide tools which are very sharp give up to 20% better finish than those ground normally.

The cutting fluid also has an effect on roughness. It has been reported that the presence of a cutting fluid, if chemically active, can reduce roughness by a factor of up to 4 at low speeds. But the effectiveness of such fluids drops off dramatically with speed and can cause an increase in tool wear.

Material inhomogeneity is also very important. The presence of microvoids, inclusions and the generation of alternative slip planes can cause non-typical surface marks. All these produce unwanted effects. In the case of slip steps these are produced because deformation is easier in certain glide planes. The form of these steps will depend on the orientation of the crystalline structure (see figure 6.18).

### 6.3.1.8  Tool wear

Obviously tool wear affects the surface roughness dramatically. For a given tool geometry there are three areas in which wear can take place, as shown in figure 6.19. Here region $A$ is the tip itself, region $B$ is the tool

The dimple formation model



The scratching of an inclusion



**Figure 6.18** Alternative slip of work material causing slip steps.



**Figure 6.19** Tool wear regions.



**Figure 6.20** Regions on surface profile affected by tool wear areas.

flank and region $C$ is the tool top. In regions $A$ and $B$ scratches and fracture are most likely to occur. In region $C$ the principal type of wear is due to cratering, which in turn affects the possibility of the formation of BUE. Each of these regions affect different regions of the profile as seen by an instrument.

On the profile cusp shown in figure 6.20 the central region is affected most by $A$ and then $C$; the outer regions of the cusp are affected more by flank wear. The interesting point here is that there is a very close correspondence between the geometrical features imposed on the tool by wear and the geometry imparted by the tool on to the surface. So, for example, if a scar is at the tip of the tool as shown in figure 6.21 it will produce a profile as in figure 6.21(*b*) where the scar and tool shape form a negative in the profile. The other important point to notice is that the tool shape and scar is repeated at the feed spacing, this has an important bearing later on when the problem of using the surface roughness to monitor the tool wear is considered; it is exactly the opposite of the situation described here.

**Figure 6.21** Transfer of tool scar to surface profile.

Other points need to be mentioned. One is that the effect of BUE and cratering tends to be non-repetitive from adjacent cusps because it is a continuous process rather than discrete, as in the tool scar case, so that the situation can change from one feed profile to another. This has the effect of producing a more random type of mark within each cut than a simple scar. Because of this difference in character the effects can, in principle, be separated by the use of Fourier analysis.

Figure 6.22 shows a relationship between tool life curves based on the wear land criterion and the surface roughness criterion [8].



**Figure 6.22** Tool life curves based on surface roughness criterion (---) and wear land criterion (—) (after Semmler) on steel.

The turning process efficiency and effectiveness is determined largely by the tool wear, so many different ways have been tried to measure or to monitor the tool for impending failure [9].

It should be pointed out that, although the tool wear is important, what is required is monitoring of the whole process, if possible the machine tool and, last but not least, the quality of the workpiece itself. It has already been shown in chapter 2 that the workpiece finish can be used to do just this. However, the sensor has to measure the workpiece remotely.

One factor crucial to satisfactory operation and yet not easy to eradicate in practice is operator error. In the reference above failure in turning and drilling can be attributed in 50% of cases to the operation. Failure due to other factors such as failure of the monitoring system are 15% or less.

Sensors for tool condition monitoring should fulfil certain requirements. Some of these are listed below:

(a) Measurement should be as near to the tool as possible.
(b) There should be no reduction in the static and dynamic stiffness of the machine tool.
(c) No reduction of the working space or curtailing of any of the cutting parameters is allowed.
(d) Easy maintenance is essential.
(e) The method should be insensitive to dirt, chips etc.
(f) Reliable signal transmission is essential from the sensor to the controller.

From the above it can be seen that the requirements are so formidable that no single sensor technique is completely suitable. It may be, however, that optical methods based on the workpiece surface could be used, especially in the dry cutting mode. Ultrasonics have been used to examine the surface by coupling through the coolant to the surface [10] as also is the use of acoustic emission to detect tool workpiece contact.

As usual when researchers use optical or ultrasonic sensors the beam is focused on the surface in an attempt to mimic the stylus. This is fundamentally wrong. Focusing the source onto the detector via the surface enables spectral methods to be applied, which are much more suitable [11].

Tool breakage detection using cutting force monitoring is well known and not discussed here. Monitoring torque via the armature current of the spindle motor is also a possibility. Tool breakage has a catastrophic effect on the texture, which could be monitored optically.

### 6.3.1.9 Chip formation and turning mechanisms

Chip formation depends on a number of factors. One is the tool material (i.e. CBN or diamond). Another main factor is the workpiece material. Obviously the surface finish gets worse if tool wear (which depends on cutting distance) occurs.

Figure 6.23 shows that the surface finish in $R_a$ varies linearly with distance machined [12].

The figure shows that there is a difference in roughness if high CBN content is used rather than low which shows itself as having a different intercept rather than a different slope.

It is interesting to note that the surface finish is related to the surface finish of the chip being produced. If the chip has segments—like a saw tooth—along its surface the same average segment spacing is also evident on the surface but in a smaller degree. The spacing of the irregularities in the chip profile is dependent



**Figure 6.23** Effect of cutting distance on $R_a$.

on cutting speed. Also decreasing the rake angle from $-27°$ to $-10°$ and decreasing the depth of cut produces shorter segment spacings on the chip and hence the surface.

### 6.3.2 Diamond turning

This is quite different from that of ordinary turning in the sense that a complicated tool shape is often involved as well as there being a complicated cutting mechanism [13]. Diamond turning is frequently used for obtaining very fine surfaces on soft materials such as aluminium and copper where abrasive methods like grinding are unsuitable. Furthermore, diamond turning is carried out at very high cutting speeds so that problems associated with BUE do not arise. Under these circumstances it should be possible to achieve the true theoretical surface roughness. On the other hand, the cutting tool can have two facets, one that cuts and one that can be used to burnish the machined surface (figure 6.24). When the burnishing follows the cutting, the surface finish can be very fine indeed. Typical roughnesses for ordinary diamond turning may be about 25–50nm $R_a$ but with burnishing they can be considerably reduced.



**Figure 6.24** Dual purpose of cutting tool.

However, the smoothing over of the tool mark by plastic flow caused by the following edge can introduce detrimental physical properties into the surface such as residual stress, so the burnishing can only be used in non-critical stress applications or where relaxation caused by the presence of stress could deform the workpiece.

The amount of burnishing also has an upper limit. If it is larger than necessary then the height levels of the crystal grains will differ. This is due to the different crystallographic orientations of the grains.

After burnishing the material springs back. The amount of spring-back is determined by the forces on the clearance face and by Young's modulus, which again is dependent on the grain orientation. Different orientations will give different spring-back and thus a difference in height levels of the surface (figure 6.18 and 6.25). If the area of contact between the tool and the workpiece at the clearance faces increases, the forces and consequently the difference in height level also increase.



**Figure 6.25** Effect of crystal orientation on recovered surface.

The essential difference between conventional turning and single-point diamond turning is the relative magnitude of the influence of some of the factors. For example, the influence of anisotropy of the material upon surface roughness may be neglected in conventional cutting but should not be in single-point diamond turning.

Also the surface roughness will not be an exact replica of the shape of the cutting tool. This is due in part to the burnishing effect on the clearance face. Other effects such as plastic side flow and the Spanzipfel should not be neglected.

The use of a diamond tool for cutting has certain advantages. The extreme hardness allows the fabrication of very sharp cutting edges that can produce mirror-like finishes after a single pass, and the strength and toughness permits the machining of difficult materials which could hardly be worked with conventional tools.

However, the properties of diamond are very complex and how they are influenced by the cutting process is by no means clear. As a result there is no specific specification for diamond tools as yet [14].

There are a number of types of tool, namely single crystal and polycrystalline, which have different properties when cutting. These have been investigated by Hingle [13].

Diamond tools are well known for cutting very fine surfaces but are restricted to some extent by the affinity of diamond to steel. Consequently, other materials have to be used. One such is tungsten carbide and the field of application is the machining of micromoulds or similar small objects to high accuracy.



**Figure 6.26** Roughness and cutting speed.

Figure 6.26 shows how the $R_a$ varies with cutting speed and hardness. The upper line represents 40 $R_c$ and the lower 60 $R_c$. Although it is hardly influenced by cutting speed the hardness does seem important. In the graph the roughness values seem high even for steel microcutting. Good microcutting of steel with tungsten carbide gets to one tenth of the values shown above.

The achievable surface roughness $R_y$ for microturning is given [15].

$$R_y = \frac{f^2}{8r} + \frac{h_{\min}}{2}\left(1 + \frac{r.h_{\min}}{f^2}\right)$$
(6.11)

where $f$ is feed, $r$ is tool radius and $h_{\min}$ is the minimum cutting depth to produce chips.

A rough guide to $h_{\min}$ is $h_{\min} = f^2/2r$, the simple spherometer formula for the parabolic mark left by the tool. In practice, the parabola is modified in microcutting to be nearer to a saw tooth because of elastic recovery of some material in microcutting [16].

### 6.3.3  Milling and broaching

#### 6.3.3.1  General

At first glance milling and broaching appear to be quite different processes. This is not altogether so because surface broaching is similar to peripheral milling having a cutter radius of infinity. In fact surface broaching is being increasingly used to replace milling as a stock removal process and, at the same time, milling is being used for finer surfaces.

Milling processes are many and varied and as a result the cutters used and their impression on the surface can be complex. However, there are two general types: peripheral milling and face milling. Other applications such as the milling of slots, forms and helices are essentially variants of the basic methods and each is specialized to such an extent that general surface considerations are not possible.

In peripheral milling the cutter teeth are machined to give cutting edges on the periphery. These may be gashed axially or, more often, spirally as seen in figure 6.27. There are two rake angles, one for the spiral and the other for the axial rake. It can easily be shown that a large effective rake can be produced on a milling cutter without unduly weakening the tooth by having a large radial rake. This is important because, as in turning, the cutting force per unit length of cutting edge reduces rapidly with an increase in rake. Another advantage of using a helical rather than a purely radial cutter is a more even distribution of cutting force, which tends to preclude chatter and hence waviness on the surface. A disadvantage is the end thrust component of the cutting force.



**Figure 6.27**  Angles of importance in milling: $\alpha$, spiral angle; $\beta$, axial rake.

There are two methods of removing metal, one known as upcut milling and the other as downcut or climb milling. The essential features are shown in figure 6.28.

Although upcut milling is most often used it does rely on very efficient clamping of the part on the bed, otherwise chatter is easily generated on the surface. Climb milling does not need this but does require an adequate backlash eliminator [16]. Another difference is that the two techniques produce different-shaped chips.

Typical pictorial representations of upcut and downcut milling are shown in figure 6.28. Peripheral milled surfaces in general are not critical or are finished with another process.



**Figure 6.28**  Milling cuts: (*a*) climb milling; (*b*) upcut milling.

The other main type of milling—face milling—is capable of producing finer surfaces, so more detail as an example will be given here [18].

The factors affecting the surface texture in face milling are as follows:

1. selection of cutter geometry;
2. the accuracy of grinding of the selected angles on the cutter teeth;
3. the setting of the teeth relative to the cutter body;
4. the alignment of the machine spindle to the worktable.

A typical roughness pattern is shown in figure 6.29.



**Figure 6.29** Surface appearance of face milling, single tooth.

To consider these points in the way described by Dickenson [18] a selection of cutter angles is shown in figure 6.30. The various tool angles on a face-milling cutter fall into two categories: (i) the cutter geometry which directly affects the surface roughness and (ii) those angles which indirectly affect the surface roughness by influencing the cutting actions. The radial clearance angle and the corner angle fall into the first of these two categories.

It seems from tests that the surface roughness values deteriorate as the radial clearance angle is increased, so to keep a good surface finish this should be kept small. If a corner tip radius is used it should be approximately 10 $\mu$m.

### 6.3.3.2 Roughness on the surface

#### (a) Setting of teeth relative to cutter body

Anyone who has examined face-milled parts realizes that the dominant mark on the surface is invariably produced by one or two teeth being non-coplanar with the rest. This has the effect of producing a periodicity corresponding to feed per revolution superimposed on the feed/tooth periodicity. In theory, the surface roughness should be determined from the latter; in practice it is more likely to be the former.

The fact is that, if one tooth is only a very small amount lower than the other teeth, it completely swamps out the other cutter marks. As an example [18], on a 12-tooth cutter with a radial clearance angle of 1/4° operating at a feed of 0.1mm/tooth, one tooth would need only to be proud by 0.5 $\mu$m for the surface roughness to change from 0.1 $\mu$m to 1.5 $\mu$m! This illustrates the sensitivity of the setting. Also the fact that the dominant periodicity is changed from feed/tooth to feed/revolution means that the cut-off wavelength of the instrument has to be chosen with care. In the example above, a cut-off of 2.5 mm would be more than adequate and 0.8 mm cut-off acceptable if none of the teeth were proud, but if one were, the 0.8 mm cut-off would give disastrously low readings of roughness. This is because the dominant periodicity would lie outside the cut-off (i.e. 1.2 mm), which means that the true roughness value is attenuated.

**Figure 6.30** Views of tool with relevant angles: (*a*) plan view; (*b*) front; (*c*) side view.

In practice there should be no problem because the inspector should examine the surface prior to measuring and determine the suitable cut-off relative to the dominant periodicity as laid down in the standards.

*(b) The pattern of roughness on the surface*
This is determined to a great extent by the alignment of the spindle to its worktable. Any misalignment will lead to an increase in the surface roughness values for cutters which are ground without a corner radius. Obviously if the cutters have a radius it is not so important.

The pattern, assuming that no one tooth is proud, can take three forms:

(1) where all the leading edges of the cutters generate the mark—called all front-cut;
(2) where all the backing edges of the cutters generate the mark—called all back-cut;
(3) where there is equal cutting front and back.

Ideally the third possibility is the one that should be produced but in practice it is difficult to avoid back cutting because of the run-out of the spindle and the elastic recovery of the material.

As with the 'proud' tooth problem the main issue in face-milling alignment problems arises with the selection of the cut-off. Unbelievably, the problem is worst when the front cutting and backcutting are more or less balanced, as shown in figure 6.31.

The spatial interaction of the forward cutting and backward cutting can produce some very strange periodic effects if only one profile is taken. Shown in figure 6.31 are two tracks. The waveform of the periodicity is quite different in both cases, as also is the variation in the spacing of the dominant features. This means that it is quite difficult to decide on the cut-off to use. Basically the only method is to take a profile graph and ensure that at least five or so maxima lie within one sampling length. The difficulty arises because in such a

roughness one profile graph cannot adequately represent the surface. The variation of $R_a$ with position across the surface is quite staggering when it is considered that the surface is essentially very deterministic and well behaved (figure 6.31).

It is not unusual for 40% variation to be found in the $R_a$ value on a regular face-milled surface, owing to the positioning of the profile track relative to the centre position of the cutter path. Thus at least three



**Figure 6.31** Multiple tooth milling—different profiles, same surface.

measurements on face-milled surfaces have to be taken. Note that if $R_q$ rather than $R_a$ or peak measurements were measured the variation would be much less because the phase effects of the different components which cause the variations would not contribute.

*(i) Effect of feed on surface roughness*
An increase in feed tends to increase surface roughness, so the feed should be as small as possible. However, the purpose of the milling operation is usually to remove metal as quickly as possible and so feeds tend to be high. As usual the two requirements go against each other.

*(c) Other factors in milling*
Lobing in radial peripheral milling caused by chatter is often due to the fact that all teeth on the cutter have the same directional orientation, which can in certain circumstances set up chatter. It has been shown, however, that some improvement in stability can be achieved if the cutters have non-uniform pitch as illustrated in figure 6.32 [19].



**Figure 6.32** Chatter-suppressing cutter.

*(i) Other types of milling*

Ion milling is something of a misnomer and will be dealt with in superfine finishing.

### 6.3.3.3   Theoretical milling finish

Because milling is not normally regarded as a finishing process, relatively little work has been carried out on attempting to predict the surface finish from the process parameters — unlike in turning, only very early work has been reported [20].

The $R_t$ value of the surface roughness has been given as

$$R_t = \frac{f^2}{8[R \pm (fn/\pi)]} \tag{6.12}$$

where $f$ is the feed per tooth, $R$ is the cutter radius and $n$ is the number of teeth on the cutter. The plus sign in the denominator corresponds to upmilling while the negative sign corresponds to downmilling. Equation (6.12) can be in good agreement with practice if the spindle of the milling machine is held to a very low value of run-out. It represents only a very simple modification of the theoretical surface roughness for turning (equation (6.3)). It seems probable that much more work needs to be done on this in the future with the renewed interest in milling.

### 6.3.4   Dry cutting

### 6.3.4.1   General

Although the use of cooling lubricants has been accepted in the past there are signs that ways are being sought to eliminate them. The reasons are mainly economic—their cost of procurement and disposal.

Adoption of dry cutting would help the use of optical methods for monitoring the process via the surface because the possibility of mist would be removed. Mist is more degrading than chips or debris because solids can be ignored in systems involving diffraction rather than imaging. Diffraction occurs when the surface is used as a mirror rather than an image. On the other hand ultrasonics benefit from having a coolant because it can be used as a coupling to the surface. Coolant has small attenuation relative to air (which reduces ultrasonic signals greatly).

### 6.3.4.2   Mechanisms

The main function of the coolant fluid is to reduce the effect of heat generation by reducing friction between the tool and the workpiece and to eliminate the exertion of unacceptable influences on the structure [21].

Another useful attribute of the coolant is its ability to help the removal of chips produced by machining. Where is dry cutting likely to be possible? The answer is when the temperatures generated during cutting are not too high.

One criterion for acceptability is that cutting times, tool wear and workpiece equality for wet cutting have to be equalled or bettered by dry cutting before dry cutting can be considered as an option.

It appears that this criterion is satisfied for cast material such as cast iron. This has short chips and low cutting temperatures brought about by the low friction coefficient. Low friction is probably due to the presence of embodied graphite in the workpiece material.

It is also believed that dry cutting has been responsible for producing lower thermal shocks but this has not yet been adequately authenticated! High temperatures invariably affect chip formation so, with dry cutting where there is a high temperature, some control over the chip formation is necessary. This is usually achieved by using cutting inserts with chip shaping grooves.

### 6.3.4.3  Pseudo dry cutting

There are ways of machining which lie somewhere between wet and dry cutting. One such method uses a minimal quantity of lubricant designated (MQL) [21]. This situation arises when coolant is blown off the surface; the coolant could have been necessary for an earlier process. The problem is that of deciding the amount of coolant still on the surface. In practice, if dry cutting is to be used it should apply to all operations in the process rather than just one. MQL operation has the capability of stopping chips and debris from sticking to the tool or workpiece, which sometimes happens in dry cutting.

Another hybrid technique uses cutting tools that have been coated to nanometre thickness with carbide/metal or solid lubricant/metal [23].

By providing numerous alternative nanolayers of hard and tough, hard and hard, or other combinations of desirable physical properties it is possible to enhance the tool characteristics, either from a wear point of view or that of thermal properties. Current trends suggest that single material tools are being replaced by complex layers of materials designed for specific applications. Thin coatings of about 5 $\mu$m have been used for some years. Coatings of TiC, TiN and Al$_2$0$_3$ were put on tools by various means usually involving sputtering or similar deposition methods. Nanolayers, however, allow a much more flexible approach involving combinations of carbides, nitrides and oxides for hard/hard bonds (e.g. B$_4$C/SiC) and hard/soft combinations (e.g. carbide/metal). It is also possible with nanolayers to enhance the chemical resistance of tools to corrosion. Often, because of the reduced friction between the tool and workpiece, the surface finish becomes more consistent. It is also possible to improve the precision of the process by dry action. One example is in dry EDM machining [23].

It is inevitable that the use of coolant liquids (CL) will reduce in time partly because of cost but also due to the presence of undesirable elements such as chlorine and lead. A great deal of research is being carried out to replace wet cutting with dry or pseudo dry processes.

## 6.4  Abrasive processes

### 6.4.1  General

In the previous section single-point and multiple-point cutting have been considered from the point of view of surface generation. In general these processes are used for the removal of bulk material to get the workpiece down to its intended dimension. Single-point cutting can be used as a finishing process, particularly if the tool is a diamond. The advantage of having one tool is that it can be positioned more or less anywhere in space and so can produce many shapes. An alternative way to produce a finish is by means of abrasive processes. The basic difference is that there are multiple random grains which carry out the machining process rather than one fixed one (or a few in the case of multiple-tooth milling). Consequently, it is more meaningful to consider *average* grain behaviour rather than the individual behaviour of the grains. Put in terms of turning, it means that there are many tools cutting with negative rake at a very small feed (which is equivalent to the effective grain width of the abrasive!). This ensures a much finer finish in itself. When coupled with the averaging effect of the whole wheel profile the resultant roughness is very low. Also, whereas in turning and milling the cutting is continuous, in most abrasive processes it is not. The abrasive action is usually analysed in terms of three mechanisms—cutting, rubbing and ploughing—and, more importantly, the angle the grain makes with the surface is negative rather than positive as in turning. This negative rake angle makes the probability of burnishing and ploughing more prevalent than in straightforward cutting. There are many factors in the abrasive process which influence roughness and are not considered here, such as hardness etc.

Figure 6.33 shows a situation in which a grain is forming a chip. It also shows the side ploughing [24]. During the material removal process the grains wear in the same way that the cutting tool does in turning, but in this case the abrasive usually fractures or is torn off. Again, there are some similarities between the deterministic cutting processes like turning and the abrasive ones, yet there are many more differences.

**Figure 6.33** Abrasive grit: reaction on surface and grain imprint.

There are three main abrasive processes: grinding, honing and lapping. The main difference between these is that in honing, although the grains are held rigidly, the speeds are low and reciprocating movement is often used. Superfinishing is similar except that a slight axial vibration is superimposed. Lapping is again slightly different in that the grains are loose, they are contained in a fluid and are therefore free to move. The type of surface produced by these methods is easy to identify, as shown in figures 6.34 and 6.35. It is in the areal view that the differences become obvious.



**Figure 6.34** Surface roughness (from DIN 4761): 1, peripheral grinding; 2, face grinding; 3, honing; 4, lapping; (*a*) and (*b*), process variations.

Figure 6.36 shows the conventional grain sizes and the surface roughness produced by the three abrasive processes [25].

The grain size in lapping is smaller than in the other processes as seen from figure 6.36. It seems that the relationship between grain size and roughness is about 200:1 and the relation between $R_z$ and $R_a$ is about 6:1, which is to be expected for a random process having a Gaussian height distribution; it corresponds to the extreme values being $-3\sigma$ to $+3\sigma$.

**Figure 6.35** Abrasive processes, same $R_z$ roughness of 4 $\mu$m: 1, grinding; 2, honing; 3, lapping.



**Figure 6.36** Graph showing range of grain size and workpiece roughness: (1), grinding; (2), honing; (3), lapping.

Energy, force and contact times all relate to the average areal roughness formed. Although they link directly to the surface roughness they also closely relate to the 'unit machining event' concept.

Figures 6.37–6.40 show how these parameters differ for various types of grinding, honing and lapping. Figure 6.41 gives a clear indication of the difference between process times and contact times for the three processes. The fixed grain methods have repetitive process times but lapping does not. This implies that there is a further random element that should be taken into consideration for lapping. Although the contact versus process times are not the main factors influencing the surface geometry, the times greatly influence temperature and subsurface effects.

It is interesting to note here that the very nature of abrasive processes allows a very close relationship to be made between the statistics of the abrasive element and that of the texture. In grinding, for example, it is

**Figure 6.37** Grain contact times as function of normal force: 1, OD grinding; 2, reciprocating grinding; 3, ID grinding; 4, face grinding and creep feed grinding; 5, honing; 6, lapping.



**Figure 6.38** Energy and average area of cut per grain: 1, lapping; 2, face grinding and creep feed grinding; 3, honing; 4, cylinder grinding and reciprocating grinding.

reasonable to expect that the envelope of the highest grains on the wheel should transfer directly to the surface geometry—as a negative. This is reasonably true in creep feed grinding where the machine is very rigid, but not usually in other forms of grinding.

In the case of bonded abrasives, forces can be normal and tangential and produce bending moments and shear stresses.

The cross-sectional area of cut perpendicular to the cutting-speed vector arises from the normal force on the grain and the workpiece properties. See figures 6.39, 6.40 and 6.41 for aspects of the mechanism. The power and energy per grain are proportional to the tangential force and the cutting speed. The cutting speeds in grinding (10-200 ms$^{-1}$) are one or two orders of magnitude higher than in honing (1–5 ms$^{-1}$). Therefore there is a particularly high energy input in grinding which results in very high temperatures compared with honing and lapping. Figure 6.42 shows the temperature build-up.

**Figure 6.39** Average cross-sectional area of cut and normal force per grain. Fields of application: 1, lapping; 2, face grinding and creep feed grinding; 3, honing; 4, cylindrical grinding and reciprocating grinding.



**Figure 6.40** Normal force per unit area for abrasive processes. Fields of application: 1, lapping; 2, face grinding and creep feed grinding; 3, honing; 4, cylindrical grinding and reciprocating grinding (23132).

The phases passed through during the contact cycle are as follows;

1. There is elastic deformation and the abrasive grain generates heat by friction along the workpiece surface.
2. With further penetration elastic and plastic deformation results. The friction between the grain and the work material increases, which results in the possibility of internal friction and the intensification of thermal stresses.
3. Elastic and plastic deformation and actual chip formation result.

In the case of loose abrasives the grains penetrate up to 5–15% of their size. These produce high surface pressures, which generate high elastic and plastic deformation. This can result in severe work hardening

**Figure 6.41** Process and process times for abrasive processes. Grain contact times or engagement times per grain $t_c$ and process times $t_p$: (*a*) peripheral grinding; (*b*) face grinding; (*c*) honing; (*d*) surface lapping.



**Figure 6.42** Thermal grain load in grinding as a function of contact time (according to Steffens): *a*, cutting edge before engagement; *b*, ploughing; *c*, chip formation; *d*, cutting edge after engagement.

and even the destruction of the abrasive grains. There is always a mixture of rolling grains and cutting grains [26–29].

### 6.4.2 Types of grinding

Whereas in most of the other manufacturing processes such as milling, turning, EDM, etc, the roughness and form of the resulting product may be important parameters, they are rarely if ever the dominant or most important parameter. Metal removal efficiency is often the most important factor. However, in grinding the situation is reversed. Grinding is usually regarded as a finishing process, that is one in which the bulk of

extraneous material has already been removed. It is the grinding process that is expected to finish the part to its intended size, shape and roughness. Very often the form and roughness requirements represent the most expensive and difficult criteria to satisfy. As in most processes there are a number of variants and grinding is no exception.

The principal ones in grinding are as follows: cylindrical external and internal grinding shown in figure 6.43(*a*) and (*b*) are similar to peripheral milling in the mode of operation of the wheel (milling cutter); vertical spindle grinding, as in cup grinding, is similar to face milling (figure 6.43(*d*)). The differences between these methods are merely relative differences in work speed and length of cut. Also, there is plunge grinding in which there is no axial feed across the component and cross-feed grinding in which an axial feed is present. Horizontal spindle reciprocating grinding is yet another—with or without cross-feed (figure 6.43(*c*)). Perhaps one of the most important is centreless grinding (figure 6.44) in which the workpiece (nominally cylindrical) is not held between centres, as in cylindrical grinding, but is allowed to float.



**Figure 6.43** Variants of grinding: (*a*) external cylindrical; (*b*) internal cylindrical; (*c*) reciprocating face; (*d*) vertical spindle (cup).



**Figure 6.44** Centreless grinding.

### 6.4.3 Comments on grinding

Recapitulating, whereas in the other cutting processes the number of cutting points of the tool is either one or a few, in the grinding process it is many. Furthermore, the cutting or abrasive grains are randomly distributed in the grip of a bonding agent. This makes the grinding process very difficult to analyse in the same way as is possible with deterministic cutting processes such as turning. In these processes bounds can be set on the value of surface roughness that could be predicted by using simple, yet reasonable, assumptions for the process. This is not possible for grinding, however, because of the random characteristic, the very small chips produced and the speed at which they are produced relative to the conventional processes.

### 6.4.4 Nature of the grinding process

#### 6.4.4.1 General

Grinding wheels are made up of a large number of grains held by a bonding agent. The abrasive grain can be aluminium oxide, silicon carbide, cubic boron nitride or diamond (either natural or synthetic). The relative merits of these materials or the bonds will not be of primary concern in this context except inasmuch as they affect the form or texture. The features of the wheel structure which are important in determining the resultant surface roughness are the density of grains/mm$^2$, the height distribution of the grains, the elasticity of the grains within the bond and the fracture characteristics of the bond and grain. Other associated issues include wheel dressing and spark-out. Because each of these has to be taken with the other parameters of importance such as cutting speed, workpiece speed, feed rate, it is more meaningful to discuss roughness and form in a number of specific contexts rather than to attempt to generalize the process. This, from the point of view of roughness and form, is at present impossible.

#### 6.4.4.2 Factorial experiment

Following Shaw and co-workers [4,5], the variety of parameters that could influence the surface are legion. The problem is to establish which of the possibilities are most likely to be important. This is where the use of factorial design of industrial experiments is useful. In a classic experiment in horizontal reciprocating grinding the question posed was: which grinding parameters were the most important for roughness? The results of this experiment are included here as a good example of how a structured statistical experiment can help to understand the manufacturing process. More details of the technique and this experiment are given in section 5.7.5.

*(a) Wheel characteristics*
These include the diameter $D$, grain type, grain spacing, grain size, tip radius, wheel grade and elastic modulus $E$, structure number, bond type, dressing method and wheel balance.

*(b) Workpiece characteristics*
The work diameter $D_w$ and workpiece constitution including elasticity $E_w$ are important.

*(c) Machine characteristics*
Spindle and table stiffness, damping vibration characteristics and vibration isolation should be considered.

*(d) Operating conditions*
The wheel speed ($V$), workpiece speed ($v$) and wheel depth of cut ($d$) can be useful and whether upgrinding or downgrinding (as in up- and downmilling) are being used.

In the experiment only two of these extensive possibilities were chosen: grain size and wheel depth of cut. The high and low values of the two parameters are given in table 6.1. $R_a$ values for the four combinations of grain size and downfeed $d$ are given in table 6.2.

**Table 6.1**

|  | Symbol | High value | Low value |
|---|---|---|---|
| Nominal grain diameter | $g$ | No 24: 0.03 in diameter | No 60: 0.01 in diameter |
| Wheel depth of cut | $d$ | 0.0004 in | 0.0002 in |

**Table 6.2**

| Grain Size | Depth of cut (in) | Designation | $R_a$ ($\mu$in) |
|---|---|---|---|
| No 24 | 0.0004 | $gd$ | 32.5 |
| No 24 | 0.0002 | $g$ | 26.5 |
| No 60 | 0.0004 | $d$ | 23.6 |
| No 60 | 0.0002 | (1) | 20.6 |

The direct effect of grain size in this limited experiment is

$$\pm \frac{1}{2^n}[(gd - d) + (g - 1)] = 3.7 \ \mu in.$$

The direct effect of depth of cut is

$$\frac{1}{2^n}[(gd - g) + (d - 1)] = 2.25 \ \mu in.$$

These results indicate that the grain size is about 60% more important than the depth of cut, so the obvious parameter to consider is the grain size rather than the cut.

This example is simplistic in the sense that a large number of interactions between the various parameters are possible, as explained in the earlier section on statistical methods. Whether or not these individual experiments are meaningful is neither here nor there; the fact is that in complex processes such as grinding these statistical methods need to be used to home in on the essential quality factors. Once these key parameters have been identified, separate experiments need to be carried out in as much detail as possible.

The factorial exercise illustrates two things. First, that the surface roughness is a suitable workpiece feature to test for with the statistical design of experiments. Second, that the depth of cut is not as important as may have been thought.

### 6.4.5   Centreless grinding (figure 6.44)

#### 6.4.5.1   General

In order to isolate the features of this very important type of grinding that are most likely to influence roughness and roundness, it is useful to be reminded of the possible factors. The factors can be roughly divided into three categories. A, B and C, as in figure 6.45 [31].

A. The system variables (which can be assumed to be inherent):

(1) rigidity of wheel, regulating wheel spindles and housings, and the rigidity of the work rest holder;
(2) specification of the grinding wheel (which determines the contact stiffness and removal rate);
(3) grinding wheel speed—usually fixed;

**Figure 6.45** Transfer of variables to workpiece geometrical fidelity.

(4) vibration of two wheels and blade;
(5) diameters of two wheels;
(6) workpiece characteristics, hardness, geometry, length, diameter, etc.

B. Process variables (easily changeable during grinding; see figure 6.44):

(1) the regulating wheel speed (controls the work speed, the through feed rate and the amplitude and frequency of regulating wheel vibration);
(2) metal removal rate—infeed per workpiece revolution;
(3) dressing conditions for grinding and regulating wheels.

C. Set-up variables—changed every set-up for a given workpiece:

(1) work support blade—material, thickness and top-face angle;
(2) work centre height above grinding wheel centre;
(3) regulating wheel through feed angle;
(4) diamond dresser angle with respect to the regulating wheel axis;
(5) dressing diamond offset;
(6) longitudinal profile of the regulating wheel dressing bar.

In view of the number of variables they have been grouped under the following headings:

(a) metal removal process
(b) system vibration
(c) workpiece characteristics
(d) workpiece stability.

The transfer from the machine tool to the workpiece is thereby simplified.

### 6.4.5.2 Important parameters for roughness and roundness

Among the variables listed above, there are some relating to the set-up of the grinder; in particular, they are concerned with the regulating wheel and the work support blade. These two machine elements have no parallel in other examples of grinding and it is not surprising therefore that they are particularly important where

roughness and roundness are concerned. Roundness especially is the feature of surface geometry which is most affected in centreless grinding.

It is most important that the regulating wheel is accurately trued. Conventionally, this is achieved by a single-point diamond dresser applied two or three times at a depth of 20 $\mu$m and then at a light cut. Run-out of the regulating wheel periphery at the truing point is typically about 8 $\mu$m, although that of the spindle itself is considerably better, as low as 0.2 $\mu$m. Alternatively the regulating wheel is ground by plunge grinding. This reduces the run-out and profile out-of-straightness to one-tenth of the conventional method. It is assumed that the slide accuracy of the truing grinder has the same order of accuracy as the machine slide or better.

Truing the regulating wheel in this way reduces its wear and reduces the surface roughness of the work-piece—say, cylinder rollers from 0.3 $\mu$m $R_a$ to 0.1 $\mu$m $R_a$. Even more dramatic is the improvement in roundness. A typical value of roundness error of 2 $\mu$m can be reduced to 0.2 $\mu$m when the regulating wheel is properly trued. All this is achieved at the cost of having a much stiffer system and a much narrower stable region [30–32].

Simulations have shown that, irrespective of the work support angle, when the workpiece centre height above the grinding wheel centre is zero, the system improvement as far as roundness is concerned is marginal, and that for workpieces with roundness values to be less than 0.5 $\mu$m, made up of high frequencies as well as lobing, it is often better to worry more about the type of wheel specified and the dressing conditioning than the aspects of machine stability.

This is not to say that chatter is unimportant. In grinding it can spoil the capability of the machine tool very quickly and measures have to be incorporated to reduce it wherever possible.

### 6.4.5.3   Roundness considerations (figure 6.46)

Furukawa *et al* [32] give simple rules for the occurrence of waviness conditions. They state that the system would be unstable if

$$\alpha = n'_e \beta \quad \text{and} \quad \pi\beta = (n_e - 1)\beta \tag{6.13}$$

where $n'_e$ and $n_e$ are even numbers. For optimum roundness it was recommended that $\beta = 7°$. Furukawa *et al* added the further recommendation that

$$\alpha = (n'_e - 1)\beta. \tag{6.14}$$

Some conclusions reached about different configurations are the following:

1. Lobed shapes relating to odd harmonics of order below the 11th (i.e. 3, 5, etc) are better removed by using a relatively large value of $\beta$ (see figure 6.47).
2. Lobed shapes relating to even harmonics of order below the 10th (i.e. 2 etc) are better removed with a small angle of $\beta$.
3. Other errors generated, which could include even or odd harmonics, are dependent on the magnitude of the infeed.

Luckily attributable roundness errors often tend to be smaller than would be expected from simulation experiments [33].

Simulations and experimental studies especially by Rowe have been used extensively to investigate roundness [34].

Note:
It is interesting to note that the angle $\beta = 7°$ is about the same as that needed to get maximum fidelity for the vee-block method of measuring roundness using an offset probe.

(a)

Grinding
wheel

Work
piece

Control
wheel

(b)

Grinding
wheel

Workpiece

Control
wheel

$\alpha$

$\beta$

Workplate

A convenient waviness condition

(c)

Grinding
wheel

$\partial_2$

$\beta$

$\alpha$

$\phi$

Control wheel

$\alpha^1$

$\partial_1$

Movement normal to wheel surface

(d)

$\partial_1$

$\beta$

$\alpha$

$\partial_2$

$\alpha$

$\beta$

The spreading and averaging of roundness errors

**Figure 6.46** Centreless grinding, roundness considerations.

**Figure 6.47** Harmonic changes with angle.

Various conflicting views have been put forward, especially for very highly accurate work on the stiffness requirements. There have been some statements that introducing compliance may even help the accuracy of the workpiece and reduce the roundness. This seems to be suspect, although it may well be that controlled damping would have a beneficial effect, especially in reducing regenerative chatter. The endless search for stiffer designs leads to the conclusion that, from the point of view of expense, much more effort should be put into alternative directions. An obvious one is the use of optimized damping [35].

Considerable expertise by Miyashita [32] has been expended over the years to carry the process of cylindrical grinding to the important position in manufacture, especially in hard brittle materials, that it holds today.

### 6.4.6 Cylindrical grinding

#### 6.4.6.1 Spark-out

In cylindrical plunge grinding the important parameter as far as the roughness is concerned is the speed ratio $q$, where this is the quotient of cutting speed $V_s$ and the workpiece peripheral speed $V_{ft}$ at contact (figure 6.48).



**Figure 6.48** Spark-out as a function of speed.

The effects of the wheel speed ratio $q$ on the workpiece roughness have to take two ranges into account:

(i) the infeed per workpiece revolution is greater than the roughness (creep feed grinding); (ii) the infeed per workpiece revolution is less than the roughness (reciprocating grinding).

Note that in normal grinding the $q$ value lies between 20 and 100 and in high-speed grinding, creep feed grinding is characterized by $q$ values of between 1000 and 20 000. Obviously the sign of the $q$ value depends on whether upgrinding or downgrinding is being used (figure 6.48).

Some approximate formulae have been derived which relate the surface roughness value to the $q$ value at least in form [36]. For creep feed grinding it has been postulated that the roughness is given by

$$R_z = K^* |q|^{-2/3} \tag{6.15}$$

where $K^*$ involves the factors associated with grinding conditions such as the contact length and another called the 'characteristic grinding value', which loosely means the length of the grinding wheel circumference which has to grind a longitudinal element of the workpiece circumference until the value of $R_z$ has been covered in the direction of the plunge motion. For reciprocal grinding the formula is modified somewhat to give

$$R_z = K^* (1 + 1/|q|)^{-2/5}. \tag{6.16}$$

It is interesting to see how the value of $R_z$ expressed as a ratio of its maximum value (over the $q$ range) varies with $q$ going from 0 to $\infty$. The radial infeed and peripheral speed of the workpiece have to be considered for these border cases.

More will be said in the summary of this chapter about the usefulness or otherwise of such weakly correlated equations as (6.15) and (6.16).

For very small $q$ the wheel is motionless in comparison with the workpiece, and the radial infeed per workpiece revolution approaches zero. The individual cutting profiles superpose themselves onto every profile of the workpiece circumference on nearly one workpiece revolution. The resulting transverse roughness $R_z$ of the workpiece corresponds to 'spark-out' roughness.

For very high $q$ the workpiece is nearly motionless relative to the grinding wheel. By superposition of the cutting profiles of the grinding wheel the spark-out roughness as a transverse or axial value is produced on a short section of the workpiece circumference.

It was this technique that was used to generate the first anisotropic random calibration specimens by PTB in the 1960s on flat specimens, so the value of the roughness approaches a low value for both boundaries of $q$.

Figure 6.48 shows the general picture. The spark-out roughness results in both limits (or at least approximately). In practice no wheel is capable of covering the range of $q$ shown so that there is usually a break at the top in the range 20–200. This break seems to be apparent for aluminium oxide wheels as well as cubic boron nitride (CBN) wheels [37].

Whilst touching on the subject of CBN it should be noted that surface roughness has been one of the main difficulties associated with the use of CBN in grinding when used in place of conventional abrasives. For a wheel in the as-trued condition the surface roughness generated could well be acceptable, but sharpening the wheel either by progressive use or by a separate treatment tends to result in much rougher surfaces than are desirable—presumably because the wheel itself has become so much rougher?

One way to overcome this is to use fine-grit CBN but this causes lower grinding ratios (ratio of metal removed to wheel wear) and so is counterproductive from the point of view of economy, especially when internally grinding. An alternative method is to use polycrystalline instead of monocrystalline CBN because the grain fracture occurs on a much smaller scale.

Surfaces ground with CBN, although tending to be rougher than conventional material, do tend to be more uniform in character, which is often a useful feature if a subsequent finishing process is to be used such as honing or polishing.

Form or roundness error produced by CBN grinding is less than for conventional materials because the wheel wear is so much less. This means that the grinding profile can be more controlled, although there may be difficulties in the reprofiling of the wheel [37].

### 6.4.6.2 Elastic effects

Probably the biggest problem associated with grinding rather than simple cutting is the difficulty of dealing with deflections. What is the real depth of cut? What is the real elastic modulus of the wheel, it being composite and difficult to assess? Is the relationship between elastic modulus and hardness significant? What is true is that local elastic deformation of the wheel and the workpiece also influences the surface roughness.

Unfortunately these questions need to be asked, yet convincing answers are not altogether forthcoming. However, exposure to the problem and its possible influence on the finish is helpful.

It is not the brief of this book to go into much detail on the relative theories but a good, plausible comment is as follows.

The reason why totally convincing answers are not forthcoming is because the machining process can be dealt with on two levels: one from the result point of view, the other phenomenological. It is only by considering the latter that advances will be made. Yet, for the practising engineer it is understandable that it is only the results which count. Production technology cannot yet bridge the gap between the enormous possible variants.

Mechanisms have been described by Shaw to illustrate the effect on surface roughness of the active cutting grains. However, the elastic effects of the grains themselves are also likely to be important; the effect is to reduce the peak-to-valley height of the surface and thereby improve the roughness. Improvement in the surface roughness due to elastic recovery of the workpiece will occur through two mechanisms [3]: (i) a change of cross-section of the grooves and (ii) a levelling of the grooves. The rotation of grains caused by a result of normal and tangential forces has a complicated effect on the surface roughness, as seen in figure 6.49.



**Figure 6.49** Grain rotation.

In general negative and positive grain rotation causes positive radial deflection. This effectively reduces the depth of cut of the higher asperities on the wheel and improves the finish. Because of lack of evidence of this mechanism many efforts have been made to simulate the effect of the wheel profile on the surface produced. Obviously this is nothing like as straightforward as in the case of single-point cutting. The equivalent of the theoretical finish worked out deterministically is now random and open to a number of possibilities. One way has been to profile the grinding wheel over many close tracks (~200) and develop a composite wheel surface made from the envelope of the maxima [38]. When compared with the corresponding profile of the surface produced it is a great deal larger. An $R_t$ value of the wheel was about 2 $\mu$m greater than that of the surface— a considerable difference, probably caused by grain deflection [39].

Of the four components of deflection in the contact zone three occur in or near the wheel surface, while only one applies to the workpiece. The nature and constitution of the wheel is therefore most important, but this does not imply that the workpiece deflection can be totally ignored.

Workpiece hardness is an important parameter of many finished components. As the workpiece hardness increases, the surface roughness improves [39] owing to greater elastic recovery and the increased normal forces which in turn produce larger grain deflection.

Wheel hardness has an effect on the surface roughness but it seems that there is some confusion as to the nature of the effect. Some workers have argued that the harder the wheel, the better the finish and vice versa [40, 41].

It seems probable, however, that the real factor again is the elastic deformation length, which is determined more by the wheel elasticity than the hardness [42].

### 6.4.6.3 Texture generated in grinding

Surface roughness is affected by the depth of cut but not perhaps to the degree expected. It is difficult to measure in practice; what is input as a cut may not be due to the elastic deformations mentioned above. This is typical of the difficulty of tying down the roughness to any simple relationship with the process parameters in abrasive processes, as shown later. Relating depth of cut to roughness is a typical problem.

In practice it is found that the effect is not very significant and can almost be considered constant for many grinding applications. Variations in surface roughness in constant-force experiments in plunge grinding by Hahn and Lindsay [40] were attributed to wheel wear (figure 6.50). For forces approaching zero there appears to be limiting roughness which is thought to be that roughness corresponding to the naturally wearing surface of the abrasive grain.



**Figure 6.50** Hahn experiment on texture produced by grinding.

It should be pointed out that much lower values of surface roughness are obtainable by spark-out. This has been examined in some detail by Nakayama and Shaw [4] who assume that each pass of spark-out multiplies the effective number of cutting grains in proportion. The problem in the grinding process is that of determining which grains are actually cutting. Assuming that this is known and is of value $C$, Nakayama and Shaw produced a relationship between the $R_t$ value of the surface and grinding parameters. If $C = A(R_t - h_0)$, where $A$ and $h_0$ are constants for a given wheel dressing, and the height of surface roughness is $R_T$ then

$$R_t = \frac{h_0}{2}\left[1 + \left(\frac{2v}{VAh_0^2\sqrt{2\rho D}}\right)\right]$$ (6.17)

where $v$ is work speed, $V$ is wheel speed, $D$ is wheel diameter and $\rho$ is grain radius. For a typical wheel (60H) $h_0$ is 1.5 $\mu$m and $A$ is $5 \times 10^{-2}$ mm$^{-3}$. In the equation, ploughing and built-up edge are ignored so that this is very much a rule-of-thumb assessment (~30%). A useful, but perhaps somewhat simplistic, approach for spark-out can then be applied, realizing from their analysis that $CR_t$ = constant. So, having estimated how $C$ varies with $h$ (which is done using a sooted glass plate), the $C$ versus $R_t$ curve in figure 6.51 can be used.

A typical reduction of finish is 30–50% for four passes. However, the reduction is not as a rule as small as predicted because many grains hit air rather than metal on the second and subsequent passes.

The wheel dressing plays an important part in the surface roughness [43]. Poor dressing, for example, can introduce a spiral wave on the surface which is anything but random and can in fact dominate the roughness. Hahn and Lindsay [40] have shown quite clearly that if the dressing lead is very low the surface roughness is considerably reduced, which seems obvious, or at least the other possibility is that high lead gives a high surface roughness (figure 6.52).

The effect of random patterned passes to reduce the possibility of periodicities has not been properly explored. This technique would certainly help to target the effect of wheel dressing to its primary purpose of regenerating the grinding wheel and not dominating the surface roughness.



**Figure 6.51** Effect on surface roughness of spark-out.



**Figure 6.52** Effect on surface roughness of dressing lead.

### 6.4.6.4 Chatter

The other main effect producing periodicities of the roughness—as well as the roundness as mentioned earlier—is chatter. It has already been indicated that this is due to a number of parameters of grinding including the process, machine tool and workpiece. A typical block diagram is shown in figure 6.53 [44]. Unlike the basic texture, chatter marks are more easily identified with specific machining parameters because of the more distinctive nature of the geometry produced.

It has been pointed out that chatter characterization is very different from that of cutting for the following reasons:

1. There is usually a slow growing rate of amplitudes.
2. Complicated waveforms are generated on the surface.
3. Grinding wheel regenerative chatter predominates.
4. There is a progressively decreasing tendency of vibration frequencies with chatter growth.



**Figure 6.53** Block diagram showing influences producing chatter in grinding.

Methods have been proposed which can reduce the incidence of chatter in centreless grinding. In particular [45], a diagrammatical coincidence method, the 'rounding effect criterion diagram', has been developed which not only defines chatter-free conditions analytically but also the effect of operational settings. It also describes the influence of the dynamic characteristics of the system on the rounding process (the removal of lobes) and has enabled guidelines for the design of machines and the selection of wheels to be defined. From the diagram a procedure has been proposed for the set-up of the grinding operation for stable grinding. This has also been achieved to some degree in earlier experiments by Trmal and Kaliszer [46] and in other significant papers by them in the grinding area. However, it is only recently that Japanese papers have included design criteria for the machine tool itself in the overall process of grinding and these criteria have been based on empirical data, not theoretical!

### 6.4.6.5 Other types of grinding

Before considering in more detail what can be expected as to the actual shape of the surface roughness as opposed to just its size as measured by $R_a$ or $R_q$, consider briefly some other types of grinding and its influence on surface roughness.

The basic problem with the grinding process and surface roughness is that if the process is automated, which is now possible for conventional grinding, the surface roughness can be relatively poor when compared with high-quality lapping. But lapping is difficult to automate, so a balance between good texture and efficiency is needed. One such method is called low-temperature precision grinding (LPG). This yields good surface roughness and no thermal damage. It has been used effectively for workpieces where a high load-carrying capacity is needed. Under these circumstances a bearing ratio of about 90% and not 50% at the mean line of the surface roughness is required.

Such a change in bearing ratio from conventional grinding has been managed by Loladze [47] and is shown in figure 6.54. In this method low grinding speeds were used, producing temperatures of less than 370°C.

Almost the opposite effect is true for heavy-duty grinding, such as the type used in steel-snagging foundry work and deburring. The surface roughness seems to be dominated more by the wheel type rather than the grinding speed, which is reminiscent of the results of Nakayama and Shaw [4].

Typical differences are shown in figure 6.55 when using four identical grinding conditions, such as a maximum load of 1500 N, a wheel speed of 4000 rpm and a traverse speed of 60 mm s$^{-1}$. In the figure, wheel type SR is a high-purity sintered aluminium wheel used for heavy duty purposes with a grit size of 12 and ZS represents a wheel of 25% zirconium and 75% aluminium [48].



**Figure 6.54** Stable grinding and its effect on bearing ratio curve.



**Figure 6.55** Influence of wheel type on bearing ratio curve.

### 6.4.7 General comments on grinding

Grinding as a process is very complicated, so it is to be expected that there is no fixed relationship between the grinding parameters and surface texture. Furthermore, because of the relatively recent growth in instrumentation for measuring roughness, many of the classical experiments have limited use because only $R_a$ or $R_z$ was measured. Nevertheless, some general results have been obtained. The implication of this will be discussed in the summary at the end of the chapter. It is hardly encouraging! The variables of importance could conceivably be listed as grinding time, wheel hardness grade, grit size and depth of cut. Some indication of their importance has been shown in the factorial experiment. For cylindrical grinding, for example [47], there has been an attempt to link the parameters in a non-dimensional way albeit in a very specialized experiment with alumina wheels.

It is postulated that

$$R_a = \left( \frac{v}{V} \frac{x}{W} \frac{1}{rn_c} \frac{d}{D_e} \right)^{n/2} \varphi(a) \tag{6.18}$$

where $v$ is the workpiece speed, $V$ the wheel speed, $x$ the traverse rate, $W$ the wheel width, $d$ the effective depth of cut, $D$ the wheel diameter, $D_w$ the workpiece diameter ($D_e = 1/D_w \pm 1/D$) and $n_c$ the number of active grinding grits per unit surface area of the wheel; $r$ is the average width-depth ratio of a groove produced on the surface. Values of $n_c$ in equation (6.18) above vary from 1.2 to 1.4, although in high-speed grinding the value of $n$ is approximately 3.0.

The ratios involved depend on the metal removal rates. At low rates high grinding temperatures result and the wheel becomes glazed (type I mechanism). At increased rates (type II) the grits acquire flats as a result of attritional wear and some material from the workpiece adheres to the wheel. At high removal rates (type III) much more metal adheres and the wheel becomes loaded. Types I, II and III are the basic mechanisms of macrogrinding.

For type I and II grinding (at least in plunge grinding) the roughness increases with time at least for a short period. Figure 6.56 illustrates the opposite of what happens in creep feed grinding and figure 6.57 shows the effect on roundness.

For both types of mechanism the roughness is nearly independent of the wheel hardness and tends to increase with grit size. To give a rough guide, if the grit increases in size by 25% the roughness increases by a slightly lower percentage.



**Figure 6.56** Cylindrical grinding — effect of grinding time.

**Figure 6.57** Roundness as a function of workpiece resolution — creep feed.

In type III grinding, for a given table speed and wheel speed, the finer the grit and the higher the wheel hardness the better the texture. This type of grinding happens when the number of grits being exposed is about equal to the number being lost in the process of machining. Of the three grinding mechanisms the general conclusion is that type III is the most efficient from the point of view of obtaining a satisfactory surface roughness at the highest removal rate and yet producing the minimum of deleterious effects on the surface (e.g. workpiece burn).

### 6.4.7.2 Slow grinding

Invariably associated with grinding is high speed working with its associated high temperatures. However, this is not always the case. In gear making, for example, a new method of finishing which improves the surface on the resultant gear and allows a small degree of gear correction has emerged. This is called 'shave grinding' which is a slightly contentious name because alternative names are 'gear honing' and 'power honing'. It is safe to follow Brinksmeir's nomenclature which is 'shave grinding' and is the result of a rolling and sliding movement between the geared workpiece surface and abrasive grains bonded in a resin, vitrified or metal tool matrix shaped similar to an internal gear [49].

The tool is shaped like a gear. The material removal is a result of the kinematically imposed relative movement (sliding) between the grains of the tool and workpiece. The surface texture produced has a complex lay because the curvature of the workpiece and the tool are continually changing. Also because the grinding wheel has a relatively low wear resistance the topography changes rapidly [50].

Because of the relative shapes of the tool and workpiece very low grinding speeds are produced of about 0.5 m/sec.

### 6.4.8 Nanogrinding

Part of this section is produced in chapter 8 on nanotechnology.

Traditional precision finishing of hard and brittle materials such as a single-crystal wafer and magnetic heads consists of grinding, lapping and polishing (figure 6.58). These depend too much on skilled labour. The process is based on traditional optical production technology. What is happening is that the traditional lapping and polishing methods are being replaced by nanogrinding. The idea is to produce damage-free surfaces

equal to those produced by lapping and polishing on brittle materials in roughness and 'form'. Also the technique is meant to be suitable for mass production. The conventional production of optical form is based on pressure copying where the predominant parameters are (i) grain size, (ii) uniformity of grains and (iii) machining pressure. Material removal $M$ is based on Preston's formula $M = pvk$, where $p$ is pressure, $v$ is velocity and $k$ is a constant. This refers to loose abrasive grains such as are found in lapping. There are three basic modes of material removal in this regime, as shown in figure 6.59: (i) brittle mode machining; (ii) microcrack machining; (iii) ductile mode machining. In traditional optical manufacturing, all these are combined to produce the required figure and texture. The predominant parameter is the load applied on every grain.



**Figure 6.58** Optical manufacture, computer-controlled three-process finish.



**Figure 6.59** Mechanisms of microgrinding.

What the figure describes is the amount of material removed. What it does not give is the amount simply moved. The regime where plastic flow occurs around the grain is where ductile machining fits.

The critical question is how can ductile machining, whether grinding or turning, be guaranteed? The answer to this in part is that it has already been achieved in diamond turning. Optical parts have been made using diamond turning; the problem now is to carry the techniques to the grinding process.

The removal of brittle materials has conventionally been considered to be caused by the indentation and scratching by an abrasive powder as in lapping and polishing. Hence studies of material removal in the case of brittle parts have focused on fracture mechanics mechanisms. Indentation and scratch tests have been extensively used to help in this study. In the case of lapping and polishing the pressure distribution between

the workpiece and tool dominates. This has led to the concept of 'pressure copying' which has been the traditional method of surface generation to a high finish. The value of pressure on each grain (or average grain $p_c$) has been found.

There is another copying method. This is called the 'motion-copying' technique in which the tool is forced to follow a given path determined elsewhere in the machine tool. The figure on the workpiece becomes an exact replica of this motion. Diamond-turning methods on soft materials generally work on a motion-copying technique. However, to do the same on brittle materials using grinding is not simple. This has meant the determination of a critical distance $d_c$ between the workpiece and tool which has to be maintained in order to guarantee ductile conditions. In an important paper Bifano *et al* [51] found experimentally the value of $d_c$ for germanium. Subsequently, the $d_c$ values were found for a variety of materials. (This work is the basis of what follows.) It was found that the $d_c$ values could be used to tighten the specification of grinding (or cutting) machines in terms of accuracy, feed, resolution and stiffness. From this requirement a new generation of grinding machines has been developed mainly based on the work of Miyashita and co-workers [52].

Figure 6.60 sums up the problem. If the extraneous displacements, flexures, strains, or whatever, can be kept below $d_c$ for a given material, then ductile cutting can be achieved; if it is not, then the brittle regime takes over and the roughness and also the 'figure' to a lesser extent deteriorate. Figure 6.61 shows how the nanogrinding domain relates to conventional grinding in terms of grain size. Notice the difference in the height distribution of active grains. The strategy for ductile machining is shown in table 6.3.



**Figure 6.60** Accuracy related to material removal mechanism.

Table 6.3 gives the possible design criteria for grinding machines and wheels that need to be satisfied to realize the two functions of motion copying in surface generation and ductile mode in material removal. Type I specifies the case of applying a fine abrasive grinding wheel, of which the grain size is not more than the $d_c$ value, to avoid the risk of chipped or dropped grains from the wheel resulting in cracks in the machined surface. Type II specifies the case of applying a coarse abrasive wheel where the grain size is more than the $d_c$ value. This does not assume that any chipping or dropping of abrasive grains results in cracks on the workpiece surface.

**Figure 6.61** Conventional machining processes (*a*) versus nanogrinding process (*b*).

**Table 6.3** Criteria for the grinding of brittle materials in the ductile mode of material removal and the motion-copying mode of surface generation (after Miyashita and Yoshioka [52]).

| | |
|---|---|
| *Type I:* Fine abrasive wheel, grain size | < critical distance $d_c$ |
|     Wheel run-out | $< d_c$ |
|     Feed resolution | $< d_c$ |
|     Work and wheel support stiffness | Sufficiently high for setting depth of cut less than $d_c$ under grinding load |
| | |
| *Type II:* Large abrasive wheel, grain size | $> d_c$ |
|     Wheel run-out | $< d_c$ |
|     Feed resolution | $< d_c$ |
|     Height distribution of cutting points | $< d_c$ |
|     Work and wheel support stiffness | Sufficiently high for setting depth of cut less than $d_c$ under grinding |

Table 6.4 gives the design specifications of grinding machine and abrasive wheel in terms of the $d_c$ value for motion accuracies, feed resolution, truing accuracies and height distribution of cutting points and stiffness. In the worked footnote for stiffness, for $d_c = 100$ nm and a grinding load of 10 N, a stiffness of 1 N nm$^{-1}$ is required for depressing the resultant deformation of the wheel head feed system by less than 10 nm.

In summary, therefore, there are a number of well-defined requirements for a machine tool capable of achieving ductile grinding. These are, or at least appear to be, according to present information as follows:

1. Size error is related to the feed resolution and should be less than $d_c$.

**Table 6.4** Design specifications of machine tools and abrasive wheels.

| | |
|---|---|
| Feed resolution | $d_c/10$ |
| Straightness of carriage | $d_c$/travel range |
| Work and wheel support stiffness | No more deformation than the feed resolution under grinding load |
| Vibration level of work and wheel support systems | $< d_c$ |
| Truing accuracies | $< d_c$ |
| Height distribution of cutting points | $< d_c$ |

For example, a critical distance ($d_c$) of 100 nm and a grinding load of 10 N produces motion errors under the grinding load of less than 10 nm, and a stiffness of the work and wheel support system greater than the ratio of the grinding to feed resolution, that is 1 N nm$^{-1}$ (10 N per 10 nm).

2. Shape error is related to the geometry of the wheel and the error motion of the work and wheel support systems.
3. Surface roughness is related directly to the height distribution of the active cutting points on the wheel. Ideally, of course, one should be the mirror image of the other.

The criterion is shown diagrammatically in figure 6.62. There is a profound difference in the texture between brittle and ductile grinding as seen in figure 6.63. Even taking into account the different surface parameters used, the improvement is more than twenty to one.

One other point is that the $d_c$ obtained from scratch tests depends to some extent on how the tests are carried out. As an example, the way in which the $d_c$ value for silicon depends on rake angle is as shown in table 6.5.

The real problem with ductile grinding is that it is not fully understood. The actual mechanism for ductile flow is not as yet known. The very small dimension involved ($\sim 10^2$ nm) suggests that it may be linked to



**Figure 6.62** Ductile grinding criteria: wheel run-out $< d_c$; feed resolution $< d_c$; stiffness of machine tool-sufficiently high for setting depth of cut less than $d_c$ under grinding load; truing cut-ductile mode truing; wheel wear mode-ductile mode wear (no dropped grit).

**Figure 6.63** Comparison of roughness in and out of ductile mode: (*a*) $R_t$ with brittle mode grinding; (*b*) roughness in ductile mode.

**Table 6.5**

| Rake angle (deg) | $d_c$ value (nm) |
| --- | --- |
| 0 | 80 |
| −30 | 120 |
| −45 | 180 |
| −60 | 260 |

the density of dislocations in the material or perhaps faults. At these levels it is clear that explanations using conventional mechanisms of cutting are questionable. It is probable that atomic-scale mechanisms need to be used. So far all the advances have been made possible using empirical data.

### 6.4.9  *General comments on roughness*

What type of surface roughness and form can be expected from grinding and other abrasive processes, and what if anything can the texture reveal about the grinding and the machine tool? In what follows some possible mechanisms will be considered.

 The purpose of this is not to set down proven practical relationships—they hardly exist! It is to explore methods linking the surface to the process so that if a surface is shown to be suitable for a given function there will be some rules whereby the most suitable process for generating the surface can be chosen. This approach, from function to process, is the reverse of what usually happens and therefore needs to be tried out in theory in order to get a better understanding.

The centrepiece in both functional and manufacturing considerations is the surface, so it is important that it is understood. There are two basic factors which determine the surface from the process viewpoint. One is how the cutting element—the tool or grain—is presented, in space, to the surface; the other is how the cutting element affects the surface when cutting. Some idea of the former can readily be obtained from the kinematics of the machine tool and the nature of the process. The latter is much more concerned with the physical properties of the workpiece and the cutting element. In effect, one factor influencing the generation of the surface is geometrical and the other is physical.

The result is also geometrical and physical on and in the surface. It is not unreasonable to presume that the purely geometric properties are therefore carried through from the process dynamics to the surface geometry. The most effective surface function for identifying this connection is the autocorrelation function. How the surface is generated and how this shows itself in the shape of the autocorrelation function is of fundamental importance and will be considered in an elementary way here. The way this is approached might be revealing.

A number of factors might be expected to enter into understanding how the particular autocorrelation function shape is generated. Although at this stage quantitative explanations of the statistical relationship between the process and the surface roughness are difficult to achieve, there is no doubt that a relationship does exist. It is a great pity that a more concentrated attempt to use the valuable information contained in the correlation function of the roughness has not been carried out. This might have produced basic links between the process parameters, the roughness and even the surface integrity! The breakdown given here is an attempt to clarify the situation, belatedly! (Note that because abrasive processes are random it is better to consider the autocorrelation function rather than the power spectrum, where valuable visual correlation can be lost. For the same reason, periodic surfaces such as obtained in turning are best approached by using spectra rather than correlation. The fact that one can be obtained from the other is irrelevant. Every opportunity to get maximum correlation should be taken!)

The following detail can be revealed by the correlation function of the roughness profile:

1. The randomness of the process—in other words, how the grains are distributed both horizontally and vertically.
2. The shape of the grains—whether burnished, loaded or angular.
3. How each grain interacts with the surface—the efficiency of the cutting and the hardness of the workpiece material. Furthermore, indications of subsurface damage may well be visible.

Taking the first issue first, if the wheel were perfectly random with no dressing problems associated with it there would be a uniform density of grains per unit distance across the whole surface at any point along its circumference and across its width. This implies that the statistics governing the process are of a Poissonian nature. The nature of the randomness is the factor that determines the envelope of the autocorrelation function of the surface texture. Factors 2 and 3 basically determine the fine detail of the autocorrelation function (figure 6.64).

Obviously factors of the unit event (e.g. interaction effects and shape) are not necessarily independent of each other, but they are sufficiently separate to allow some information to be gleaned from the respective region on the autocorrelation function.

Figure 6.64 shows that the 'envelope' factors can be considered to be independent of the 'event' parameters.

In what follows such an investigation linking process parameters to the surface will be given, and some idealized examples will be examined.

Thus if the situation is as shown in figure 6.65, the arrows show where a grain might impinge on the surface. This is meant to be random positioning. Then the probability of $K$ such events occurring in the space $l_1 \rightarrow l_2$ is (if $\lambda(x)$ is the density of grain hits per unit length) here expressed in the general Poissonian case as a function of $x$ is:

**Figure 6.64** Relation of autocorrelation shape to process.



**Figure 6.65** Poissonian point process for grain hits.

$$p(K) = \exp\left(-\int_{l_1}^{l_2} \lambda(x)\mathrm{d}x\right)\left(\int_{l_1}^{l_2} \lambda(x)\mathrm{d}x\right)^K \bigg/ K!. \tag{6.19}$$

Obviously the usual or desired case is when $\lambda(x)$ is a constant. Under these circumstances and allowing $l_1 = 0$ and $l_2 = \tau$

$$p(K) = \exp(-\lambda\tau)(\lambda\tau)^K \big/ K! \tag{6.20}$$

which is the familiar form of the Poisson distribution.

For $K = 0$ equation (6.20) becomes

$$p(0) = \exp(-\lambda\tau) \tag{6.21}$$

and for $K = 1$

$$p(1) = \lambda\tau \exp(-\lambda\tau)$$

and so on.

Similarly if the distribution of even widths (rather than the positions at which they interact with the surface) is uniform, the probability of $K$ events occurring of width $\mu$ in an interval $W$ is given by

$$p(K) = \frac{(\mu W)^K \exp(-\mu W)}{K!}. \tag{6.22}$$

Formulae (6.14) to (6.22) all indicate that the very character of abrasive processes (e.g. that uniform distributions are likely to occur in space) points to the presence of exponential terms in the autocorrelation function of the surface. Unfortunately a simple exponential autocorrelation could only happen for hypothetical surfaces. However, it does indicate that the autocorrelation could be expected to have an exponential envelope as its natural shape. The exponential shape seems to be as natural for the spatial characteristics of the surface as the Gaussian one is for heights on the surface!

The autocorrelation function embodies the typical (average) impression and the distribution by means of conditional expectation. This is the expectation of mechanism $\gamma$ given multiple events $\beta$.

$$E_\beta \left( E\left( \frac{\gamma}{\beta} \right) \right) = \sum_{i \text{ over } \beta} E\left( \frac{\gamma}{i} \right) \dots p(i) \tag{6.23}$$

where $\gamma$ is the correlation procedure. The autocorrelation $C_p(\tau)$

$$C_p(\tau) = \sum_{i=1}^{\beta} (f(x).f(x+\tau)/f(x), f(x+\tau) \text{ in } i \text{ events}).\left( \lambda_p^i \frac{\exp(-\lambda_p \tau)}{i!} \right)$$

$$= E(f(x).(f(x+\tau) \text{ given } f(x), f(x+\tau)) \text{ within one and only one event}$$

$$+ \quad - \quad - \quad - \quad - \quad - \quad \text{two} \quad - \quad - \quad -$$

$$+ \tag{6.24}$$

Actually only one case is needed because for $\tau > L$, $A(\tau) = 0$. Also the one and only condition requires that grain centre is within $\tau - L/2$ and $L/2$ and that no other grain encompasses either the origin or $\tau$ positions.

Taking the probability and constraints into account

$$C_p = E(f(x).f(x+\tau))\exp(-\lambda_p(L-\tau)).\lambda_p(L-\tau).\exp - (\lambda_p \tau).\exp(-\lambda_p \tau)$$

$$= A_G(\tau).(-\lambda_p L).\lambda_p(L-\tau)\exp(-\lambda_p \tau)$$

$$= A_G(\tau).\exp(-\lambda_p \tau).k \tag{6.25}$$

where $K$ is a constant dependent on $\lambda_p$ and $L$, both of which are closely related to specified process parameters.

Hence the profile autocorrelation is the typical grain correlation modulated by the density of grains $\lambda_p$. Equation (6.25) encapsulates the Poisson, hence Markov feature, so the shape and size of the typical grain impression and the exponential envelope character of the correlation function explicitly demonstrate the dependence on actual process parameters. Furthermore if the typical correlation length is $C_{cp}$ then the depth of cut $d$ is approximately given by $d = \dfrac{A_{cp}^2}{8R}$ where $R$ is the grain size. For example, for an $A_{cp}$ of 20 $\mu$m and $R$ of 100 $\mu$m, $d \sim 0.5$ $\mu$m.

A mathematical model enables investigations and simulations to be carried out on a number of processes.

Thus $C(\tau) = \dfrac{K'}{4\zeta w_0^3}\exp(-\zeta w_0|\tau|)\left[ \cos(w_a \tau) + \zeta\left( \dfrac{w_0}{w_a} \right)\text{sgn } \tau \sin(w_a \tau) \right]$ \hfill (6.26)

Where the average frequency is $w_a$. $\zeta$ is a measure of damping in the cutting process.

Consider the next feature, the impression the grain has on the surface. It is informative to consider what the autocorrelation function of a single impression might be.

Thus if $C(\tau)$ is the ACF of the single impression on the surface, it is given by

$$C(\tau)=\int_0^{L-\tau} z_1(x)z_1(x+\tau)\mathrm{d}x/(L-\tau). \tag{6.27}$$

Taking a simple case of a square impression (i.e. square cut) to illustrate the point, this yields

$$C(\tau)=\sigma^2(1-|\tau|/L). \tag{6.28}$$

where $\sigma^2$ is the (RMS)$^2$ of the cutting shape and $L$ is the width. Letting $\tau/L=\bar{\tau}$ and $\sigma^2=1$, then

$$C(\tau)=(1-|\bar{\tau}|/L). \tag{6.29}$$

If the impressions are not interacting with each other but represent a succession of indentations across the surface as in figure 6.66, these produce correlation shapes which essentially consist of the convolution of the grain impression shape with the impulse train representing where the grains hit the surface (figure 6.67).



**Figure 6.66** Independent grain impressions.



**Figure 6.67** Convolution of grain impression shape with position.

This exercise gives a clue as to what happens in practice when the impressions left by the grains interact as the process removes metal.

When interactions between the unit machining events occur, the surface obviously becomes more complex but is not intractable! In general if $C(\tau, L)$ is the autocorrelation function associated with an independent grain impression and $p(L)$ is the probability density of an impression of size $L$ then a useful expression is

$$C(\tau)=\int_0^{L_{\max}} C(\tau,L)p(L)\mathrm{d}L \qquad \text{where } C(\tau,L)=0 \text{ for } L<\tau. \tag{6.30}$$

For machining processes in which the unit event may be considered to be additive or subtractive without complex effects occurring, such as material movement, the situation is different from that in which machining takes place. This might happen in spray deposition for the addition of material and as in electrodischarge machining for removal of material. If the particles are of about the same order in size and the interactions linear then the criterion for separated events as such need not be a constraint. The simple convolution explains the process. Hence, the build-up of a surface layer or the erosion of it can be represented by the profile $z'(x)$ given by

$$z'(x)=z(x)*h(x) \tag{6.31}$$

where $h(x)$ is the shape of the individual impression (which might be a composite shape as in equation (6.30)) and $z(x)$ is the spatial impulse train at which the events hit the surface).

Under these conditions [53]

$$C(\tau) = \lambda^2 \int_{-\infty}^{\infty} h(\alpha)h(\alpha + \tau)\mathrm{d}\alpha \qquad (6.32)$$

where $\lambda$ is the density of events per unit distance. Note that the integrand in equation (6.32) represents the autocorrelation function of the *unit event*, for example the single sputtered particle or the individual crater. The interaction between them completely disappears! In other words, in this special case, shown in figure 6.68, the autocorrelation function of the unit event is preserved throughout the whole process. At the same time, for this case the height distribution is classically Gaussian. The situation is not so straightforward in grinding because the interaction between furrows cut by the grain is complex. However, it is useful to demonstrate the principle above because, even for grinding and other abrasive processes, a remnant of the average unit cutting event is always present in the ACF, together with other features representing the aspect of the machine tool [54].



**Figure 6.68** Build-up of surface by superposition of unit events.

An important message that emerges from this part of the section on grinding is that in very complicated processes (and in simple ones for that matter) simulation and theoretical methods can be very helpful in pointing out certain characteristics; features can be explored which are difficult to investigate practically. From the exercise here, for example, it has been possible to show how the ACF remembers the detail of the unit event of machining, even after many revolutions of the wheel. This illustrates the potential in the idea that the surface roughness remembers what is happening in the process. It can therefore realistically be used as a basis for process investigation (e.g. for wheel loading and consequent automatic dressing).

There have been many attempts to link grinding parameters with the surface texture. They have usually ended up with products of non dimensional factors.

Very little correlation has been observed between roughness parameters and grinding parameters. Recent work [55] suggests that there may be some correlation between waviness parameters and grinding parameters especially the waviness height $W_t$ which appeared best for distinguishing process variables. A scale sensitive, fractal-based parameter smooth-rough crossover (SRC) seemed to be best for relating the wheel and workpiece topographies but, in general, a direct relationship seems impossible.

It has been mentioned above that the statistics governing abrasive processes are basically Poisson/Markov rather than fractal or scale sensitive fractal. In fact, if the grains of a grinding wheel are shown on circles representing the path of the grains, the $f(x)$ direction corresponds to Markov whereas, at any one position of $y$,

**Figure 6.69** Machine tool and process kinematics.

the profile height is the envelope of all the high cutting grains on the wheel. So $f(y)$ at any time is the resultant effect of all high grains (in that $x$ position on the wheel) which produce the surface. This is not Markov nor is it fractal. It is a 'Martingale', which is a statistical 'game' condition in which all events up to now (the present) produce just one value. This is the equivalent of the envelope, which is the profile. Plunge and creep feed grinding are good examples of this condition. Hence in very general terms the value $f(y)$ is a Martingale type of statistic in which the present is decided by all that went before and the $f(x)$ is the Markov process in which what occurs depends only on the present value.

The two statistical statements for $f(y)$ and $f(x)$ have a certain symmetry even if they are difficult to deal with. Variations in $f(y)$ represent the movement of the centre of the wheel in the $y$ direction i.e. the waviness perpendicular to the profile $f(x)$ which is a machine tool characteristic. $K$ variations in $f(x)$ are concerned with the process.

### 6.4.10 Honing

Honing is a finishing operation in which stones are used to provide a fine roughness of the order of 0.5 $\mu$in $R_a$ often on internal cylindrical bores. The stones, which are usually made of aluminium oxide (in the form of what is called sticks), comprise grains of size 30–600 $\mu$m joined by a vitrified or resinoid bond.

Honing is usually employed to produce an internal bearing surface for use in the cylinders of diesel or petrol engines. This involves the use of the hone as a finishing operation on the surface that has been previously ground. For internal honing the amount of material removed is small (~2 $\mu$m). Traditional honing methods employ three or more stones (six is typical) held in shoes mounted on a frame, which can be adjusted to conform to the correct bore size, the shoes being arranged to follow the general shape of the bore. In operation the cutting speed is about 1 m s$^{-1}$ with a pressure of 0.5 N mm$^{-2}$. In practice, the honing tool is allowed to dwell for a small amount of time at the end of its traverse before retraction. This causes all sorts of problems of interpretation of the surface roughness of cylinder bores because the lay changes direction rapidly.

Before more critical evaluation of honing is undertaken it should be mentioned that another similar process is superfinishing. This is very similar to honing but the stones are given a very small axial vibration, and it is used on external cylindrical surfaces. In honing and superfinishing the process is often used to remove the effect of abusive grinding. This is the situation in which the machining—in this case grinding—produces thermally induced stresses in the surface left by the main finishing process (grinding).

A picture of a typical honed surface is shown in figure 6.70.

An areal picture of honing can be shown in another way using diffraction by a laser throwing light at an angle through the cylinder (figure 6.71). $R_a$ values ranging from 5 $\mu$m down to 0.3 $\mu$m can be generated by the honing method.

**Figure 6.70** Plateau honing.



**Figure 6.71** Diffraction method to resolve honing angle.

The picture of honing as shown in figure 6.70 is somewhat idealized. This picture is probably typical for the middle of a cylinder bore, for example. There are, however, big differences at the end of the stroke because of the change in direction of the hone.

Usually the hone describes a figure of eight at the stroke end. This has the effect of destroying the continuity of the oil channel scratches and roughening the load-carrying lands. For this reason it does not always pay to allow the piston to utilize all the bore.

### 6.4.11 Polishing (lapping)

Polishing involves the use of free grains usually suspended in some liquid rather than bonded into a wheel as in grinding. It is invariably a process associated with generating nominally ideally smooth surfaces. This process is usually defined as one in which the microroughness is commensurate with the diameter of the molecules, or of the same scale as the lattice spacing. One could argue that the smoothest surface is the surface of a liquid in an absolutely undisturbed state. A surface approaching this ideal is that of glass when congealed in a state of rest. The faces of monocrystals grown from salt solutions present another example of a surface with a fine structure.

Of all the methods of treating a surface already mentioned, such as turning and grinding, the smoothing of the surface is at the expense of plastic deformation. The essential mechanism of fracture in plastic materials is the removal of shaving and, together with this main phenomenon, the plastic deformation of adjacent particles. However, the main phenomenon in the machining of friable materials is the development of cracks within the mass of material, which penetrate to some depth below the surface and intersect, thereby producing a mechanically weakened layer easily fractured by the repeated action of the abrasive. This is the microfracture mode mentioned earlier.

To a large extent the laws operating in the case of plastic materials, usually metals, are different to those operating during the treatment of glass and ceramics, certain rock crystals and to a lesser degree metals such as cast iron, germanium, etc.

In polishing there is a problem of determining the true mechanism of material removal and a number of theories have been advanced. It is not the intention here to go into great detail. However, it is informative to see how one outcome can result from a variety of wholly different mechanisms. The earliest polishing theories apply to that of glass polishing. French [50] suggested that an abrasive particle produces a pressure on the surface of glass which has a very small resistance to fracture in comparison with its resistance to compacting, and that this causes cracks to develop along planes of maximum shear (which are inclined to the axis at an angle of about 45°). Cleavage commences along the planes of maximum shear, but due to the fact that the particles act like a wedge, the shear plane is diverted upwards, which leads to a conchoidal fracture. Preston [57] considered the action quite differently. In his opinion the grain rolls between the tool and the glass. Owing to the irregular shape of the grains they produce a series of impacts one after the other, as a result of which conical cracks are generated. These cracks he called 'chatter' cracks. These intersect each other and, with repeated action by the abrasive grains, pieces are pulled out. Therefore he concluded that there would be an outer layer showing relief and under it a cracked layer. The direction of the cracks and the resulting stresses in crystals differ for differently orientated faces and differ markedly from glass in, for example, lithium fluoride.

The pronounced difference in the smoothness of ground and polished glass was responsible for the theory concerning the different mechanisms of these two processes and for the development of different theories of polishing.

The real question is whether a polished surface is qualitatively different from a finely ground surface or whether it is possible to develop conditions—as in ductile grinding—under which a transition from a ground surface to a polished surface with a gradually decreasing value of roughness can be detected.

It is now more or less held that mechanical removal plays an important role in polishing, as does a thermal surface flow (as discussed by Bielby [58]) and formation of a silica-gel surface (in glass) by hydrolysis.

A polished surface is characterized by the absence of cavities and projections having dimensions in excess of the wavelength of light. In the case of polished surfaces the projections are not greater than one or two molecular layers [59].

It is clear that while splintering, as described by Preston, plays a large part in the grinding and polishing of glass (and ceramics), ploughing or ductile movement of material generally only plays a small part because of the very hard and brittle nature of the material, so that in glass grinding, for example, even at light loads and stiff machines there must be a considerable danger of splintering. These remarks do not apply to ductile grinding.

It basically reduces to the question as to what extent the surface of polished material is different from the bulk. It can be asserted in accordance with French [56] that this is definitely a factor to be considered.

One pointer to this is the fact that smooth-sided scratches called sleaks do exist on glass surfaces and are quite different from the deeper scratches which exhibit the Preston-type fractures.

It has been pointed out that apart from the property of reflecting light, polished surfaces exhibit other characteristics that differ from those of the bulk [60]. Thus polished surfaces sometimes exhibit higher strengths than the bulk (although it depends on how this is defined) and increased mechanical resistance. Polished surfaces often have different heat and electrical surface properties, reinforcing the belief that polishing is more than a simple mechanical process.

Surface melting to produce the very smooth surface was shown by Bowden and Hughes in 1937 to contribute to the polishing process. The essential factor here is not necessarily the relative hardness of the lap and the material, it is just as likely to be the difference in melting point.

So far, however, there still does not appear to be a universal view of a mechanism to explain the smoothness often obtained in a few angstroms.

The terms polishing and sometimes smoothing are nominally the same as lapping. They are all basically the same process—the essential feature being the presence of free grains in a slurry and the lack of

heavy loads and high speeds. The mechanical movement needed is traditionally provided either by a machine or by hand as in the case of polishing small telescope mirrors. It is to be hoped that the more controlled ductile grinding will replace this method. More will be said later on the philosophy of the machining method used to optimize a particular function. The question is whether the measurement methodology should follow the functional requirements of the surface or whether it should reflect the manufacturing process. Or should it do both?

One of the problems with complex processes such as lapping is that sometimes there are effects which have little to do with the actual process. For example some work carried out in 1996 [61] attempts to relate removal rate and surface finish to the lapping process, yet the paper finds a very significant effect of the facing of the lapping plate on removal rate. The facing or trueing of the plate acts as a turning operation in which the feed to the facing tool is an important factor. This is hardly a typical lapping parameter!



**Figure 6.72** Removal rate and surface finish [61].

Perhaps the most significant finding is the different effects of the fixed grains with respect to the loose grains.

It is clear that loose grains dramatically change removal rate but have little effect on surface finish.

It is conceivable that the loose grains present sharp facets to the workpiece at all times, thereby enhancing removal, whereas the fixed grains having high asperities soon get smoothed and therefore have a limited effect on material removal. Quite why the surface roughness is constant is a mystery—it should decrease as the lapping sets increase! It may be that these results only relate to the materials involved (i.e. Mn—An ferrite and diamond grains) but it is more likely to be symptomatic for the fine finish lapping process.

So far the main processes have been discussed—the conventional machining methods. It should be emphasized that this chapter is not meant to be an exposition of the processes themselves, but a discussion of the way in which processes affect surfaces. In what follows, some of the unconventional methods and the geometry produced by them will be briefly considered. The word unconventional though is a misnomer as many of the processes are often used. However, for clarity they will all be grouped together in this way. Within the term 'unconventional' there are many divisions, some of which are shown simply in figure 6.73. The pictures speak for themselves [62]. Others are also included. It can be seen that they can be divided into mechanical, chemical and thermal methods. In all cases the concept of a tool producing a chip is dispensed with. The surface is literally machined by contact with free bodies, either mechanical, chemical, electromagnetic or atomic.

These methods have been developed usually to machine hard or otherwise unconventional materials. As a result the texture has been very much a secondary consideration. Hence only a simple description of the techniques will be given here.

**Figure 6.73** Mechanical and electrothermal processes.

## 6.5 Unconventional machining

The typical types are shown schematically in figure 6.73.

### 6.5.1 Ultrasonic machining

A process that similarly uses an abrasive grain in a slurry is ultrasonic machining. In this the surface generated is random as in the case of polishing and grinding, but the significant difference is that there is no easy escape route for the abrasive (or the debris particles).

The machining action is produced by the vibration of a shaped tool tip in an abrasive slurry which forms a cavity of sorts—hopefully of the required shape for a bore or shoulder in a bore in the workpiece. A light static load of about 10 N is applied to hold the tool tip against the workpiece. Material is removed by the abrasive action of the grains in the slurry which is trapped in the gap between the tool and the workpiece. The tool tip imparts the necessary energy for metal removal by means of an ultrasonic vibrator which shakes the tool tip. It is arranged that the vibration mode of the tool ensures an antinode at the tip (figure 6.74).



**Figure 6.74** Ultrasonic machining.

In this figure, the order of the antinode is optional. The vibration is fundamentally working in the cantilever mode ($m$ is an integer). This has exactly the same effect as the wear mechanism that occurs in fretting in which small vibrations trap debris between two mating surfaces.

The frequencies used tend to be in the range 20–40 kHz and the type of vibrating transducer can be piezoelectric, in which the vibration is determined by knowing the velocity of sound $C$ in the piezoelectric material given by its density $\rho$ and its elastic modulus $E$ [63]. Thus

$$C = \sqrt{E/\rho}. \tag{6.33}$$

Knowing $C$ and, given the exciting frequency $f$, the half wavelength of the wavelength of the vibration can be found as

$$C/2f \tag{6.34}$$

to ensure that the maximum amplitude vibration ensues. Other forms of transducer include magnetostrictive devices which incorporate a piece of ferromagnetic material such as nickel.

### 6.5.2 Magnetic float polishing

Polishing can also be achieved using a similar slurry [59] but there is a difference as seen in figure 6.75.



**Figure 6.75** Magnetic float polishing.

In this method the workpiece rotates (or reciprocates) the abrasive material, usually SiC up to about 40% in weight, contained in a magnetic fluid (ferricolloid). It has been verified that 0.04 $\mu$m $R_t$ is possible to achieve using 4 $\mu$m SiC particles with a uniform force provided by the magnets. The polishing force on each grain is small due to the fact that the abrasive floats. Also the temperature of each polishing point is reduced since the magnetic fluid is made thermally conductive. The potential use of this method is in making aspheric shapes as well as flat surfaces.

### 6.5.3 Physical and chemical machining

#### 6.5.3.1 Electrochemical machining (ECM)

The components are immersed in the electrolyte as the anode and the burrs or the surface skins are removed by means of an electric current. The rate of removal depends on the electrochemical equivalent of the metal and the current density.

In practice the workpiece, contained in a bath, is the anode and is connected to a dc supply. Also in the bath is the cathode, which is profiled as the negative of whatever shape is intended on the workpiece. As the cathode is lowered, the material of the workpiece is removed at a rate which depends on the cathode feed rate. One of the advantages of such a technique is that, in theory at least, the tool does not wear due to electrolysis, although it can sometimes corrode. The electrolyte depends on the application but often a dilute acid is used in preference to a salt. A typical configuration is shown in figure 6.76.

The surface roughness generally is not very good. The factor which influences it most is polarization, that is the flow and direction of flow of the electrolyte. This effect is usually present at quite a high velocity. The reason for this is set out in any production process book, but basically the requirement is that the gap between workpiece and tool should be as near as possible constant and preferably small so as to be able to maintain dimensional accuracy. This equilibrium gap $g_e$ is determined from the relationship

$$g_e = EVK/J = \rho f/J \tag{6.35}$$

where $E$ is the electrochemical equivalent, $K$ is the conductivity of the electrolyte, $J$ is the equilibrium current value, $V$ is the applied voltage, $\rho$ is the work material density and $f$ is the tool feed rate.

**Figure 6.76** Electrochemical machining.

Equation (6.35) shows that the equilibrium gap is proportional to the applied voltage and to the feed rate and that the equilibrium current density is inversely proportional to feed. In practice the value of $K$ is not constant but varies as the temperature of the electrolyte changes. Hence the gap increases as the electrolyte passes through it. So to keep the gap constant the temperature of the electrolyte needs to be constant, which implies that the flow is as high as possible. This produces flow marks on the surface, particularly if there is any debris in the fluid. Furthermore, in the case where small holes or gaps are being formed, high velocity of electrolyte is difficult to achieve. This can often lead to boiling of the electrolyte which in turn detrimentally affects the surface roughness [65]. Usually the velocity of the electrolyte is about 50 ms$^{-1}$.

The form on the workpiece is very important, perhaps even more so than the surface roughness. In die sinking by ECM it is possible to produce a shape on the workpiece that is very nearly the negative of the tool. However, the relationship between tool shape and workpiece shape is not always straightforward [65]. For example, if the tool is curved, as in figure 6.77, the gap at 1 is different to that at 2 because the gap (equation (6.35)) is proportional to the feed normal to the surface, which is $f \cos \alpha$ at 2 rather than $f$ at 1.

Hence form of an unwanted nature can easily be produced on surfaces that have been machined by ECM. As far as the roughness is concerned, very often flow marks show where the electrolyte has been. The surfaces are usually rough, of the order of 2–3 $\mu$m $R_a$

### 6.5.3.2 Electrolytic grinding

A variant on ECM is electrolytic grinding, which is used on very hard materials. In this the tool is a metal wheel which has abrasive grains in it. An electrolyte, which also plays the part of cutting fluid, is circulated



**Figure 6.77** Relationship between workpiece and tool shape.

between the wheel and the workpiece. When the wheel contacts the workpiece grinding occurs in addition to ECM. The result is a very high rate of removal and a very fine surface roughness in which the abrasion reduces the effect of polarization. Surfaces as fine as 0.02 $\mu$m $R_a$ have been produced by this method. Also, because the electrochemical action is taking most of the metal away the wheel does not wear as much as in conventional grinding.

Apart from the flow line effect often found, experimenters [66] have found that very high flow velocities considerably reduced the presence of short-wavelength components on the surface.

In some systems, especially in hole making rather than die sinking, the tool of copper or steel is insulated on the outside and just has a thin conductive tip. This tooling design affects the surface integrity and hence the fatigue and corrosion resistance of the material. Also the relatively rough surface on the side walls is a result of the low current density on the side gap caused by the insulation of the tool. Under certain circumstances an improved tool having no insulation is used with an electrolyte of $NaCl_2$. This can considerably improve the surface roughness as well as the overall reliability of the system [67].

### 6.5.3.3 Electrodischarge machining (EDM)

This is a method used to machine very hard materials and uses the discharge of an arc from one conductor to another. Basically when a voltage is applied between two conductors immersed in a dielectric fluid, the fluid will ionize if the potential difference is high enough. A spark is produced between the conductors, which will develop into an arc if the potential difference is maintained. The property of the arc that causes the removal of metal is its temperature. Typical values are between 5 and 10 000°C, well above the melting point of most metals. Local melting occurs, especially in the conductor connected to the positive terminal (figure 6.78).



**Figure 6.78** Electrodischarge machining.

The tool is lowered towards the workpiece in the liquid dielectric, usually paraffin or white spirit. When the gap is small the dielectric ionizes and the spark jumps across the gap. If the potential difference falls the arc decays. A normal gap of about 25–50 $\mu$m is maintained by a servo motor whose demand signal is actuated by the difference between a reference voltage and the gap breakdown voltage. The sparking rate is limited by the pulse rate from the generator or, for very short pulse intervals, by the de-ionization rate of the fluid.

The tool electrode is often made of copper or any other conducting material. However, for long life and high form accuracy tungsten carbide is sometimes used.

In appearance the surface roughness looks like shot blasting because of the apparent cratering. This is because the spark does not always retain one position on the surface but wanders slightly, as is the nature of electrical discharges. Because of this action the surface is isotropic in character. The form obtained in EDM can be quite accurate (~ 5 $\mu$m).

The surface roughness value is very dependent on the rate of machining: the two are conflicting. If a good surface roughness is required then the metal removal rate should be small. Surface roughness values of 0.25 $\mu$m have been achieved but only at low rates of material removal.

The use of water as opposed to oil as the dielectric in EDM [64] produces a poor surface roughness and poor subsurface conditions. These are mainly due to the fact that a thick thermally affected layer is produced; water is not efficient in getting rid of the molten metal on the surface. There is also the fact that the removal of metal can be eased somewhat by the addition of organic compounds such as sugars, polyhydric alcohol and their polymers to the water. However, with water the surface roughness is degraded, for example a removal rate of 3 mm$^3$ min$^{-1}$ with oil or paraffin produces a surface roughness of 15 $\mu$m $R_t$ whereas with water a figure of 40 $\mu$m results—a considerable worsening. Recent attempts to develop 'dry' EDM is reported in 6.3.4.2.

### 6.5.4   Forming processes

### 6.5 4.1   General

Forming as a means of improving the surface roughness seems an improbable method, yet it is possible and has been used effectively in a number of different applications. One such method is called ballizing; another is swaging, as, for example, is used in cleaning up the internal bore of a barrel in armaments.

### 6.5.4.1   Surface texture and the plastic deformation processes

Many of the processes used in engineering do not involve the removal or the depositing of material. They rely on the distribution of material from one shape to another. The first shape is easy to make and to distribute; the second shape is useful in function. An example is sheet steel rolls being transformed into car bodies by a forming process. The mechanisms involved are usually a compression together with some element of lateral flow. Sometimes the process is carried out under high pressure and temperature, sometimes just one. The main point is that contact and flow are the essential mechanisms involved and surface texture is a key feature of both. Some of these issues will be discussed in chapter 7 and the areal possibilities (3D) have already been examined in chapter 2.

The matter of surface characterization has been highlighted by the need to deal with textured surfaces such as are now produced by lasers.

One recent attempt to explain the regimes in forming [68] is based on earlier work by Stout [69].

Basically the mechanism of contact is broken up into three variants of the material ratio (Figure 6.79).



**Figure 6.79** Cross section of contact zone.

During deep drawing, stretch forming and similar processes the contact can be material/material and liquid/material.

At any level the sum of the areas has to be equal to the nominal area.

Areas in which lubricant is trapped are subject to hydrostatic pressure in the same way that liquid is trapped in elasto-hydrodynamic and plasto-hydrodynamic lubrication in ball bearings and gears.

In the figure an area of void which touches a boundary of the nominal area can theoretically allow liquid to move out to the contact zone so any pressure developed is hydrodynamic. Surrounded voids that could entrap liquids are subject to hydrostatic pressure.

The basis of the classification at any level is 'material area ratio + open void + enclosed void'. Hence three curves describing this relationship can be plotted as shown in figure 6.80.



**Figure 6.80** Apportioning of areas.

This type of map is useful because it highlights the areal properties of real surfaces rather than relying on profiles. The feature which becomes revealed is the closed void which is exactly the same type of feature as the 'col' found when areal (3D) summits are being identified (see chapter 2). The closed voids are usually referred to in terms of volume enclosed.

Areally patterned surfaces have constant values of the contact and void areas as a function of height $z$.

The problem with this development is that it does not go far enough. Ideally the contacting surface should be included because the map of areas can be misleading. Areas of voids from both surfaces can interfere with each other.

Another problem is that there is no thought given to initial contact. This is determined by using elastic theory within which surfaces' asperities are important. The height properties of the surface under purely plastic deformation do not include the surface parameters usually specified.

It should be remembered that, simple as the breakdown given above is, in terms of the number of parameters, each parameter is a curve. To get numbers some decision regarding height has to be made. In practice this should be determined by the actual application.

Values taken from the material ratio curve purporting to estimate the voids which carry lubricant have been made [70]. At the same time bearing parameters have been postulated to indicate load carrying capacity. These have been applied to a number of surfaces. One shown below in figure 6.81 is the plateau honed surface used in the automotive industry.

Some of these parameters are shown in table 6.6. How they are defined is given in table 6.7. They have been designated 'functional parameters' because they are specific to an application.

**Figure 6.81** Multi-processed surfaces (e.g. plateau honed surface).

**Table 6.6** Functional parameters.

- Surface material ratio $S_q$
- Void volume ratio $S_{vr}$
- Surface bearing index $S_{bi}$
- Core fluid retention $S_{ci}$
- Of these even a basic set has fourteen parameters

**Table 6.7** Some functional parameters.

- Surface material ratio
  $S_{bc} = S_q/Z_{0.015}$
  Where $Z_{0.015}$ is the height of the surface at 5% material ratio

- Core fluid retention index
  $S_{ci} = (V_v (0.05) - V_v (0.08))/S_q$ (unit area)
  Where $V$ is valley

  If $S_{ci}$ large then there is good fluid retention
  $<S_c <0.95 - (h_{0.05} - h_{0.8})$

- Valley fluid retention index
  $S_{vi} (V_v (h = 0.8))/S_q$ (unit area)
  $0<S_{vi} <0.2 - (h_{0.8} h_{0.05})$

Two void definitions are given below:

*Core fluid retention index $S_{ci}$*
Ratio of the void volume per unit area at the core zone over the RMS deviation.

$$S_{ci} = \left( \frac{V_v(0.05) - V_v(0.8)}{unit\ area} \right) \Big/ S_q$$

(6.36)

Large $S_{ci}$ indicates good fluid retention. For a Gaussian surface $S_{ci} \sim 1.5$. As the surface wears the index decreases.

Relation: $0 < S_{ci} < 0.95 - (h_{0.05} - h_{0.08})$

*Valley fluid retention index $S_{vi}$*

Ratio of void value in valley zone over the RMS deviation.

$$S_{vi} = (V_v(h0.8)) \big/ S_q \tag{6.37}$$

A large value of $S_{vi}$ indicates a good fluid retention. For a good random surface — Gaussian for example — this index is about 0.1.

$0 < S_{vi} < 0.2 - (h_{0.8} - h_{0.05})$

The levels 0.8 and 0.05 are arbitrary height values.

A bearing type parameter is the surface bearing index $S_{bi}$.

*Surface bearing index $S_{bi}$*

$$S_{bi} = \frac{S_q}{Z_{0.05}} = \frac{1}{h_{0.05}} \tag{6.38}$$

where $Z_{0.05}$ is the height of the surface at 5% bearing ratio.

A large surface bearing index indicates a good bearing. For a Gaussian surface the surface bearing index is about 0.6. For a wide range of engineering surfaces this index is between 0.3 and 2. As the surface wears the index increases (i.e. $Z_{0.05}$).

Although the idea of functional parameters as given above is attractive and useful, the parameters are not definitive—there is a lot of detail left out. Areal (3D) parameters are principally important in flow and should reflect the 'lay' or machine tool path. Two examples given below do not benefit from the functional parameters above.

The instrument stylus must mimic the tool path. Deviations mean errors.

The milling and honing examples cannot be easily solved. Conventional proposed areal (3D) parameters listed below in table 6.8 fare no better at describing figures 6.82 and 6.83.

The example given above is useful because it serves to illustrate the need to consider tribology in the mechanism of manufacturing processes. In the maps described in chapter 7 tribology is fundamental and has been included in prediction of functional behaviour for many years. The same comment cannot be made for manufacture.



**Figure 6.82** Milling.

**Figure 6.83** Hone reversal in plateau honing.

**Table 6.8** Conventional Parameters.
EUR 15178EN Provisional

- RMS of surface $S_q$
  Ten point height $S_z$
  Skew $S_{sk}$
  Kurtosis $S_{ku}$
- Density of summits $S_{ds}$
  Texture aspect ratio $S_{tr}$
  Texture direction $S_d$
  Shortest correlation length $S_{al}$
- RMS slope $S_{\Delta q}$
  Mean summit curvature $S_{sc}$
  Material surface area ratio $S_{ar}$

One of the problems in trying to understand the mechanisms of a process is the considerable difficulty of seeing what happens. This is especially true in contact situations and lateral movement. Some useful attempts have been made [71] in this and other experiments.

One of the elements in the test apparatus is made transparent thereby enabling some visual record to be made usually via a CCT camera backed by a computer. In the two references above one of the dies of the forming process is quartz. This enables the contact points to be examined in running conditions.

Exactly the same principle was used years ago in running gear simulation exercises like the two disc machine. However, this experiment was taken one step further by having both discs made from Perspex chosen to have the same E/H value as some hardened steels.

In the first reference different viscosity lubricants were used ranging from 80 to 23800 centi Stokes; the surface roughness changes as the degree of compression expressed in % height. As compression proceeds asperities start to grow with bulk plastic deformation. What is not clear is the mechanism for growth.

Indirectly finding the effect of the surface finish on the coefficient of friction encountered in sheet metal forming showed that the roughness influenced the pressure dependence of the coefficient of friction: the higher the percentage of real contact area the lower the pressure dependence. By means of new methods of examining the micro spots under load it is possible to see lubricant actually permeating the real contact zones.

The behaviour of the hydrodynamic pockets of lubricant (i.e. those with escape routes) and the hydrostatic were shown to behave quite differently under pressure. In the former the area of escape depended upon the pocket volume whereas in the latter it did not. Various shapes were used for the pocket [72, 73]. Although the roughness of the surfaces as usually defined were not used, the $T_a$ value of the hydrodynamic trace was used as an estimate of the film thickness. The interesting point about this work is that the areal pictures covering the die surfaces were lumped into the topographic examination. Also the experimental

**Figure 6.84** Roughness as a function of compression with trapped lubricant.

work (incorporating a transparent die) was backed up by theory giving a high degree of confidence in the work. The bringing together of the practical and theoretical aspects such as shown in this work is a good omen for future developments.

Although not much quantitative information is available on resulting surface finish to forming drawing extrusion processes etc, there are some obvious examples where it is important.

Figure 6.85 shows an extrusion. Examination of the extruded material gives a way of monitoring the process. The scratches on the material can give information on process conditions, for example, friction, texture of the die, pressure and speed.



**Figure 6.85** Injection moulding.

The shape of the die as well as its texture is important, as figure 6.86 shows.

It could be argued that shaping the corner and reducing the surface finish in figure 6.86(*b*) are the only geometrical features that could be varied. This is not completely true. The die itself should be machined with longitudinal surface lay to reduce the axial pressure gradient and increase flow. In all probability the die would be turned about the axis for convenience. This lay is detrimental to flow yet is likely to be cheaper! A balance between cost and performance has to be made before this choice is made.

Figure 6.87 shows that surface finish can hinder air venting. Skew of the mould surface finish determines venting efficiency. However, it has to be remembered that there are two aspects to the importance of the texture. This shows one example where the particular shape of the profile can help produce the mould as well as be good for the moulding. As flow is involved the 'lay' of the surface is also important as seen in the figure.

(a)

Die

Tensile stresses
greater than shear
= melt fracture

Non-streamlined flow

(b)

Shape (form)

Friction (roughness)

Smooth flow

Land

**Figure 6.86** Extrusion die shape and texture.

(*a*) Venting

(a)

? Air is trapped

(b)

Air can escape

Moulded piece in (b)    Good adhesion

Moulded piece in (a)    Poor adhesion

Profile in (b)    'Do not polish the mould'

Profile in (a)

(*b*) Runners 'resistance'

Lay
circumferential
high friction

Lay
axial
low friction

'Do not get
circular lay'

The surface finish not only affects the quality of the mould,
it affects the performance of the piece

**Figure 6.87** Examples in moulding.

### 6.5.4.2  Ballizing

In this a slightly oversized ball is forced through a premachined hole. The ball system is crudely the same as used in broaching only it is a finishing process. A schematic diagram is shown in figure 6.88.



**Figure 6.88**  Ballizing.

The passage of the ball (which is usually of very hard material such as tungsten carbide) through the hole improves the roughness and the roundness of the hole. It also achieves the important function of closing up the tolerance in size.

In the context of this book the roughness and roundness can be regarded as the surface parameters that are improved. For example, a typical roundness profile can be changed from 8 $\mu$m LSC to 1.5 $\mu$m by one pass, but obviously the final degree of improvement depends on the texture and the roundness of the original ball [74] (figure 6.89).

### 6.5.5  Micro- and nanomachining

### 6.5.5.1  General

This heading fundamentally covers processes such as when a single-point diamond is used to machine soft materials or the ultra-precision polishing of hard and brittle materials. It also encompasses 'atomic-bit' machining.



**Figure 6.89**  Effect on roundness.

As Miyamoto and Taniguchi [75] point out, the accuracy of form and dimension and even surface roughness cannot be achieved by a simple extension of conventional machining processes and techniques. As the need grows for nanotechnology—that in which accuracy is measured in nanometres rather than micrometres—the system engineering demand increases rapidly. Closed-loop systems—especially as mentioned before—closed around the workpiece become the norm.

Taking first the machining of soft materials to nanometre accuracy, this was developed mainly by J Bryan of Lawrence Livermore Laboratories who achieved the astonishing positioning of the tool to one-millionth of an inch in 1972 [76]. This was achieved by having an independent metrology loop within the machine tool.

The fundamental problem of machining in these circumstances is that small chips have to be produced in order to get the desired finish. However, as the chip size gets smaller the shear stress in the cutting tool becomes very large. In essence for depths of cut of 1 $\mu$m or less, the material in effect gets 'harder'. This is similar to the phenomenon found in using stylus instruments, that is the surface suffers very little damage despite the sometimes heavy pressure on it.

The reason for both effects is that when the indentation or cut is of the order of fractions of a micrometre, it is small compared with the average distances between dislocations.

In polishing this effect has already been mentioned. In principle of course single-point diamonds can be used to cut glass and ceramics, but the diamond tools wear quickly and so the conventional way is to use free abrasives. The depth of cut on soft materials can be less than 5 $\mu$m producing a surface of 20 nm $R_a$ and an effective depth of stress of about 1 $\mu$m.

The machining of hard and brittle materials has taken on extra importance in recent years because of the need to make a variety of components in difficult materials, for example mirrors for high-output lasers. X-ray mirrors for telescopes and microscopes, infrared lenses in germanium, scanners for laser printers and so on.

### 6.5.5.2 Micropolishing

A number of new polishing processes have been developed, mostly in Japan. These are (according to Taniguchi [32]): Elastic emission machining [77]; mechano-chemical etching; hydrodynamic polishing; and



**Figure 6.90** Elastic emission machining (EEM).



**Figure 6.91** Mechano-chemical machining.

selective chemico-mechanical polishing. Figures 6.90–6.93 show the basic mechanisms [73]. There are variants on well understood processes so they are repeated here rather than in chapter 8 on nanotechnology.

Tables 6.9–6.11 show the current processes with the associated tolerances and surface characterization.

The methods mentioned above are basically mechanical with an ingredient of chemical action included. The idea of micromachining utilizing the fracture of the material is to some extent attractive. Any process so employed must make use of some of the inherent faults or flaws already in the workpiece material. The fracture of material is generally discussed in the context of point defects, dislocations or sometimes simply cracks. The order of density respectively within the material is $10^8$, $10^5$ and $10^3$ mm$^{-3}$.

From the point of view of elastic emission machining the basic idea is to introduce the fracture preferentially in the minute dislocation-free areas. This has been attempted with EEM by using ultra-fine powder particles. The powder particles of about 10 $\mu$m size are accelerated via the rotation of a plastic sphere shown in figure 6.90 to such an extent where they hit the surface. According to Mori and Sugiyama the atomic-size fracture results and more importantly the finished surface is quite undisturbed in its physical and crystalline state. This is definitely not 'abusive' machining in the conventional sense.

As an example of an alternative approach to abrasion this method is interesting. In order to induce elastic fracture on the atomic dimension, ultra-fine powder particles (of alumina or zirconia) are mixed with water and the mixture introduced between the sphere and the workpiece. The interaction occurs in the elastic region between the spinning sphere and the workpiece. Obviously the hydrodynamic fluid film has to be smaller than the powder size. This is arranged by altering the speed of the sphere and the pressure on it.

A surface roughness of the order of 20 nm $R_a$ can be expected by this method. The problem of selective machining in the semiconductor industry using this method is still under investigation.

One fast growing process is chemical mechanical polishing (CMP). There are a number of materials besides silicon which CMP is especially suited to. They are gallium arsenide GaAs, gallium phosphide GaP and indium phosphide InP. Also it is now possible to use CMP techniques without the need of toxic agents.

Another stratagem has been to reduce or to eliminate the polishing component of the process altogether using electrolytic in-process dressing techniques (ELID).



**Figure 6.92** Mechano-chemical machining with hydrodynamic pressure.



**Figure 6.93** Chemico-mechanical machining.

**Table 6.9** Atomic-bit materials processes.

| Type | Processing mechanism | Processing method | Dimensional accuracy control |
|---|---|---|---|
| Micro-cutting mechanics | Shearing or tensile rupture microbit fracture, solid | Fine cutting with single-point diamond tool | Accuracy of shape and surface roughness (machine tools) |
| Atomic-bit separation (removal) | Elastic failure (solid, atomic bit) | EEM, magnetic fluid machining (ultra-fine polishing) | Surface accuracy (polishing tool profiling), depth accuracy (micro) |
| | Chemical decomposition (gas, liquid, solid) | Chemical etching, reactive plasma etching, mechano-chemical machining, chemico-mechanical machining | Surface accuracy (polishing tool profiling), depth accuracy (macro), pattern accuracy (mask) |
| | Electrochemical decomposition (liquid, solid) | Electrolytic polishing, electrolytic processing (etching) | Surface accuracy (polishing tool profiling), depth accuracy (macro), pattern accuracy (mask) |
| | Vaporizing (thermal) (gas, solid) | Electron beam machining, laser machining, thermal ray machining | Linear accuracy (control), surface accuracy (macro), depth accuracy (macro) |
| | Diffusion separation (thermal) (solid, liquid, gas, solid) | Diffusion removal (dissolution) | Surface pattern accuracy (mask), surface accuracy (macro), depth accuracy |
| | Melting separation (thermal) (liquid, gas, solid) | Melting removal | Surface pattern accuracy (mask), surface accuracy (macro), depth accuracy (macro) |
| | Sputtering (solid) | Ion sputter machining, reactive ion sputter machining | Surface pattern accuracy (mask), surface accuracy (macro), depth accuracy (micro) |
| Atomic-bit consolidation (accretion) | Chemical deposition and bonding (gas, liquid, solid) | Chemical plating, gas phase plating, oxidation, nitridation, activated reaction plating (ARP) | Surface pattern accuracy (mask), thickness accuracy (macro) |
| | Electrochemical deposition and bonding (gas, liquid, solid) | Electroplating, anodic oxidation electroforming | Surface pattern accuracy (mask), thickness accuracy (micro) |
| | Thermal deposition and bonding (gas, liquid, solid) | Vapour deposition, epitaxial growth, molecular beam epitaxy | Surface pattern accuracy (mask), thickness accuracy (micro) |
| | Diffusion bonding, melting (thermal bonding) | Sintering, blistering, ion nitridation, dipping, molten plating | Surface pattern accuracy (mask), thickness (depth) accuracy (macro) |
| | Physical deposition and bonding (kinetic) | Sputtering deposition, ionized plating, cluster ion epitaxy, ion beam deposition | Surface pattern accuracy (mask), depth accuracy (micro) |
| | Implantation (kinetic) | Ion implantation (injection) | Surface pattern accuracy (mask), depth accuracy (micro) |
| Atomic-bit deformation | Thermal flow | Surface tension (thermal, optical, laser, electron beam), gas (high temperature) | Surface accuracy (macro) |
| | Viscous flow | Liquid flow (hydro, polishing) | Surface accuracy (macro) |
| | Friction flow | Fine particle flow (polishing, burnishing, lapping) | Surface accuracy (micro) |

**Table 6.10** Tolerances of products.

| Accuracy/ tolerance mean value ($\mu$m) | Mechanical component |
|---|---|
| 200 | Ordinary mechanical equipment |
| 50 | Gear, screw, parts of typewriters, automobile engine parts, sawing machine parts |
| 5 | Watch parts, precision gear/screw (ball), machine tool bearing, rotary compressor parts |
| 0.5 | Ball bearing, rotor bearing, air bearing, needle bearings, flapper servo values, gyrobearing |
| 0.05 | Block gauge, diamond tool, precision *XY* table guide, microtome |
| 0.005 | Surface roughness |
| < 1 nm | |

**Table 6.11** Tolerances of products.

| Accuracy/ tolerance ($\mu$m) | Electric (electronic) component | Optical component |
|---|---|---|
| 200 | General purpose electrical equipment | Camera body |
| 50 | Electronic packages, micromotor transistor, diode | Shutter of camera, lens holder |
| 5 | Electric relay/resistor, condenser, disc memory, colour mask, video tape cylinder | Lens, prism, optical fibre connector |
| 0.5 | Magnetic head, magnetic scale, CCD elements, quartz vibrator, magnetic bubble, magnetron | Precision lens/prism, optical scale laser mirror |
| 0.05 | IC, video disc, LSI | Optical flat, optical disc |
| 0.005 | Super LSI | Precision diffraction grating |
| < 1 | Synthesized semiconductor | |

The driving force behind these polishing techniques has been and still is silicon machining so that the workpiece is flat or has a small curvature. Problems involved in having to provide a number of degrees of freedom do not arise, consequently the process is somewhat specialized.

The model of the action is not yet worked out although there have been a number of tries. Figure 6.94 shows the polishing pad as a series of asperities. Material is removed by the rubbing action of the asperities together with attrition by the slurry.

The mechanism of polishing is a fairly simple equating of the attrition to damping.

$$\frac{dT}{dt} = K.P.\frac{ds}{dt} \tag{6.39}$$

where $\frac{dT}{dt}$ is removed layer per unit time, $K$ is Preston coefficient [79]. $P$ is the polishing pressure and $\frac{ds}{dt}$ is the relative speed between pad and wafer [80].

This process is developing quickly but is expensive. The surface finish is typically 2–3 nm and the total thickness variation of the wafer is 200–400 nm which is probably more significant than the surface finish.

Despite some real disadvantages the method is suitable for the non-silicon materials which is why it has been persevered with.

**Figure 6.94** CMP model [78].

### 6.5.5.3 Three dimensional micromachining

The driving force behind the various methods of machining wafers has been the semiconductor industry. There is, however, a growing need to accurately machine small objects of the order of millimetres. This is necessary for the growing fields of micromechanics, MEMS, and microdynamics.

There is a difficulty in keeping the surface finish proportionately small relative to the dimension. One way of approaching this problem is to consider only processes which impart a shape onto the workpiece by means of a tool; (the definition of which is quite wide), following the breakdown suggested by Masuzawa and Toenshoff [81] in terms of the shape generator (shape specification element).

These processes differ from the polishing actions because the processes in table 6.12 can in principle machine three dimensionally—a factor important for miniature motors, activators and similar devices. Also, whereas three dimensions is readily possible with tools, it is not so obvious with masks. Nevertheless 3D holographic masks are possible so masks should be included. Lithography-producing surface structures do not produce 3D geometry and so are not included here.

The problem of miniaturization is not simply that everything becomes smaller and so gets more difficult to see and to handle. The control over different aspects of small workpieces is not uniform. Control of size tolerance is different from that of the surface finish. The ratio between size and roughness at one metre is quite different from that at one millimetre and at one micrometre.

**Table 6.12**

| Shape generator | Tool | | Mask |
|---|---|---|---|
| | Fixed | Controlled | |
| Process | Moulding | EDM | Etching |
| | | Turning | Electro forming |
| | Coining | Milling | |
| | | Punching | |
| | Electro forming | Grinding | |
| | | Laser beam machining | |
| | | Ion beam machining | |

The following are some extra factors for introducing errors which have to be taken into account:

(i)   Mechanical deformation
(ii)  Thermal deformation
(iii) Interface between the tool and workpiece need not be contacting — extra control is needed if there is a non-contacting situation.

Toenshoff [81] brings out the point that the co-ordinate system throughout the operations of tool making tool use and part assembly should be maintained. The same co-ordinate system could be used for tool making and use, the problem being that either one or the other is in use—there is no overlap.

One system in which the small size of the tool and workpiece can be utilized is when the worktable of the machine is used for toolmaking, for micromachining with the tool and for assembly of the workpieces made using the tool—all on the one table. Clearly the co-ordinate system can be utilized at least in part for the whole operation.

Most tools are convex to allow for a maximum flexibility of workpiece shape which is fortunate because miniature concave tools and parts pose extra problems not just of handling and machining but mainly of assembly. Key parts are often supported and constrained to move in the centre of an assembly rather than on the outside.

One of the most useful methods of producing small tools or small parts of the order of 100 $\mu$m is WEDG (wire electrodischarge grinding) (figure 6.95(a)).



**Figure 6.95** WEDG.

In this the wire is continuously moved tangentially when the workpiece is rotated and moved axially downward. Figure 6.95(b) shows some shapes.

The principle is the same as EDM so it can be used to machine materials such as tungsten. Moving the wire relatively slowly ensures a fine surface finish.

Grinding and turning can be used for micromachining. In grinding the speeds are high and depths of cut small so in many cases the mode of cutting is ductile. In turning diamond cutting is used. In both methods care has to be taken to ensure that the workpiece does not deflect during the machining so a strong work material is essential. 50 $\mu$m diameter tools have been produced.

Machining of concave shapes such as holes down to 5 micrometre have been achieved using various techniques. Some processes are electrochemical, ultrasonic and laser beam methods. [82, 83, 84].

Detailed description of such techniques is not within the scope of this book. It is becoming more difficult to consider surface properties such as geometry and finish independently. Parameters chosen to ensure a sharp edge for example usually produce a good finish.

Attempts to get a small 'machining unit' such as by using very short laser pulses or high frequency ultrasonics produce reasonable surface with little subsurface damage (i.e. high integrity).

It is now considered feasible to mix processes using tools and masks to get both workshape flexibility and a unit of machining near to the atomic level. Current thought believes it is feasible to get 3D parts at 10–50$\mu$m dimensions. These invariably have an axis of rotation that enables the dimension and finish to be achieved. However, prismatic parts at about 1 mm dimension are still difficult because the tool has to move tangentially without the benefit of a rotation!

### 6.5.6 Atomic-scale machining

#### 6.5.6.1 General

The alternative to the concept of machining using physical particles is to use even smaller particles than powder. This means the use of particles of atomic size. Various names have been given to machining processes doing this: energy beam processing, ion sputtering forming, ion sputtering machining, and so on. Basically in these methods a high-energy beam of electrons or ions is made to impinge on the surface. Often such methods are used, not to make very fine surfaces, but to machine fine patterns onto the workpiece. Another is to remove abused layers of surface produced by, for example, EDM.

#### 6.5.6.2 Electron beam machining

##### (a) Electron beam methods

At one time it was thought that electron beams could be used to achieve ultra-fine precision milling of surfaces using the argument that it is difficult to get anything smaller than $10^{-12}$ mm, and although the mass is small (~$10^{-28}$ g) the energies can be high (~100 keV). Furthermore electron beams are capable of being focused down to small spots of the order of 1 $\mu$m at high speed.



**Figure 6.96** Electron beam machining. $V$ is acceleration voltage, $\rho$ charge density, $R = 2.2 \times 10^{-12} V^2/\rho$ cm. For $V = 50$ kV, $R = 7$ $\mu$m; for 10kV, $R = 0.3$ $\mu$m.

However, there is a basic problem concerning the depth of penetration. At 50 kV for example, according to Taniguchi [63], the depth of penetration in aluminium is about 10 $\mu$m, which is very large. The situation is shown in figure 6.96. The action is often referred to as thermal machining.

The concept of using atoms or ions to machine in the nanometre range is satisfying. It obeys the metrologist's inclination, which is to match scales of size in every functional undertaking. Following this philosophy the metrology involved should also be atomic in its nature. X-rays have been used for this very reason. In addition, as will be shown, some of the atomic repercussions of atomic machining can be used to monitor the process at the same time, for example electrons emitted from the surface can be used to image it. General atomic type mechanisms are shown in table 6.13.

**Table 6.13**

| Type | Processing mechanism (atomic-bit size) | Processing method (energy particle beam) | Direct (without datum surface) dimensional accuracy control system (depth (ID), pattern (2D), shape (3D)) |
|---|---|---|---|
| Microcutting | Shearing slip or tensile rupture microbit fracture (solid) | Fine cutting with single-point diamond tool (for soft materials) | Accuracy of shape and surface roughness (machine tool dependent) |
| Micropolishing | Elastic atomic failure (solid, atomic-bit failure) | EEM, magnetic fluid machining (ultra-fine polishing) (for hard or brittle materials) | Surface accuracy (polishing tool profiling), depth accuracy (micro) |
| Atomic-bit separation (removing) | Chemical decomposition (gas, liquid, solid) | Chemical etching (isotropic, anisotropic), reactive plasma etching, mechano-chemical machining, chemico-mechanical machining (ion, electron, laser beam photoresist etching) | Surface accuracy (preformed profiling), depth accuracy (macro), pattern accuracy (mask) |
| | Electrochemical decomposition (liquid, solid) | Electrolytic polishing, electrolytic processing (etching) | Surface accuracy (preformed profiling), depth accuracy (macro), pattern accuracy (mask) |
| | Vaporizing (thermal) (gas, solid) | Electron beam machining, laser machining, thermal ray machining | Linear accuracy (position control), surface accuracy (macro), depth accuracy (macro), pattern accuracy (macro) |
| | Diffusion separation (thermal) (solid, liquid, gas, solid) | Diffusion removing (dissolution) | Surface pattern accuracy (mask), surface accuracy (macro), depth accuracy |
| | Melting separation (thermal) (liquid, gas, solid) | Melt removing | Surface pattern accuracy (mask), surface accuracy (macro), depth accuracy (macro) |
| | Sputtering (solid) | Ion sputter machining, reactive ion sputter machining (reactive ion etching, RIE) | Surface pattern accuracy (position control mask), surface accuracy (macro), depth accuracy (micro), shape accuracy (preformed profile) |

**Table 6.13** (*continued*)

| Type | Processing mechanism (atomic-bit size) | Processing method (energy particle beam) | Direct (without datum surface) dimensional accuracy control system (depth (ID), pattern (2D), shape (3D)) |
|---|---|---|---|
| Atomic-bit consolidation (accreting) | Chemical deposition and bonding (gas, liquid, solid) | Chemical plating, gas-phase plating, oxidation, nitriding, activated reaction plating (ARP) | Surface pattern accuracy (mask), thickness accuracy (macro) |
| | Electrochemical deposition and bonding (gas, liquid, solid) | Electroplating, anodic oxidation electroforming, electrophoresis forming | Surface pattern accuracy (mould, mask), thickness accuracy (micro) |
| | Thermal deposition and bonding (gas, liquid, solid) | Vapour deposition, epitaxial growth, molecular beam epitaxy | Surface pattern accuracy (mould, mask), thickness accuracy (micro) |
| | Diffusion bonding, melting (thermal bonding) | Sintering, blistering, ion nitriding, dipping, molten plating | Surface pattern accuracy (mask), thickness (depth), accuracy (macro) |
| | Physical deposition and bonding (kinetic) | Sputtering deposition, ionized plating, cluster ion epitaxy, ion beam deposition, ion beam mixing | Surface pattern accuracy (mask), depth accuracy (micro) |
| | Implantation (kinetic) | Ion implantation (injection) | Surface pattern accuracy (position control mask), depth accuracy (micro) |
| Atomic-bit deformation | Thermal flow | Surface tension (thermal, optical, laser, electron beam, gas high temperature), flattening | Surface accuracy (macro) (preformed profile), surface depth, pattern accuracy (mould) |
| | Viscous flow | Liquid flow (hydro) polishing, injection moulding | Surface accuracy (macro) |
| | Friction flow (shearing slip) | Fine particle flow polishing (rubbing, burnishing, lapping, coining, stamping) | Surface accuracy (micro) |
| | Molecular orientation (ion rubbing) | — | Surface pattern accuracy (stamping) |
| Surface treatment | Thermal action (electron, treatment photon, ion, etc) | Hardening, annealing (metal, semiconductor), glazing, diffusion | (Macro): Open-loop control of macroscopic processing conditions or states of processing bit tools, i.e. atom, molecule of gas, liquid, ion, electron, photon |
| | Chemical action (reactive) | Polymerization, depolymerization | (Micro): Closed-loop control of microscopic states of workpiece by means of feedback control of state of work |
| | Electrochemical action (electron, ion, photon beam assisted) | Surface reactive finishing | |

The penetration effect is a serious disadvantage of electron beam machining as such but it does not preclude the use of electron beams for other applications concerned with the surface. One such application is electron beam lithography in which integrated circuit marks are patterned by exposing a photoresist to the electron beam. This electron beam initiates polymerization of the photoresist. But because the wavelength of the electron beam is very small compared with the wavelength of light—a few orders of magnitude—the resolution can be very much higher when compared with exposure due to visible light. By this means gate densities of many orders of magnitude better than ever expected are being achieved in the semiconductor industry. But the method is not as yet relevant to fine finish. The question arises as to what is finish at this scale.

### 6.5.6.3 Ion beam machining

This is an atomic-bit machining process capable of very high resolution. A simplified form is shown in figure 6.97.



**Figure 6.97** Ion beam sputter — machining mechanism.

Ion beam machining is characterized by the rather unique ability to machine most non-organic materials whether they be metals, ceramics or even semiconductors. The method differs from electron beam machining in that the mechanism is not basically 'thermal' but it is a sputtering process — atoms are literally knocked off the surface. This is achieved by bombarding the surface with argon ions or other inert gas ions. These ions are accelerated to energies of about 10 keV and elastically collide with the surface atoms and knock them out. The penetration depth of an ion at 1 keV is said to be about 5 $\mu$m — not so deep as high-energy electrons.

Generally the method does not generate heat or mechanical strain in the surface layer. However, some of the ions are retained. The ratio of atoms removed to the number of incident ions is called the sputtering rate. The characteristics are shown in figure 6.98 [63, 66]. The sputtering rate $S$ varies with incident angle of the ions.

There are a number of different types of source, namely the duo-plasmatron type, the high-frequency plasma type and the ion shower type. The latter two are most often used. Basically, in the ion shower type the ions are generated in the plasma source of argon gas of low vacuum and extracted to the machining chamber of

Perpendicularly incident Ar ions

(a), (b) Sputtering ratio S atom/ion

(c) Rate of machining depth or speed $V_s$ = $S_r \cos\theta$, cm/s/(A/cm$^2$)

(a), (c) 10kV, BK·7 glass
(b) 50kV, (SiO$_2$)

$\times 10^{-5}$

Ion incident angle (0°, 30°, 60°, 90°)

I: Ion beam current (A)
a: Cross-sectional area of ion beam (cm$^2$)
A: Ion projected area (cm$^2$)

Machined surface

Machined volume (cm$^2$)
r = Machining time (s)

S: Sputtering rate
$$= \frac{\text{Number of sputtered atoms}}{\text{Number of incident ions}} \quad \left(\frac{\text{atom}}{\text{ion}}\right)$$

$V_s$: Rate of machining depth (speed)
$$= \frac{\text{Sputter machining depth}}{\text{Ion beam current density}} \quad \left(\frac{\text{cm/s}}{\text{A/cm}^2}\right)$$

$$= \frac{(d/r)}{(I/a)} \frac{\left(\frac{\text{cm}}{\text{S}}\right)}{\left(\frac{\text{A}}{\text{cm}^2}\right)} = S \cdot \frac{M}{NZ_{pe}} \cos\theta = S_r \cos\theta$$

$$S: \left(Ad\rho \frac{N \cdot Z}{M}\right) \Big/ \left(\frac{I_\varsigma}{\varepsilon}\right) \qquad a = A \cos\theta$$

$\rho$: Density (g/cm$^3$)
N: Avogadro's number $6.02\times10^{23}$ $\left(\dfrac{\text{Molecule}}{\text{mole}}\right)$
e: Elementary charge $1.6\times10^{-19}$ (C/ion)
Z: (atom/molecule)
M: Molecular weight of specimen (g/mole)
$S_r$: Sputter machining rate (cm$^3$/C)

**Figure 6.98** Characteristics of ion beam sputtering.

high vacuum through the accelerating electrode and on to where the specimen is set up. The high-frequency or radio-frequency plasma method is that in which the target is laid directly on the cathode of a parallel-plate electrode in between which is a plasma of argon gas built up by low-vacuum and RF excitation. In the latter method the target is bombarded by ions and electrons alternately, but in the ion shower the target is only hit with ions.

Some applications of sputtering are given below. Most often they are concerned with the removal of material to provide a given shape of form or pattern rather than a specific surface texture.

Consider one case in which diamonds are sharpened [85]. An ion shower method was used but the apparatus had to be fitted with a neutralizer to stop charge being built up on the diamond. Using sputter machining it was found that after 8 hours the tip was submicrometre (figure 6.99) [86, 87].

It was found that using this technique the surface roughness of the diamond tip was the same as that obtained by mechanical polishing. Evidently this method is very controllable and useful for making instrument styluses, microtome knives and the like. Styluses of tip dimension 0.01 $\mu$m have been made by this method. Another forming application is the use of ion milling in making aspheric lenses (figure 6.100).

Ion sputtering can be used to make fine surfaces on single and polycrystalline alumina ceramics, for example. It need not necessarily be used for smoothing the surface and can indeed be used to roughen up the surface to prepare for biological implants [87]. It is especially useful for the machining of very hard and brittle materials such as tungsten carbide, diamond and ceramics, but can also be used for iron and steel [88].

**Figure 6.99** Machining of diamond tip.



**Figure 6.100** Aspheric machining.

The surface roughness in the latter materials depends on whether or not the specimen is moving. Typical results are shown in table 6.14.

The angle of incidence of the ions is important because this affects the rate of removal of atoms (figure 6.101). There are two effects of angle. As the angle from the normal increases the ease with which atoms are displaced increases, but the actual rate decreases because the ion current density decreases by the $\cos \theta$ factor and the number of reflected ions increases rapidly [86]. Another factor is whether or not the angle is changed during bombardment (figure 6.102).

In general, although for brittle materials like silicon, glass, etc, the surface can be smoothed by ion sputtering quite effectively, for many metals this is not so. The surface can become slightly rougher than the initial roughness even if the original texture is mirror like. There are exceptions to this, however. For example, in the ion sputtering of tungsten carbide the surface texture can be maintained.

This is very important because it means that the fine dimensional adjustment of gauge blocks can be made using ion sputtering at the same time as maintaining the finish [87]. Another major use of ion sputtering is to remove abusive layers of the surface damaged by processing (e.g. by EDM).

**Table 6.14**

| Apparatus | Heat treatment | Stationary | Normal | Uniformly changing | |
|---|---|---|---|---|---|
| Ion shower 1.5 keV | Quenched | 0.52 | 0.10 | 0.77 | < 0.04 |
| | Annealed | 0.47 | 0.12 | 0.74 | < 0.06 |
| RF plasma 1.4 keV | Quenched | 0.52 | 0.10 | 0.77 | < 0.05 |
| | Annealed | 0.47 | 0.15 | 0.74 | < 0.07 |
| | | Depth of sputtering | $R_t$ | Depth of sputtering | $R_t$ |

**Figure 6.101** Machining rate as a function of angle.



**Figure 6.102** Effect of changing angle during machining.

### 6.5.6.4 *General comment on atomic-type processes*

Note that there is a fundamental difference between conventional machining and atomic-bit machining. For processing methods in which relatively large chips are removed there will always exist some defect in the workpiece material where fracture failures can be initiated at a low specific processing energy to allow material to flow hence chip to form. This enables the use of solid tools. However, it becomes difficult to use such tools in atomic-bit processing because chipping and wear of the tools become too large to allow the workpiece to be machined with sufficiently high accuracy. It becomes necessary to supply processing energy directly to the processing point on the workpiece using the high-power energetic energy beam. For example, the threshold of specific processing energy ($Jm^{-3}$) necessary to remove a unit volume of stock material from a workpiece atom by atom is extremely large.

Some idea of the density of flaws and dislocations can be gained from table 6.15 for iron.

Table 6.16 shows the types of energetic beam machining currently being considered. The processing energy beam has to be a directionally oriented flow of energy flux consisting of very fine particles. In principle these can be of chemical or electrochemical reactive energy. Alternatively, photons, electrons or ions can be used.

Recently another type of beam, such as thermally activated or accelerated neutral atoms or molecular beams, has been used for the deposition and growth of semiconductor layers.

Table 6.15 Threshold of specific machining energy ($\delta$ J cm$^{-3}$) for removal of iron (from Taniguchi [63])

| Processing mechanism | Defects or uniformities (and processing unit, m) | | | | Remarks |
|---|---|---|---|---|---|
| | Atom/molecule ($10^{-10}$ – $10^{-9}$ m) | Point vacancy (atomic cluster) ($10^{-9}$ – $10^{-7}$ m) | Movable dislocation, microcrack ($10^{-7}$–$10^{-5}$ m) | Crack, cavity, grain boundary ($10^{-5}$ – $10^{-3}$ m) | Deposition/consolidation, deforming |
| Chemical, electrochemical decomposition (dissolution) | $10^5$ – $10^4$ | $10^4$ – $10^3$ | | | Chemical/electrochemical plating, bonding, reacting (photoelectron) |
| Brittle or tensile fracture (cleavage) removal | Brittle materials (glass, ceramics) | $10^4$ – $10^3$ (elastic cluster tensile breaking) | $10^3$ – $10^2$    Microcrack fracture | Brittle crushing | |
| Plastic, shearing deformation slip (microcut, polishing) removal | Ductile materials (metal, plastics) | $10^4$ – $10^3$ (elastic cluster shear slip) | Dislocation slip    Plastic deformation slip<br>$10^3$ – $10^1$<br>($\omega = 10^6$ – $10^1$) | | |
| Melting, diffusion separation (thermal) | $10^5$ – $10^4$ | $10^4$ – $10^3$ (cluster) | | | Diffusion hardening, nitriding, bonding |
| Evaporating separation (thermal) | $10^6$ – $10^4$ | $10^4$ – $10^3$ (cluster) | | | Vacuum evaporation, molecular beam deposition/epitaxy, vapour deposition |
| Lattice atom separation (ion sputter, ion etching) | $10^6$ – $10^4$ | $10^4$ – $10^3$ (cluster) | | | Atomic consolidation, ion, deposition, implantation, sputter deposition |

$\omega$ = specific stock removing energy, J cm$^{-3}$.

$\delta = \omega$ except plastic deformation.

$\delta$ for elastic deformation = $Y^2/2E$. ($Y$ = elastic limit).

$\delta$ for elastic shearing failure = $\tau_{th}^2/2G < \sigma_{th}^2/2E$

$\sigma_{th}$, $\tau_{th}$ = theoretical tensile and shear strength.

Glass and ceramics with microcracks undergo brittle fracture (processing unit over about 1 $\mu$m), but with no defects (processing unit under about 1 $\mu$m) elastic shearing failure or shearing slip occurs.

**Table 6.16**

| | Beam type | Focusing | Environment | Processing mechanism | Application |
|---|---|---|---|---|---|
| a | Photon (laser) | Focused or broad beam | Air | Thermal (photon reactive) | Removing, consolidating |
| b | Electron | Focused or broad beam | Vacuum | Thermal (electron reactive) | Removing, consolidating |
| c | Electrodischarge current | Total guided beam | Insulation oil | Thermal, microwave for dielectrics | Removing |
| d | Ion | Focused or broad beam | Vacuum | Dynamic (ion reactive) | Removing, deposition |
| e | Reactive ion (cold plasma) | Broad beam | Vacuum | Chemical and sputter etch (dynamic) | Removing, deposition |
| f | Atomic/molecular (cold molecular and hot molecular) | Broad beam | Vacuum | Thermal and dynamic | Deposition, diffusion |
| g | Plasma (hot) | Focused beam | Air | Thermal | Removing, consolidating |
| h | Chemical reactants | Broad beam | Chemical liquid | Chemical etch, reactive process | Deposition, removing |
| i | Electrolyte reactants atom | Tool-guided beam | Electrolyte | Electrochemical etch, reactive process | Deposition, removing |

There are other advantages of using beam methods which do not immediately reveal themselves. One is that it is not necessary to use solid tools and their associated fitments — which means in effect that the machine tool for this purpose can be thought of as having infinite stiffness!

### 6.5.6.5 *Molecular beam epitaxy*

Molecular beam epitaxy is a process which can grow a crystalline thin layer having a thickness of molecular or even atomic dimension. The method has been used to prepare new semiconductive materials but the mechanism by which a flat surface results is still not fully understood. In what is called homo-epitaxial processing, a molecular beam of an element—say silicon—can be impinged onto a silicon lattice having flaws.

In taking up positions of low potential energy the silicon molecule beam in effect corrects the crystalline state. The result can be remarkably flat and smooth. Surfaces of 1 or 2 nm can be produced which are substantially fault-free. Furthermore, the overall orientation of the lattice can be improved over that of the substrate [90].

One of the problems with this technique is the care with which the operation has to be conducted. Very high vacuum is needed and the substrate has to be capable of being heated to 1000°C to allow Si to be deposited.

At present this method cannot be widely used but it has its uses as a corrective process. It is in effect the first example of a finishing process within nanoprocesses.

### 6.5.7 *Structured surfaces—engineered surfaces*

There has been a recent development in the functional significance of surface finish. This is the use of periodic patterns on the surface to enhance a specific use. Evans and Bryan have compiled a comprehensive set of applications ranging from road signs to mouse mats. One important contribution has been a attempt to clarify the distinction between 'engineered' surfaces and 'structured' surfaces.

A consensus from CIRP Group S is:

Structured surface—surfaces with a deterministic pattern of usually high aspect ratio, geometric features designed to give a specific function.

Engineered surface—surfaces where the manufacturing profess is optimized to generate variation in geometry and/or near surface material properties to give a specific function.

A proviso is that deterministic patterns need not be of uniform spacing (for example a Fresnel lense is structured but not evenly spaced). Also structured surfaces need not be isotropic. The whole problem of characterizing structured and engineered surfaces has not been solved. A possible approach has been indicated in chapter 2 with the use of space frequency functions such as the Wigner function. Structured surfaces have generally simple amplitude form and complicated spatial form which are much easier to deal with in areal space domain. Traditional surfaces tend to be the opposite (i.e. complex height form and simple spacing form).

It is not usually realized what a wide range of uses these surfaces have. They include the following [91]:

**Table 6.17** Functions of structured surfaces.

| Function | Example |
|---|---|
| Optical | Gratings |
| | Fresnel lenses |
| | Diffractive optics |
| | Reflective road signs |
| | Filters |
| | Wavelength-specific mirrors |
| Mechanical Contact | Vacuum chucks |
| | Seal surfaces |
| | Diesel injectors |
| | Piston rings/cylinder liners |
| | Synchro rings |
| | Hard disk surfaces |
| | Velcro |
| | Grooved roadways |
| Hydrodynamics | Tyre treads |
| | Drag reduction film |
| | Deck shoes |
| | Golf balls |
| Metrology artefacts | Distortion test artefact |
| | Roughness artefacts |
| | Tactile test artefacts |
| Friction and wear | 'Undulated' surfaces |
| | Abrasives, tools, files |
| Biological | Cell culture systems |
| | Capillary electrophoresis |
| | Breast implants |
| | Bio-MEMS, fluidics |
| Adhesion | Water seals |
| | Epoxy dental fillings |
| | Surface preparation for paint |
| Thermal | Heat exchanger fins |

For many of the above uses it makes sense to use a patterned surface. Perhaps the most unlikely use is in mechanical contacts such as seal surfaces.

The idea that the contact between surfaces should occur with a small vertical movement is attractive—giving a very high initial stiffness. The actual structure—leading inevitably to escape paths of liquids or gases seems less convincing. For best credibility some aspects of flow or scatter are the obvious choice of application.

It has been suggested that having the choice of patterned surfaces gives the designer an extra degree of flexibility when specifying a surface for a given function. This may well be true when the properties of structured surfaces become better known.

Clearly in terms of light scatter, the theory of Fresnel lenses, diffraction gratings and reflectors is known and used but applications involving friction and wear and normal contact are less well known. Bearings with spiral grooves have been used for many years to increase the pressure build-up in the gap between journal and shaft for example [92].

Structuring reduction rollers in steel mills was advocated by Schneider in the USSR for many years before it was convenient to do so with lasers [93].

Some ways of making structured surfaces are listed below. According to Evans and Bryan there are four ways:

Moving/removing material:
    Machining
    Knurling, burnishing
Etching
    Plating
    Evaporation
Replication
    Hot embossing
    Infection moulding
    Casting
Material modification
    Laser texturing
    Solidification of liquids (e.g. polymer curing)

An equally convenient classification of fabrication process is according to:

• Those used for direct fabrication of the part of interest or of a mould;
• Those used to replicate.

One of the basic problems in dealing with structured surfaces that are different from ordinary machined surfaces is the high slopes often encountered.

Each unit of pattern can be almost vertically attached to a substrate, making it very difficult to make and to measure. Such high slopes are often required for example in antireflection panels. Ways of measuring these surfaces are discussed in chapter 4.

## 6.6   Ways of making structured surfaces [94]

### 6.6.1   General

In general structured surfaces are made using conventional processes like cutting with a single point tool such as diamond turning. The actual process depends on the application but (for very fine detail the method adopted is usually diamond turning (for example in Fresnel lenses). Crude structure is obtained by

plunge cutting the diamond onto the surface. More refined shapes are obtained by profiling with a sharp tool as in making sinusoidal surface texture standards. The problem with diamond turning is tool wear. This is not by attrition but is usually in the form of chipping produced by hard spots in the material being turned.

Laser machining is an attractive way of making structured surfaces. Usually a short wavelength ultraviolet laser is used of 157, 193, or 248 nm and the laser is applied in very short pulses. This has the advantage of having no heat damage to the surface. Using laser machining enables structure to be imparted on many different and difficult materials [95].

Sometimes masks are used in association with the laser beam which gives a wide variety of possibilities of form—the so-called ablation geometry.

Solid state lasers can be used with transposed wavelength to put them into the UV range. Typical values of 266 nm with nanosecond duration are used which can cause instant vaporization of material at the focus of the laser. Excimer lasers have also been used with wavelengths of 193 nm.

Electron beam and ion beam milling have also been used for structuring. Because the equivalent wavelength is very much smaller than the laser, finer detail is possible, although absorption is less controlled and the machining rate is very slow.

X-ray lithography is another possibility using a collimated beam from a collimated source. Such a method is LIGA which can produce some very deep grooves with aspect ratios of 100:1. Polymer resist films have been structured, developed and then plated to produce a mould from which many replicas can be made.

Because of the relative ease of production of structured surfaces by the methods briefly touched upon above and others, structured (and engineered) surfaces will be an increasingly useful tool for the designer.

### 6.6.2 The interface

The surface itself has been considered in sections up to 6.5. However, it is well known that this is only part of the story. Just as important is what happens under the surface. In what follows the mechanism of the machining will be briefly discussed. This will be followed by an account of sub-surface investigation.

Unfortunately there does not appear to be a comprehensive model for material removal.

The criterion for the surface, the interface and the chip, however formed, is that of energy balance. It is the energy in creating the chip and the surfaces that determines how much heat and stress gets into the subsurface.

Peklenik tried to develop parameters of the process which reflect the status of the process itself [96].

The basic idea is that as the process becomes unstable, for whatever reason, this will show up as variations in the output power spectrum, which can be measured or estimated by the entropies of the process, which in turn is given in terms of the standard deviation of the energy quanta $\sigma E$ and that of the frequency



**Figure 6.103** Surface interfaces.

limits of the energy quanta $\sigma_{\Delta f}$. The energy quanta are taken to be the discernible peaks in the spectrum. If $K_1$ and $K_2$ are contacts the entropy $S$ is given by equation (6.40).

$$S = K_1 \ln(K_2 \sigma_E \sigma_{\Delta f}) \qquad (6.40)$$

Values of $K_1$ and $K_2$ are constants for the process and are estimated. The value in this work is the partitioning of the available energy into distinct interfaces.

(a) The propagating surface interface
(b) The friction surface interface $\Big\}$ shown in figure 6.103

Thus, if $U_i(t)$ is the input energy, $U_p$ the propagation energy and $U_p$ the frictional energy, then, neglecting the surface free energy

$$U_i(t) = U_p(t) + U_f(t) \qquad (6.41)$$

How this form (6.41) characterizes the process and takes into account the entropy is not relevant to the discussion; the important point is the realization that the energy is split up between the surface interfaces which determine the actual surface texture produced and the subsurface properties.

It would be interesting to see if the so-called energy quanta had any effect on the surface finish 'along the lay' or even across the lay. This could then link the new suggested process parameters with the functional output.

The problem of energy partition as laid down in equation (6.41) is that it does not take into account the scale of size. In conventional cutting the propagation term is paramount but in nanocutting the balance changes and what is neglected—the surface free energy—is now important. This factor makes the difference between fracture mechanics and ductile machining.

One redeeming factor is that nanoprocessing can be checked realistically using molecular dynamics (MD). The process is concerned with the movement of finite numbers of molecules which can be accurately reflected in the simulation. It has to be remembered, however, that large number simulations are difficult if not impossible to check.

Energy partition has been used often e.g.[97].

One recent example uses a two-colour infra red detector to measure temperatures within the cutting zone. Given the measured temperatures and the temperatures calculated from the thermal Fourier equations, the energies can be evaluated. In one case it was possible to show that about 20% of the heat is transported to the workpiece in wet grinding whereas it is nearer 60% for dry grinding thereby allowing a quantitative measure of coolant efficiency to be obtained [98].

### 6.6.2.1 Brittle/ductile transition in nanometric machining.

The transition from brittle to ductile machining at the nanoscale is accepted. What is not so clear is the mechanism that produces the transition. It is usually described in fracture mechanics in terms of the energy balance described above (equation (6.41)) between strain energy and surface energy [99].

It has also been pointed out that it is probably the presence of defects in the material which allows the ductile condition. Invariably the idea is that defects transform into cracks. The inference is that no ductile effect will be observed in a defect free material.

As usual nowadays the hypothesis is tested by molecular dynamic simulation (MD) or RMD (renormalized molecular dynamics) in which several scales of size can be used by lumping the unit blocks of material successively together. This means that the first simulation can be in terms of clusters of molecules, then molecules, them atoms. Quantum effects cannot be simulated yet.

A typical result is that a defect-free monolith of, say, silicon can be machined in a ductile mode only in a vacuum but exhibits brittle/ductile transitions in a normal atmosphere. The atmospheric gases can, it seems,

make a big difference in the possibility of cracking—the atmospheric molecules in effect look like defects when taken into crevices caused by the plastic flow of the base material [100].

Computer simulations using MD are reasonably straightforward for the cutting mechanism such as grinding with the numerous grains of the wheel making contact in a random way. One approach has been to simulate what is in effect the 'unit event' of grinding mentioned earlier. It is the 'average' grain impression.

In one application this unit event is taken to be the behaviour of a hardness indentor acting normally onto the surface followed by a movement of the indentor laterally across the surface to simulate the grinding scratch (figure 6.104).

It is claimed that looking at the MD pictures of the indentation enables the onset of ductile properties to be detected. This information can then be used to find the minimum depth of cut before the ductile regime takes over. [101].

It has to be said that pictures obtained using molecular dynamics are not always clear visually as well as intellectually. Often the boundary conditions are not clear or the criteria for molecule movement are absent.



**Figure 6.104** (*a*) Indentation, (*b*) movement of grain.

There is one example of using MD which shows an effect very clearly. It is the effect of crystallographic orientation [102]. The MD approach shows clearly the effect of the crystal orientations and the cutting directions (Fig. 6.105). In this case it is sufficient to see the directions of the stress fields; there is no need for quantitative values. Under these circumstances the molecular dynamic simulation is valuable.

The alignment of the crystal and the direction of the tool determine the mode of deformation in front of the tool. For example an orientation of [111] and a cut in the [110] direction produces chip formation which is mainly shear and some compression. There is little subsurface deformation as the atoms recover elastically. If the crystal is orientated in the [110] plane and the cut in the [001] direction the dislocations are generated normal to the cutting direction which leads to more compression than shear.

### 6.6.3   General design points

There are a number of different concepts relating to workpiece manufacture and design. One point concerns the workpiece tool interface. It has always been taken for granted that the chip area and the length of the active cutting edge are straightforward. Recently the chip thickness, tool life and cutting forces have been



**Figure 6.105**  [010] on [111].

**Figure 6.106** Modes of deformation in the shear zone of cutting a single crystal.

more difficult to assess because of exotic tooling. Usually leaving out rake angles and inclination angles does not affect chip calculations or the shape of the interface but now there are rotating tools giving anticlockwise cutting which complicate the situation. For this reason more rigorous mathematical methods have had to be developed, which should be incorporated into the design strategy [103].

Also the concept of 'the machining surface' has recently been introduced into design strategies. This surface includes all the information necessary for the driving of the tool so that the envelope surface of the tool movement sweeping the 'machining surface' gives the expected free-form surface [104].

In practice there are a number of design tools which help visualization of any free-form prior to machining. One technique is the use of visual reality at the design stage.

The CAD/CAM systems today have a very limited capability to allow the designer to perform conceptual designs. This is because most CAD/CAM systems require exact geometric specifications including shapes and dimensions. In a new development virtual reality is now presented [105]. A haptic device can be worn by the designer to generate virtual free-form surfaces with virtual tools in a virtual environment presenting the designer with a degree of freedom previously considered impossible.

## 6.7 Surface integrity

### 6.7.1 Surface effects resulting from the machining process

It is well understood that geometry changes in the workpiece are not the only effects produced on the surface. Most workpieces are produced by material removal processes such as single-chip removal methods as in turning, milling or drilling, or abrasive processes such as in grinding, polishing, honing, etc. Quite generally, the metal removal produces a characteristic surface alteration. This surface makes the surface layer quite different from the interior or bulk material. The surface alteration may be considerable and quite deep or it may be superficial and very shallow; in fact confined to the outermost skin. Surface integrity is a relatively new term that covers the nature of the surface condition. The term was introduced by Field and Kahles [105].

The reason for introducing surface integrity into this chapter is because it cannot realistically be omitted. It is virtually impossible to separate out the geometric properties from those of the surface. In many contact phenomena this has already been realized, for example in the use of the plasticity index. This is a

non-dimensional parameter made up from both types of geometric and physical phenomena. This and other examples will be illustrated in the next chapter on function.

The identification of surface layer alterations is important when the workpiece operates under high stress, especially alternating, or at high temperatures or in a corrosive atmosphere.

The basic idea and model for machinery used here is derived from that of Von Turkovich [107].

When a sharp, hard tool as, for example, in cutting is moved through the surface, a small plastic deformation zone is formed. When that zone reaches a stable state, a chip is usually formed. According to reference [107], there are two basic considerations: the mechanical parameters and the material parameters. Obviously, both need to be defined. The definition will clearly depend on the actual process—in fact it will be different for single-point cutting and, say, grinding. However, as a start, consider cutting. Any specification of mechanical parameters would have to include tool rake, wedge angle, type and shape of the cutting—whether it is double or single facet in diamond turning, for example — feed, depth of cut and speed and sometimes the length of the chip-tool interface.

Specifying the material properties is difficult and should include two factors: first, the cutting forces and chip morphology; second, the basic metallurgical state of the new surface generated (which can include mechanical effects).

The overall dynamics of the process and the metallurgical responses are overwhelmingly influenced by the cutting forces and the mechanism of chip formation. Since the cutting forces represent the sum of what acts on the boundaries of the plastic zone, they provide a direct link with the material properties in terms of stress/strain rate and local temperatures. Material crystalline structure, grain size, chemical composition inclusions and impurities are the main parameters to influence the deformation response. This response manifests itself in subsurface dislocation movement. The dislocation interacts with all types of defects, grain boundaries and voids. These micromechanisms constitute the basis of the understanding of strain hardening and fracture. Since in almost all materials the internal structure can be effectively modified by heat treatment, it is used in a number of ways to change the mechanical properties under the surface.

### 6.7.2   Surface alterations

When a surface is produced by conventional or non-traditional methods, many surface alterations are produced. These include:

(1)  plastic deformation
(2)  plastically deformed debris
(3)  laps, tears, cracks
(4)  microcracks
(5)  selective etch
(6)  intergranular attack
(7)  untempered martensite
(8)  overtempered martensite

and other effects dependent on the materials.

The principal causes of surface alterations are:

(1)  high temperatures or high temperature gradients developed in the removal process;
(2)  plastic deformation and plastically deformed debris;
(3)  chemical reactions and subsequent absorption into the surface.

The effects described above in the surface are shown in table 6.18 for a number of materials. This also shows microhardness alterations and residual stress introduced into the surface. Other properties are given in tables 6.19 and 6.20.

**Table 6.18** Summary of possible surface alterations resulting from various metal removal processes (after Von Turkovich [107]).

| Material | Conventional metal removal methods | | Non-traditional removal methods | | |
|---|---|---|---|---|---|
| | Milling, drilling or turning | Grinding | EDM | ECM | CHM |
| Steels: | | | | | |
| non-hardenable | R | R | R | R | R |
| 1018 | PD | PD | MCK | SE | SE |
| | L&T | | RC | IGA | IGA |
| hardenable | R | R | R | R | R |
| 4340 | PD | PD | MCK | SE | SE |
| D6ac | L&T | MCK | RC | IGA | IGA |
| | MCK | UTM | UTM | | |
| | UTM | OTM | OTM | | |
| | OTM | | | | |
| tool steel | R | R | R | R | R |
| D2 | PD | PD | MCK | SE | SE |
| | L&T | MCK | RC | IGA | IGA |
| | MCK | UTM | UTM | | |
| | UTM | OTM | OTM | | |
| | OTM | | | | |
| stainless (martensitic) | R | R | R | R | R |
| 410 | PD | PD | MCK | SE | SE |
| | L&T | MCK | RC | IGA | IGA |
| | MCK | UTM | UTM | | |
| | UTM | OTM | OTM | | |
| | OTM | | | | |
| stainless (austenitic) | R | R | R | R | R |
| 302 | PD | PD | MCK | SE | SE |
| | L&T | | RC | IGA | IGA |
| precipitation hardening | R | R | R | R | R |
| 17-4PH | PD | PD | MCK | SE | SE |
| | L&T | OA | RC | IGA | IGA |
| | OA | | | | |
| maraging (18%Ni) | R | R | R | R | R |
| 250 grade | PD | PD | RC | SE | SE |
| | L&T | RS | RS | IGA | IGA |
| | RS | OA | OA | | |
| | OA | | | | |
| Nickel and cobalt based alloys* | | | | | |
| inconel alloy 718 | R | R | R | R | R |
| Rene 41 | PD | PD | MCK | SE | SE |
| HS31 | L&T | MCK | RC | IGA | IGA |
| IN-100 | MCK | | | | |
| Titanium alloy: | | | | | |
| Ti-6Al-4V | R | R | R | R | R |
| | PD | PD | MCK | SE | SE |
| | L&T | MCK | RC | IGA | |

*continued*

Table 6.18 (*continued*)

| Material | Conventional metal removal methods | | Non-traditional removal methods | | |
|---|---|---|---|---|---|
| | Milling, drilling or turning | Grinding | EDM | ECM | CHM |
| Tungsten (pressed and sintered) | | | | | |
| TZM | R | R | R | R | R |
| | L&T | MCK | MCK | SE | SE |
| | | | | MCK | MCK |
| | | | | IGA | IGA |

Key: R   - roughness of surface
PD  - plastic deformation and plastically deformed debris
L&T - laps and tears and crevice-like defects
MC  - microcracks
SE  -selective etch
IGA  -intergranular attack
UTM -untempered martensite
OTM - overtempered martensite
OA  - overaging
RS   - resolution or austenite reversion
RC   - recast, respattered metal or vapour-deposited metal

**Table 6.19** Surface hardness changes that may result from various metal removal processes

| Material | Conventional (milling, drilling, turning or grinding) | | EDM | | Non-traditional process processes | |
|---|---|---|---|---|---|---|
| | | | | | ECM or CHM | |
| | Surface alteration | Hardness change | Surface alteration | Hardness change | Surface alteration | Hardness change |
| Steels | PD | Increase | RC | Increase | None | Decrease |
| | UTM | Increase | UTM | Increase | | |
| | OTM | Decrease | OTM | Decrease | | |
| | RS | Decrease | RS | Decrease | | |
| | OA | Decrease | OA | Decrease | | |
| Nickel- and cobalt-based superalloys | PD | Increase | RC | Increase | None | Decrease |
| Titanium alloys | PD | Increase | RC | Increase | None | Decrease |
| Refractory alloys (TZM, tungsten) | | | RC | No change | None | Decrease |

Key: PD    – plastic deformation and plastically deformed debris
UTM – untempered martensite
OTM – overtempered martensite
RS    – resolution or austenite reversion
OA    – overaging
RC    – recast, replattered metal or vapour-deposited metal

**Table 6.20** Comparison of depth of surface integrity effects observed in material removal processes (after Von Turkovich [107]).

| Property and type of effect | Condition | Maximum observed depth of effect* (in) | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Turning or milling | Drilling | Grinding | Chemical machining | Electro-chemical machining | Electro-chemical grinding | Electrical discharge machining | Laser beam machining |
| Mechanical altered roughing material zones: | | | | | | | | | |
| Plastic deformation (PD) | Finishing† | 0.0017 | 0.0008 | 0.0003 | ... | ... | ... | ... | ... |
| | Roughing‡ | 0.0030 | 0.0047 | 0.0035 | ... | ... | ... | ... | ... |
| | Finishing | 0.043 | 0.020 | 0.000 | ... | ... | ... | ... | ... |
| | Roughing | 0.070 | 0.110 | 0.000 | ... | ... | ... | ... | ... |
| Plastically deformed debris (PD$^2$) | Finishing | ... | ... | 0.0005 | ... | ... | ... | ... | ... |
| | Roughing | ... | ... | 0.0013 | ... | ... | ... | ... | ... |
| | Finishing | ... | ... | 0.013 | ... | ... | ... | ... | ... |
| | Roughing | ... | ... | 0.003 | ... | ... | ... | ... | ... |
| Hardness alteration§ | Finishing | 0.0005 | 0.0010 | 0.0015 | 0.0010 | 0.0014 | 0.007 | 0.0010 | — |
| | Roughing | 0.0050 | 0.0200 | 0.0100 | 0.0031 | 0.0020 | 0.0015 | 0.0080 | — |
| | Finishing | 0.013 | 0.025 | 0.038 | 0.025 | 0.000 | 0.018 | 0.000 | — |
| | Roughing | 0.127 | 0.503 | 0.254 | 0.079 | 0.051 | 0.008 | 0.200 | — |
| Microcracks or macrocracks | Finishing | 0.0005 | 0.005 | 0.0005 | ... | 0.0003 | 0.0000 | 0.0005 | 0.0005 |
| | Roughing | 0.0015 | 0.0015 | 0.0090 | ... | 0.0015 | 0.0010 | 0.0070 | 0.0040 |
| | Finishing | 0.013 | 0.013 | 0.013 | ... | 0.003 | 0.000 | 0.015 | 0.015 |
| | Roughing | 0.038 | 0.000 | 0.229 | ... | 0.0000 | 0.000 | 0.170 | 0.100 |
| Residual stress‖ | Finishing | 0.0060 | — | 0.0005 | 0.0010 | 0.0000 | 0.0000 | 0.0020 | 0.0002 |
| | Roughing | 0.0140 | — | 0.0125 | 0.0010 | 0.000 | 0.0000 | 0.0030 | — |
| | Finishing | 0.152 | — | 0.013 | 0.025 | 0.000 | 0.000 | 0.001 | 0.000 |
| | Roughing | 0.358 | — | 0.318 | 0.025 | | 0.000 | 0.076 | — |
| Metallurgical altered material zones: | | | | | | | | | |
| Recrystallization | Finishing | — | — | 0.005 | ... | ... | ... | ... | ... |
| | Roughing | — | — | — | ... | ... | ... | ... | ... |
| | Finishing | — | — | 0.013 | ... | ... | ... | ... | ... |
| | Roughing | — | — | | ... | ... | ... | ... | ... |
| Intergranular attack (IGA) | Finishing | ... | ... | ... | 0.0003 | 0.0003 | 0.000 | ... | ... |
| | Roughing | ... | ... | ... | 0.0060 | 0.0015 | — | ... | ... |
| | Finishing | ... | ... | ... | 0.0000 | 0.000 | 0.000 | ... | ... |
| | Roughing | ... | ... | ... | 0.152 | 0.000 | — | ... | ... |
| Selective etch, pits, protuberances | Finishing | 0.0004 | — | 0.0002 | 0.0006 | 0.0004 | 0.0001 | 0.0005 | — |
| | Roughing | 0.0010 | 0.3000 | 0.0004 | 0.0015 | 0.0025 | 0.0005 | 0.0016 | — |
| | Finishing | 0.010 | — | 0.005 | 0.015 | 0.010 | 0.000 | 0.013 | — |
| | Roughing | 0.025 | 0.076 | 0.010 | 0.036 | 0.000 | 0.013 | 0.004 | — |
| Metallurgical transformations | Finishing | 0.0004 | 0.0015 | 0.0005 | — | 0.0000 | 0.0001 | 0.0006 | 0.0006 |
| | Roughing | 0.0030 | 0.0200 | 0.0060 | — | 0.0002 | 0.0003 | 0.0050 | 0.0015 |
| | Finishing | 0.010 | 0.030 | 0.013 | — | 0.000 | 0.000 | 0.016 | 0.010 |
| | Roughing | 0.070 | 0.500 | 0.152 | — | 0.000 | 0.000 | 0.127 | 0.030 |
| Heat-affected zone (HAZ) or recast layers | Finishing | 0.0001 | — | 0.0007 | ... | ... | ... | 0.0006 | 0.0006 |
| | Roughing | 0.0010 | 0.0030 | 0.0125 | ... | ... | ... | 0.0050 | 0.0015 |
| | Finishing | 0.003 | — | 0.018 | ... | ... | ... | 0.015 | 0.015 |
| | Roughing | 0.025 | 0.070 | 0.318 | ... | ... | ... | 0.127 | 0.038 |

A dash (—) in the table indicates no or insufficient data.

An ellipsis (...) in the table indicates no occurrences or not expected.

* Normal to the surface.

† Depth to point where hardness becomes less than ±2 points *R* (or equivalent) of bulk material hardness (hardness converted from Knoop microhardness measurements).

‡ Depth to point where residual stress becomes and remains less than 20 ksi (138 MPa) or 10% of tensile strength, whichever is greater.

§ Finishing, 'gentle' or low-stress conditions.

‖ Roughing, 'off-standard' or abusive conditions.

### 6.7.3 Residual stress

#### 6.7.3.1 General

This is introduced into the surface by the material removal process. The presence of residual stress can be useful or disastrous depending on the application. The stress may be tensile or compressive, high or low.

Some more recent results are shown in table 6.21

**Table 6.21**

| Workpiece | Process | Surface alteration | | Reference |
|---|---|---|---|---|
| | | Nature | Extent | |
| **Metals** | | | | |
| S50C steel | Grinding | Amorphous, martensitic and over-tempered layer | 30 µm | Semba (1989) |
| Case hardened steel | Hard turning | White layer | 4 µm | Tönshoff (1995) |
| Hardened steel | High speed milling | Damaged layer | 4–6 µm | Elbestawi (1997) |
| Impax, C35, Armco | EDM | White layer | 80–100 µm | Kruth (1995a,b) |
| High strength steel | EDM | Residual stress | 400 µm | Mamalis (1988) |
| OFHC Cu, brass, 6061 Al | Diamond turning | Plastic deformation | 1–1.5 µm | Evans (1987) |
| Al-Mg alloy | Diamond flycutting, 5 µm | Residual stress | 1.8–17 µm | Horio (1992) |
| AlSi9 | Dry turning | Pitting of surface | | Byrne (1997) |
| **Ceramics and glasses** | | | | |
| $Si_3N_4$ | Magnetic float polishing | Roughness <10 nm $R_a$ | | Komanduri (1996) |
| Alumina | Grinding | Residual stress | | Tönshoff (1996) |
| $Al_2O_3/ZrO_2$ ceramic | Grinding | Residual stress | ~ 25 µm | Wobker (1994) |
| Alumina | Grinding | Roughness, microcracks, twin/slip bands | < 30 µm | Xu (1996b) |
| Silicon nitrides | Grinding | Median type cracks | 20–40 µm | Xu (1996a) |
| SiC | Ultraprecision grinding | Roughness <5 nm RMS | | Suzuki (1995) |
| Ceramics | ELID grinding | Roughness 3.2 nm $Ra$ | | Bandyopadhyay (1996) |
| Optical glasses | Surface grinding | Roughness 0.2 nm RMS | | Namba (1993) |
| Optical glasses | Float polishing | Roughness 0.2 nm RMS | | Namba (1987) |
| **Single crystal materials** | | | | |
| (111)CdTe | Lapping 47 µm, 28 µm | Microcracks | 50, 15 µm | Weirauch (1985) |
| (100)Si | Lapping 10–16 µm | Microcrystallites | 1 µm | Verhey (1994) |
| Sapphire | Grinding | Cracks, dislocations, strained regions | 10 µm | Black (1997a,b) |
| {111}Si | Rotation grinding | Plastic deformation | 2 µm | Tönshoff (1994) |
| Si | Ultraprecision grinding | Dislocation layer | 500 nm | Abe (1993) |
| Si | ELID grinding | Cracks | 400 nm | Ohmori (1995) |
| (100)Si | Diamond turning 1.27–50.4 µm | Dislocation loops | 1–3 µm | |
| | | Slip planes | 1–3 µm | Kunz (1996) |
| (001)Si | Scratching w/stylus | Amorphous structure | 150 nm | Minowa (1992) |
| (001)Si | Diamond turning 100 nm, 500 nm | Amorphous structure | 150 nm | |
| | | Dislocations, microcracks | 2–3 µm | Shibata (1994) |

*(continued)*

**Table 6.21** (*continued*)

| Workpiece | Process | Surface alteration Nature | Extent | Reference |
|---|---|---|---|---|
| (001)Si | Diamond turning 120 nm | Dislocation loops | 100–400 nm | Puttick (1994) |
| (0001)CdS | Diamond turning 1.25 μm | Lattice disorder | 200–360 nm | Lucca (1996) |
| (111)CdTe | Polishing 5 μm, 0.3 μm abrasive | Dislocations | 9 μm,3 μm | Weirauch (1985) |
| CdS, ZnSe | Mechanical polishing 0.25 μm diamond | "Obstruction-type" defects | 105 nm, 377 nm | Lucca (1997) |
| SiO$_2$ | CMP | Chemical/structural modification | 100–200 nm | Trogolo (1994) |
| Sapphire | CMP | Dislocation loops | 100 nm | Black (1997a,b) |
| CdTe | CMP 50 nm alumina | Dislocations | 30–500 nm | Nouruzi-Khorasani (1990) |
| CdS, ZnSe | CMP | "Obstruction-type" defects | 49nm, 135nm | Lucca (1997) |
| (100)InP | CMP | Dislocations and cracks | 50 nm | Laczik (1996a.b) |

In general, the following comments can be made regarding the degree of residual stress and the level to which it is introduced. The issues are (i) the sign of the stress, whether tensile or compressive; (ii) the depth in the surface to which it extends; and (iii) its maximum value. Some examples follow.

When using conventional grinding, a resultant tensile-stressed layer is introduced. If the grinding is forced, as in 'abusive grinding', the resultant stress of the surface can be very high. The tensile stress can extend to 0.05mm in depth and peak out at a value of 690 MPa. For conventional grinding stress the tensile stress is still present but at a much lower level. In fact the total extent of the grinding would extend to a depth of 0.01 mm. That grinding does produce tensile stress seems to be accepted; why it does is still open to question but seems to be mostly determined by temperature effects. For steels, the effects produce either untempered or overtempered martensite in the surface layer.

Another typical example of a mechanical influence rather than the thermal or chemical shock effects is found in the shot peening process which always tends to produce compressive stresses at the surface, at least under the direct shot position [109, 110], for example. The penetration of the residual stresses depends on



**Figure 6.107** Compressive stresses produced in shot peening.

the peening conditions and can be very deep; as shown clearly by the work of Brinksmeier *et al* [111], a depth figure of 0.5mm is not unusual. Because peening is a 2D random process with no preferential directioning, it produces not only an isotropic surface roughness, but also an isotropic stress pattern in the surface. An outline figure of compressive stress as a function of depth is shown in figure 6.107.

Another example of a mechanical shock pointed out by Brinksmeier *et al* [111] is the rolling process, which differs from peening in the sense that it is definitely preferential both in the texture and the stress distribution. Rolling produces a texture on sheet steel that is suitable for use in plating and painting. As an example of rolling it has been suggested [111] that in order to increase the stability of a circular saw blade the disc is rolled at a given radius (figure 6.108). Sometimes, such as in the normal operation of a component (e.g. a roller bearing), the very nature of its operation is compressive. This has the effect of driving the compressive stress limits lower through the skin of the surface [112].

A rather more obvious compressive process which is not concerned with function of the part but more with preparation of the part for work is in tumbling to remove burrs.



**Figure 6.108** Stresses produced in rolling.

An example of thermal impact on the surface is the use of EDM. This invariably leads to tensile stresses in the subsurface. The depth of penetration depends on the discharge energy [113, 114], as shown in figure 6.109.



**Figure 6.109** Stresses produced by EDM.

Chemical impact also has the potential for introducing stresses. However, there is little concrete evidence at this stage as to the absolute values. Note the important distinction between EDM and ECM. The latter does not produce internal stresses; this is why ion milling is very often used to remove the uppermost layers of EDM.

This breakdown of stress-forming properties into mechanical, thermal and chemical is, for many processes, a little arbitrary because, more often than not, more than one is involved simultaneously. The combination most often encountered is mechanical together with thermal.

According to Brinksmeier *et al* [111], one rather simple model due to Syren [115] enables some attempt to be made to predict the residual stress produced by machining. In this model a distinction is drawn between surfaces which are cut by the tool without a following squeezing and those which are squeezed and, therefore, plastically deformed after the cutting process. Thus, surfaces generated by a cutting operation without a following squeeze have tensile stresses. Compressive stresses result if a squeezing follows. Also, high cutting temperatures favour tensile stresses for the obvious reason of the subsequent inhibiting of contraction.

Blunt cutting tools produce more deformation and therefore sometimes help towards compression. However, this situation is often complicated by the fact that if this condition is carried too far, or not enough attention is paid to actual shape, an increase in frictional forces results, which introduces tension behind the tool movement. The model has been used effectively for turning, upcut milling and downcut milling but has not been especially suited to explain grinding.

This particular problem has been addressed by a Leuven team [116], based on the classic work of Jaeger [117]. This treatment gives a relationship between the grinding parameter and the surface temperature of the workpiece. The former is governed by the variations in the tangential grinding force which, together with the grinding speed, determines the specific energy of grinding.

From the equations derived, the force component is a very important factor: it largely determines the surface temperature, which in turn has the greatest effect on tensile stresses. Some specific examples follow, summarizing relevant points from the reviews in references [111] and [118].

### 6.7.3.2 Grinding

Grinding produces high normal forces because of the shape of the abrasive grains. These forces are very large compared with, say, turning, because the individual grain cutting in grinding is negative rake. Inevitably, very high pressures and temperatures are generated. Summarizing, the residual stresses are produced by:

(1) machine conditions—these include depth of cut $a$, workpiece speed $V_{ft}$, cutting speed $V_c$,
(2) topography of the grinding wheel—the dressing and wheel behaviour;
(3) type and specification of the wheel;
(4) coolant.

In fact, if $a$ is the depth of cut and $d_c$ the equivalent wheel diameter given by

$$aV_{ft}/V_c \tag{6.42}$$

then they can be used to help predict surface temperature produced in the grinding process.

Thus force (tangential) is given according to [115] by

$$F_t = F_0 h_{eq}^f \left(\frac{V_c}{V_{ft}}\right)^{0.1} d_e^{(1-f)} \tag{6.43}$$

where $f$ is called the force exponent and takes values between 0.3 and 1.0. If $f$ has been determined from equation (6.43) by measuring $F_t$ and knowing $F_0$ as a function of the grinding conditions, an estimation of surface temperature can be made. Thus

$$T \sim a^{f-0.235} V_t^{f-0.57} V_c^{1.1-f} d_e^{0.765-f}. \tag{6.44}$$

From (6.44) it can be seen that $f$ is very important in determining the temperature. Using this equation and putting in realistic working values for the parameters shows that temperatures in the region of 1000°C can easily be obtained. As a rule of thumb, if $f$ is small (the force exponent) then high temperatures result and high tensile stresses are produced. For carbon steel (0.15% carbon), the residual stress introduced as a function of temperature is more or less linear [116] (figure 6.110).

Different materials are affected in different ways by the grinding parameters. For example, titanium surface stresses seem to be mainly dependent on the wheel speed.

The dressing of the wheel is obviously important because the shape of the grain determines the relative compressive and tensile stresses. When the wheel has just been dressed the stresses are relatively low, but as the wheel becomes blunt the tensile stresses increase dramatically. The coarser the dressing the lower the initial stress because of the openness of the wheel topography generated. Making a fine dressing increases the tensile character of stresses, but at the same time the wheel produces a lower surface roughness, so wheel dressing is a parameter in which two functional characteristics conflict.



**Figure 6.110** Residual stress as a function of temperature for steel.

Not only is the condition of the wheel important but also the grain, the bond structure and the hardness. A soft wheel can require more dressing time but it can be advantageous because less stress is introduced. This is because a grain will simply break out if the frictional force is too high, so temperatures do not get the chance to rise. As an example of this effect of grain breakdown, cubic boron nitride wheels behave in a different way from aluminium oxide wheels. In the former the grain blunts, but in the latter it splinters—this increases the sharpness of cut, so less heat is generated and on the average compressive stresses result.

There is direct evidence of a correlation between grinding energy and residual stress. One example is shown in figure 6.111.

Grinding fluids can also have an effect on stresses since dry conditions are rarely used and it has been shown that the presence of oil reduces the friction and hence stresses more than water emulsions.

There are several well-known methods for improving the residual stress conditions of ground surfaces. These usually rely on finishing processes and are usually based on plastic *movement* of material rather than metal removal. One of these techniques is spark-out [120]. This has the dual advantage that the surface roughness is also improved at the same time as reducing tensile stresses.

**Figure 6.111** Correlation of residual stresses and grinding energy per unit area [119].



**Figure 6.112** Honing as a means of introducing compressive stresses.

Other methods include shot peening or abrasive tumbling. Honing has also been used to introduce compressive stresses (figure 6.112) [121], especially in the ball bearing industry,

### 6.7.3.3 *Turning*

In normal turning, the problem of heat is less severe than in grinding for a number of reasons, the main one being that the majority of the cutting energy is transformed into heat in the shearing zone where the chip is produced. This has the result that heat escapes readily into the bulk of the workpiece and also gets into the chip, which is an excellent way for dissipation because the thermal capacity is so much larger than for grinding chips. Tensile stresses are dominant in turning because the surface is caused by the secondary cutting edge rather than the primary [122]. A smaller chip and negative cutting angles give higher friction, which would result in higher tensile stress except for the fact that the negative rake produces more plastic flow and consequently the stress is minimized [123].

The thermal impulse increases with cutting feed so that as the cutting feed increases, the stresses tend to be driven further into the surface.

*Turning versus grinding; residual stress [124]*

Because of the interest in turning as a substitute for grinding a comparison of residual stress produced by the processes has been made on bearings both large and small [124]. Here the term hard turning is usually reserved for CBN tools cutting hard metals of high Rockwell number.

Superfinishing only changes the residual stress near to the surface and is not important for deep stresses. Steels used Rockwell 58–62 case carburised.



**Figure 6.113**

The basic conclusions were:

1. Hard turned and superfinished bearings have at least as long fatigue life as ground and superfinished.
2. Depth of compressive residual stress is the major difference. (See figure 6.113).
3. Feed rate only changes residual stress near to the surface and not in deep layers.
4. Tool edge geometry is the dominant factor. (See figure 6.114).



**Figure 6.114** Effect of tool edge.

### 6.7.3.4 Milling

Considerable plastic deformation is produced so the stresses are often compressive (e.g. [123]). Kiethe [123] shows dramatically how different stress patterns are produced in upcut versus downcut milling (figure 6.115).



**Figure 6.115** Stress patterns produced by different directions of milling.

### 6.7.3.5 Shaping

The stress distribution is the same as that in turning, the essential difference between shaping and milling being the continuous cut. Tensile stresses are produced once chip formation begins.

Of the other processes perhaps it should be mentioned that most of them are not usually finishing processes but do involve a degree of formation of the surface. Among these are rolling and extrusion. Both produce compressive stresses in the surface as mentioned before, both being directional, yet shot peening produces isotropic stresses.

### 6.7.3.6 General comment

The real issue in this section has been the generation of residual stresses in the surface by the machine process. Is this really important? Unfortunately it is. It is simply no good forgetting the presence of physical influences other than the geometrical. In particular, high cycle fatigue properties show the greatest sensitivity to the manufacturing process and, because fatigue behaviour is concerned with non-geometric aspects, it will be mentioned here.

Although the causes of fatigue strength are not fully understood, consideration of the process to ensure that 'non-abusive' conditions are avoided can make a considerable improvement.

The two main factors influencing fatigue are given below:

1. The presence of tensile stresses is detrimental and those of compressive stresses favourable.
2. Microcracks introduced by the process, such as those invariably produced in EDM, and cracks produced by built-up edge or selective etching.

A comprehensive list of high cycle fatigue strengths of a wide variety of alloys etc is produced by Metcut research [126].

### 6.7.4   Measurement of stresses

#### 6.7.4.1   General

The measurement of stress is not a main theme of this book, so for this reason it can only merit a mention. Basically there are two approaches, direct and indirect. It has been suggested that a better breakdown could be achieved by splitting the methods into destructive and non-destructive. This approach has been questioned because some interference of the part is inevitable.

For indirect methods, the equilibrium of forces has to be disturbed so that the resulting deformations of the body can be deduced. For direct methods, the idea is to measure those physical properties of the body that can be influenced by stress levels. In effect, they are both ways of looking at an indirect method.

Techniques for surface/subsurface properties. A relatively full list of techniques is shown below [108].

<div align="center">

**Techniques**

**Acoustic techniques**
Acoustic emission sensing
Photoacoustic microscopy
Quantitative acoustic microscopy
Scanning acoustic microscopy
Surface Brillouin scattering
**Electron beam techniques**
Cross-sectional transmission electron microscopy
High resolution electron energy loss spectroscopy
High resolution transmission electron microscopy
Low energy electron diffraction
Reflection high energy electron diffraction
Scanning electron microscopy
**Ion beam techniques**
Cathodoluminescence
Ion channelling
**Luminescence spectroscopy**
Photoluminescence
Rutherford backscattering spectrometry
**Magnetic techniques**
Barkhausen effect/micromagnetic techniques
Magnetic article inspection
Magnetic leakage field testing
**Microindentation/nanoindentation**
**Optical scattering techniques**
Optical diffraction
Polarized light microscopy
Polarized light scattering
Speckle interferometry
Total internal reflection microscopy
**Photothermal techniques**
**Raman spectroscopy**
**Scanning probe techniques**
Electrostatic force microscopy
Kelvin probe force microscopy
Near field scanning photoluminescence

</div>

Scanning capacitance microscopy
Scanning near field optical microscopy
**X-ray scattering techniques**
Double axis (high resolution) x-ray diffraction
Double crystal x-ray topography
Grazing incidence x-ray diffraction
Grazing incidence x-ray reflectivity
Grazing incidence x-ray topography
Reciprocal space mapping
Single crystal x-ray topography
Triple axis x-ray diffraction
Triple axis x-ray topography
X-ray diffraction
X-ray scattering techniques.

These are listed alphabetically and not in order of importance. They do not include some of the very well-known methods such as microindentation with stylus methods or optical microscopy.

The methods themselves will not be described in great detail: more emphasis will be given to what is measured. Details of the methods have been covered in chapter 4.

*6.7.4.2 Indirect methods*

The first idea was to measure the deformation of parts. This method was not taken up in practice because of the errors introduced by the stresses still in the removed body [125].

The problem is eliminated by measuring the deflection of the remaining part, the idea being that once the stressed part has been removed the remainder will resort to its stable state. The part is made single so that the stresses can be calculated from the deflections of the stressed part. Different sorts of body have been used; these include cylinders [128] and rectangular bodies [127].

For the practical application for plates and cylinders, if machining stresses have to be determined the surface layers of the specimen are removed at the machined side while the deformations of the remaining part have to be measured as a function of thickness [127] (figure 6.116).

If the principal directions of stress are known, the measurements of strain in these directions are sufficient to get the biaxial stress. Accordingly, given the strains $\varepsilon$ as a function of $x$ and $y$, the stress $\varepsilon$ can be expressed [71] as

$$\sigma_{xy} = \frac{-E}{1-v^2}\left[\frac{1}{3}z\left(\frac{d\varepsilon_{x,y}}{dz} + \frac{vd\varepsilon_{x,y}}{dz}\right) + \frac{4}{3}(\varepsilon_{x,y} + v\varepsilon_{x,y}) - \frac{2}{3}\int_z^W\left(\frac{\varepsilon_{x,y} + v\varepsilon_{x,y}}{\omega}dz\right)\right] \qquad (6.45)$$



**Figure 6.116** Strain measurement as a function of thickness.

This is in terms of the specimen thickness $\omega$ and the depth between the surfaces $(W-z)$.

When removing the surface layer, it is vital not to introduce any stresses. For this reason, electrochemical polishing is most often used, which produces little or no stress. Elaborate schemes have been adopted for better estimates of stresses using specifically shaped components [127]. However, the methods are all substantially the same.

### 6.7.4.3    Direct methods

The best-known method, apart from the standard way using fluorescent dyes, which will not be discussed here, is the use of X-rays. Other techniques which are being used but on a smaller scale are the magnetic and ultrasonic methods. Electromagnetic methods are based on the stress dependence of electrical conductivity and the magnetic effects of ferromagnetic materials [128]. The ultrasonic methods depend on the fact that the propagation of ultrasonics in a solid is influenced by the presence of strain in the material [129, 130].

The advantages of these two methods are that they are non-destructive and can be carried out in a few seconds. Notwithstanding this, they need a lot of experience to be carried out properly because, for example, in ultrasonics all sorts of things other than stresses influence the passage of the waves. Things such as local structure changes in density and in the electrical and magnetic properties of X-rays have their own influence on the outcome.

In X-ray methods, it is strain that is actually measured. Furthermore, it is *lattice* strain that is assessed rather than macroscopic strain.

The following is the method explaining the Bragg diffraction and shown pictorially in figure 6.117. Here the diffraction equation is given by

$$\lambda = 2d \ \sin\theta \tag{6.46}$$

where $d$ is the lattice spacing, $\theta$ is the diffraction angle and $\lambda$ is the X-ray wavelength.

Figure 6.118 shows what happens when a crystal is strained. The new lattice spacing in the vertical direction is $d'$. The diffraction angle is obviously changed due to the strain because $\alpha$ is not equal to $\theta$. In fact it is larger and hence by measuring the change in diffraction angle the strain change $d$ to $d'$ is observed. This single-crystal picture is due to Aksenov [129].

The situation regarding multicrystals is more complicated. An idea of the complexity can be seen from figure 6.117. The problem basically is that the grains exhibiting strain invariably are not parallel to the face of the crystal, so the grain spacing has to be determined at an angle if using the nomenclature of Tonshoff in an excellent survey [130]. For angles $-50° < \psi < +50°$, the strains for $\psi = 90°$ can be obtained from these by interpolation.



**Figure 6.117** X-ray diffraction at a single unloaded and loaded crystal.

**Figure 6.118** Effect of stress on polycrystalline lattice spacings: (*a*) unstressed polycrystalline specimen; (*b*) stressed polycrystalline specimen.

By setting the lattice strains equal to those which would result for biaxial surface stress, it can be shown that the lattice strains are a linear function of $\sin^2 \psi$. A number of assumptions have to be made and calibrations carried out before meaningful results can be found (e.g. [131, 132]).

However, having carried these out, the lattice strains $\varepsilon$ can be calculated using Bragg's equation. From this, very exact calculations of the peak positions from the diffracted intensities can be measured.

In practice it must be appreciated that for many manufacturing processes in which a high level of plastic deformation takes place, the linear dependence of $\sin^2 \psi$ does not occur, owing to triaxial stress states. It has also been suggested that the geometrical surface texture can affect the law, although up to the present nothing qualitative has been proved. The discrepancies are probably due to local surface slopes of the workpiece.

Measurement of stress as a function of depth into the surface are laborious because a certain amount of material removal has to take place—at least with today's methods. At present, layers are etched away step-wise and stress levels found at each step. Plotting the results as a function of depth enables the maximum stress and its depth to be found (figure 6.119).

An alternative method to the time-consuming step-by-step method is to use a differential etch method in which the specimen is placed at an angle in the electrochemical polishing bath relative to the cathode. The metal is then removed in a wedge, therefore enabling all the relevant metal to be removed in one go. Progressively moving the X-ray equipment along the slope obviously probes different depths in sequence (see [130]).

Comparison between the indirect and direct methods is possible although there are bound to be differences. For example, the indirect method is incapable of measuring local stress levels. This is obviously no problem using the X-ray method, although in terms of cost the indirect method is cheapest. It is definitely a possibility for the future that the X-ray direct method or even the ultrasonic method will be made 'in process'.



**Figure 6.119** Stress pattern as a function of depth.

### 6.7.5 Sub-surface properties influencing function

#### 6.7.5.1 General

Summarizing, apart from the geometrical aspects of surfaces, which influence the function dramatically in many respects, as will be shown in the next section, it is important to remember that there is another factor that substantially affects the function of the part. This factor is the physical characteristics of the surface. Prime amongst these characteristics is the residual stress introduced into the surface by the machining and the material processing [136].

Although this book is concerned with geometry, the function of the part is so integrally determined by both the geometry and stresses that some brief comments should be made about how they originate and their importance. As mentioned earlier, the residual stress is just one element in what has been referred to as 'surface integrity' defined by Field.

#### 6.7.5.2 Influences of residual stress

Residual stress can enhance or impair the functional behaviour of mechanical parts [110]. A simple way to express the formation of stresses is as a result of shock, be it mechanical, thermal or chemical, as shown in figure 6.120.

The residual stress comes in the physical category and figure 6.121 shows its main influences.



**Figure 6.120** Formation and influence of residual stress.

Economic and functional pressures demand that, in critical cases, not only should the surface geometry be specified but also the physical and, hence, residual stresses. In the latter case, the residual stress condition in the surface is of most importance so that its description does come under the umbrella of surface properties if not strictly as surface geometry.

For many applications, the surface is fundamental (figure 6.121). This figure is a crude breakdown of some of the influences. These are not always bad; sometimes there can be a good influence. Each one in the block diagram will now be described in turn.

*(a) Deformation*
Residual stresses act in a body without forces or moments, they internally form a self-contained system in equilibrium. If, for some reason, the body is altered in some way, for example by machining, the equilibrium

**Figure 6.121** Effect of residual stress.

of forces is disturbed and, in order to restore the balance, the body distorts, as for example in castings. If the casting is not heat treated for stress release—which allows the stresses to dissipate without changing the final shape of the body—the body can distort, especially if not machined symmetrically, when it is released from the machine tool clamp. It has been pointed out that this is the reason why the experienced operator releases the clamps after the roughing operation so that the deformed parts can be removed before finishing, otherwise what results is a component having a good finish but poor shape!

The amount of deformation of the part is usually assumed to be roughly proportional to the removed cross-section of material. The effect of deformation resulting from the removal of metal is not the same as the inputting of stress into the workpiece *by machining*. These residual stresses are produced by metallurgical transformations or plastic deformation and are restricted to the outer skin of the part, not all through as in casting.

Another way in which stresses are introduced other than by asymmetrical machining and metal removal in casting is by heat induction; they all produce a different pattern of stress in the material. Of particular importance is residual stress in sheet metal forming [133, 134].

*(b) Static strength*

From a macroscopic viewpoint, residual stress acts like a prestress state of the material. For materials that are deformable and have a characteristic yield point, the residual stresses influence the yield strength. Assuming a criterion for plastic flow, the elastic limit can be calculated if the stress distribution is given [135]. It can be demonstrated quite readily that the elastic limit can often be decreased by up to 40% from the unstressed case, which is a considerable reduction.

In pure tension or compression loading, the elastic limit has to be lowered by the presence of uniaxial residual stresses independently of their distribution. But in bending or multiaxial stress distribution, an *increase of* the elastic limit can be realized [110].

Increasing the strength of a component by prestressing is used in many technical applications, for example, extrusion dies. Usually the loaded component is an assembly in which prestress is used and which is made of a plastically deformable material. If the material is brittle, residual stresses may lead to catastrophic cracks if the resulting stresses exceed the strength of the material at any point. Such effects have been seen in crankshafts, for instance.

An example of a case where the residual stress can be an advantage is in the case of disc springs. In these a stress distribution is generated to prevent relaxation effects. What happens is that the spring is meant to transmit high forces at low spring deformations. The shape is shown in figure 6.122.

**Figure 6.122** Stress spring.

These springs are designed to deform elastically by a factor of 75% of the theoretically possible value of *h*. The acting stresses exceed the yield strength, which causes the spring to relax [137]. This relaxation can be avoided by imposing a residual stress into the spring.

*(c) Influence of residual stress on dynamic strength*
The influence of residual stress on fatigue strength has long been known (e.g. [138]) and has been proved many times since. The relationship between residual stress and fatigue strength is shown in figure 6.123 [139] which is typical for steel.



**Figure 6.123** Fatigue strength as a function of residual stress.

A very interesting relationship between fatigue strength, residual strength and surface texture has been shown by Field and Koster [138]. This, in simplified form, is shown in figure 6.124. All relate to longitudinal stress.



**Figure 6.124** Effect of surface finish on fatigue strength.

The figure shows clearly that, as expected, the surface roughness value affects the fatigue strength somewhat, that is the rougher the surface, the poorer the strength. The effect of the surface stress is much more important. Unfortunately, what it does not show is whether the surface parameter $R_a$ is the most suited to characterize the surface for the fatigue function. More will be said about this in surface function.

It should be added here that it is not only the residual stress or the surface roughness that influences fatigue strength; heat treatment plays a part, especially for high-carbon steels, but it is not so important for mild steel [114].

### *(d) Residual stress effect on chemical corrosion*

Components having residual stress and subjected to a corrosive environment can fail prematurely. The conditions [111] under which stress corrosion occurs are a result of a specific sensitivity of the metal, the presence of tensile stresses in the surface and the presence of a corrosive environment. There are different theories about the dominant process during stress corrosion. One, called the electrochemical hypothesis, assumes potential differences between the precipitations at the grain boundaries so that the precipitates are dissolved anodically. Then tensile stresses in the surface open up the cracks and cause rupturing of the surface film in the root of the crack, thereby bringing pure metal into direct contact with the hostile environment.

### *(e) Magnetic effects*

The magnetic properties of materials are directly affected by residual stress. Any mechanical strain inside the crystal causes flux changes. Indeed, the Barkhausen effect is used to measure residual stress.

For definitions of stresses and strains the reader is advised to consult any general textbook on the properties of materials (e.g. [139, 140]).

## 6.8   Surface geometry—a fingerprint of manufacture

### 6.8.1   General

The output of a machine tool is the workpiece. Hopefully this workpiece will have been produced to a specification which will fit a function. In many cases the geometry of the part—as made—will comply with the requirements of the function. Small deviations from it, in general, will be tolerated.

Such deviations are those which result from the inability of the machine tool and the process to generate exactly a continuous smooth shape. For this functional purpose they can therefore be regarded as irritants which can sometimes impair performance.

However, there is a positive use to which these geometrical deviations can be put. They can be used to monitor the manufacture, thereby recouping some of the cost of making the workpiece.

Any degeneration of the geometric deviations can have two implications. One is that the workpiece is no longer within the tolerance laid down by the function, that is the part is unsatisfactory. The second is that the manufacture is degenerating either due to process faults or machine tool problems. It is the second of these implications of degeneration that will be addressed here. Some indication of what can be done was seen in section 6.4.9.

The workpiece is the ideal place to look for problems in manufacture because it is where the most concentrated effort to control manufacture has been put. It is the output of the system where all forces are meant to balance. Small changes in either the manufacturing process or the machine tool will show themselves on the workpiece usually before they are apparent anywhere else. Unfortunately the actual size of the perturbation on the surface is small, but it is there. Also it has to be extracted from within the normal deviations produced by well-behaved machining.

In what follows some examples will be given of the potential use of random process analysis to extract this small but critical information in single-point cutting and grinding. After this the use of Wigner functions and other space-frequency functions will be explored. Finally, some thoughts on the use of non-linear dynamic systems will be touched upon.

## 6.8.2  Use of random process analysis [141]

### 6.8.2.1  On turned parts—single-point machining

Because the mark left on the surface is basically periodic it is best to use the power spectral density rather than the autocorrelation function. This choice makes the relevant deviations more 'visible'.

*(a) Tool wear*
Similarly the best 'unit' to work from is the spatial feed frequency of the cutting tool. Figure 6.125 shows the power spectrum of a turned part.

The impulse in the power spectrum representing the feed frequency is accompanied by a succession of harmonics at twice, three times the frequency and so on. These represent the deviations of the tool shape from a true sine wave. In the figure they are shown rather larger than actual size for visibility. The basic



**Figure 6.125**  Tool wear monitoring.

premise is that the surface profile acts as a negative replica of the tool tip in 'perfect' cutting. As the tool wears each chip repeats at the feed wavelength in the profile. This repetitive scar signal has the spectrum of a 'Dirac comb' which adds equally to the harmonics as well as the fundamental. The result is that the ratio of fundamental to harmonic amplitudes should change as a function of tool wear (figure 6.125).

Thus, wear on the tool tends to affect the spectrum above the feed frequency. Also, these effects are periodic in nature. There are, however, other small-wavelength effects, which are not periodic but random. One cause of these could be microfracture of the surface, which occurs sporadically within the tool feed mark. This could be an indication of workpiece material inhomogeneities or problems in chip formation. This random component shows itself as a continuous band of frequencies, above the feed frequency, in the spectrum.

*(b) Subharmonic effects*
So the tool wear potentially changes the harmonic ratio. These are all above the feed frequency.

Another phenomenon occurs in the subharmonic region. This is mainly due to machine tool faults such as bearing wear, or chatter due to low stiffness or, at the lowest frequencies, the slideway error (figure 6.126).



**Figure 6.126** Long-wavelength effect and total signature. A, subharmonic machine tool; B, harmonic tool wear; C, high-frequency material process parameters.

The nature of the longer wavelengths occurring in the subharmonic region of the power spectrum of the surface is complex. Some will be related to the feed frequency only lower, such as bearing vibration, and some will be of lower frequency and still periodic but not necessarily related to the feed. These could be due to the slideway traverse screw mechanism. Other frequency components could be completely unrelated to the feed owing to vibration external to the machine tool and yet visible because of the lack of stiffness problems. This could also allow regenerative chatter to be initiated. Hence, spectral information below the feed frequency is a valuable source of machine tool and environmental problems. Taken together with the spectral information above the feed frequency, it can be seen that a tremendous amount of information about the process and machine tool is available (figure 6.126). Furthermore, it is available on the surface. How to get it from the surface is another matter!

Spectral investigation can be carried even further in some ways if the actual Fourier spectrum is used rather than the power spectrum or the Hartley transform, in which case the phase information which is then preserved might be able to differentiate different types of tool wear (i.e. flank wear from cutting edge wear or cratering). Unfortunately the lack of reliability of such data probably rules out this possibility.

### 6.8.2.2  Abrasive machining

The question arises as to whether the 'fingerprint' concept can be used in abrasive processes. The answer is yes. The autocorrelation should be used here because it is best for revealing random detail. This applies most to the 'unit machining events' which are the key for the autocorrelation function in the same way that the tool feed is for the power spectrum.

Clean cutting from a single grain could be said to produce an impression in the surface as shown in figure 6.127.

The autocorrelation function can be used to get an idea of the size of that most elusive piece of information—the effective chip width. It should give a good idea of the effectiveness of the abrasive process. Unfortunately machine tool problems tend to add a repetitive factor which can mask the detail. The same can usually be said for grinding wheel dressing marks in the case of grinding. These tend to be periodic, although they should not be. The safeguard in abrasion therefore is to regard any major periodicity far from the origin as likely to be a fault.

### Comment

The surface roughness has always been regarded as being important in grinding because it is one of the most desired attributes of the component in terms of quality. However, over the years the accomplishment of this 'surface quality' has been judged by using one parameter, $R_a$ or worse, $R_z$ or $R_t$, none of which are known for effective characterization. This is a pity because it limits the realistic use of grinding models advocated by many workers.

That it has not been greatly successful is realized by seeing how many variants of the models have been proposed linking the texture to grinding parameters; equation (6.18) is one, others include

$$R_z = A_1 + A_2 a_e^{e_1} \left( \frac{1}{q} \right)^{e_2} V_c^{e_3} V_w'^{e_4} \tag{6.47}$$

and from a basic model

$$R_z = A + C_{wp} C_{gw} a_e^{e_1} \left( \frac{1}{q} \right)^{e_2} V_c^{e_3} V_w'^{e_4}. \tag{6.48}$$

**Figure 6.127** Abrasive correlation

All contain many constants that depend heavily on very specific conditions, that is $A_1$, $A_2$, $e_1$, $e_2$, $e_3$, etc. The common parameter wheel ratio $q$, the cutting speed $V_c$, the specific material removal $V'_w$, engagement $a$, wheel constant ($C_{wp}$), component constant ($C_{gw}$), etc, are all included and yet undergo changes in emphasis from author to author. That it is possible to predict a surface $R_a$ or $R_z$ from a great list of such specific data is highly probable. This data for many applications is obtainable from references [142–153] which include Koenig's excellent school of process investigators. However, it is disappointing that a more fundamental relationship between the surface texture, the wheel topography and the basic grinding parameters has not been found. This failure cannot have been helped by relying on the somewhat elementary concepts used for surface characterization. No use can be made of equations (6.16), (6.18), (6.47), etc, to investigate process or machine tool problems. *The surface texture data is being wasted*! More could be done, and should be done. The surface texture has to be produced on the surface anyway. This observation applies to many processes, not just grinding.

### 6.8.3 Space-frequency functions (the Wigner function)

In chapter 2 a number of transforms were introduced with the idea of improving on the Fourier transform as a means of extracting information from the surface geometry. Most are interesting but, on balance, do not seem to result in much better extraction. Some have potential benefits in speed but, in view of the fact that in some circumstances the Fourier transformation is carried out at the speed of light in the in-process methods illustrated in chapter 4, this does not seem to be a real advantage. What is required is a method which improves on the quality of the information obtained from the surface.

The only information lacking is that concerned with changes in the roughness over the workpiece. These changes can be in the areal statistics or in the presence and characterization of flaws. These requirements point to functions that have more than the one argument (say, spatial frequency). It seems logical to introduce position as the other variable—not by means of switching to the autocorrelation function but in addition to frequency. Immediately, the space-frequency functions—the Wigner function, the ambiguity function and possibly the wavelet function—come to mind. They have as a basis the Fourier kernel and involve space and frequency variables. The properties of these have been given in chapters 2 and 3, so they will not be given here. It should suffice to indicate that all the advantages of random process analysis are preserved [154, 155]. Either of these functions are valuable. It might even be argued that both the Wigner and the ambiguity functions could be used to advantage. However, to date only the Wigner function has been used on a practical case. This was to ascertain whether chatter could be identified at an early stage of turning using the roundness profile of the workpiece, $r(\theta)$.

It turns out [142] that the typical roundness signal in the presence of chatter in the $y$ and $z$ directions of $a_y, f_y, a_z, f_z$ is given by

$$r(\theta) = \frac{a_z^2}{2r_c} \cos^2\left[\frac{f_z\theta}{n_\omega} + \frac{f_z a_y}{n_\omega r_\omega} \sin\left(\frac{f_y\theta}{n_\omega} + \varphi_y\right) + \varphi_z\right] + \frac{a_z f_m\theta}{2\pi r_c} \cos\left[\frac{f_z\theta}{n_\omega} + \frac{f_z a_y}{n_\omega r_\omega} \sin\left(\frac{f_y\theta}{n_\omega} + \varphi_y\right) + \varphi_z\right] + \frac{f_m^2\theta^2}{8\pi^2 r_c}. \quad (6.49)$$

Putting in digital terms and realizing that the second term is the dominant one gives the digital roundness signal $f(n)$

$$f(n) = a\left\{\exp\left[j\left(\frac{2\pi k_0 n}{N} + \psi_z + b\,\sin\left(\frac{2\pi k_m n}{N} + \psi_y\right)\right)\right]\right\}. \quad (6.50)$$

The interesting point is that by taking the frequency moments of this, as explained in chapter 2, the main chatter signal in the $y$ direction is obtained easily, giving

$$a_y = \frac{n_\omega r_\omega b}{f_z} \qquad f_y = \frac{n_\omega L k_m}{N} \qquad (6.51)$$

where the symbols are feed rate $f_m$, workpiece rotation $n_\omega$, radius of the workpiece $r_\omega$, $a = f_m a_z/2r_c$, $b = f_z a_z/n_\omega r_\omega$, $k_0 = f_z N/n_\omega L$ and $L$ is data length.

Equation (6.51) has been verified with practical results [142], so this method appears to be viable for the future. In addition to chatter, the steady-state solution is easily obtained, that is

$$r(\theta) = f_m^2\big/8\pi^2 r_c \qquad (6.52)$$

which is the geometry without vibration.

Consequently, both sorts of information—steady state, corresponding to the normal cutting and transient, corresponding to chatter—can be obtained. This fills a major gap in the information required for machine tool monitoring.

The ambiguity function does not use position as its second argument; it uses size. Hence, although it cannot pinpoint the position of a scratch, say, as shown in chapter 2 for the Wigner function, it can scan the surface with a 'window' of given size and hunt for data of that size — it is a type of cross-correlation, but slightly different.

The difference can be seen by comparing the formulae

$$A(\chi,\overline{\omega}) = \int_{-\infty}^{\infty} f\left(x - \frac{\chi}{2}\right) f*\left(x + \frac{\chi}{2}\right) \exp(-j\overline{\omega}x)\,dx$$

$$W(x,\omega) = \int_{-\infty}^{\infty} f\left(x - \frac{\chi}{2}\right) f*\left(x + \frac{\chi}{2}\right) \exp(-j\omega\chi)\,d\chi$$

(6.53)

In this equation $\chi$ and $\overline{\omega}$ are the 'size' feature in the spatial and frequency domains respectively, whereas for the Wigner function $x$ and $\omega$ are 'position' features.

Usage of these interesting functions will make it possible to decide which is the best one.

Some pictures of the use of this technique were given in chapter 2. The one shown in figure 6.128 is a typical Wigner function for a surface with a changing frequency, a 'chirp', which could result from an axial vibration of the tool in its holder.

### 6.8.4  Non-linear dynamics

So far the characterization of surfaces for the control of manufacture has tried to use random process analysis. The power spectrum and correlation functions have been obtained from the surface and these have been piecewise identified with various process and machine tool problems. It seems that it may be possible to go even further and to use the surface geometry to quantify and not just identify the actual parameters of the dynamic system making up the machine tool. One example of this has been investigated by Scott [156] and has been reported in chapter 2 in the section on characterization. The idea is to get the coefficient of the second-order equation (i.e. stiffness and damping) directly from the surface profiles. The idea is at present tentative but it is briefly mentioned here.

The first step is to consider basing a machining operation on a dynamic system theory. The basic approach is to obtain values of 'attractors' that appear as solutions to various models of turning. The 'attractor' is the long-term point in state space to which the system locus converges. (Here state space is taken to be a plot of $\dot{z}$ against $z$.)

Figure 6.129 shows the basic model adopted for turning. It is an approximation to the model for milling cutting used by Tlusty and Ismail [157]. Thus

$$m\ddot{z} + T\dot{z} + kz = \sin\alpha.\text{force.}$$

(6.54)

The cutting force excites a vibration $z$ in the oscillator mechanism according to the rule

$$\text{force} = bc_s c_t$$

(6.55)

where $b$ is the chip width, $c_t$ is the chip thickness and $c_s$ is the cutting stiffness.

So assuming that one cut is periodic in time $L$ the instantaneous chip thickness can be modelled as

$$c_t = \begin{cases} z_{\min}(t) - z(t) & \text{if positive} \\ 0 & \text{otherwise.} \end{cases}$$

(6.56)

**Figure 6.128** The application of the Wigner distribution function to machine tool monitoring.

As Scott [156] points out, the fact that $c_t = 0$ means that a chip is not being produced, that is the tool has left the surface. This lack of trackability is the non-linear element in the dynamic system theory.

The rearranged equation of motion (6.55) becomes

$$\ddot{z} + C\dot{z} + Az = Bc_t \tag{6.57}$$

where $A = k/m = \omega_n^2$, $B = \sin(\alpha)bcs/m$ and $C = c/m = 2\zeta\omega_n$, the usual terms in a second-order system.

As long as the tool tip remains in contact with the surface, the above is a linear system. Two types of 'attractors' could be solutions. One is the periodic and the other the quasiperiodic or toroidal. The former is

**Figure 6.129** Turning as a dynamic system: (*a*) turning process dynamic model; (*b*) quasiperiodic attractor envelope of motion.

typical of a system that has an oscillator (i.e. a repeating element). The quasiperiodic attractor is typical of a dynamic system that has two weakly interacting oscillators whose frequencies have no common period. The former is a solution of the linear equation (6.57). The latter seems to be a solution to the linear part of Tlusty and Ismail's model. Hence it is plausible to assume that both types of solution are possible dynamic motions of the tip of the cutting tool.

In the second example figure 6.129(*b*) shows a periodic attractor. In the state-space diagram the path of the tool lies on two intersecting planes, one plane corresponds to when the tool is in contact with the surface, the other when it has left the surface. In this state-space diagram the path of the tool is in free space can be taken into account. This ties in exactly with the criteria of trackability used for the stylus instrument in chapter 4. Thus, the whole concept of dynamic systems expressed in this way pervades many aspects of metrology and manufacture.

In the second example the same values of *A*, *B* and C are used as in the first example (figure 6.130). Because the tool has a different starting point the convergence is to a different 'attractor' which diverges and is symptomatic of a chaotic attractor. This is equivalent to the situation in which the tool or stylus is bouncing along the surface haphazardly.

The early examples occur when the natural frequency of the system is less than the rotational speed of the spindle. The next attractor is found in the full model when the natural frequency of the system is higher than the speed of the spindle—the usual case. Under these circumstances, while the tool tip is in contact a quasiperiodic mode is followed. When the tip is not touching it undergoes conventional harmonic motion. In other words, both modes of operation of the cutting tool can be represented.

From the surface point of view the cutting tool path will be reproduced in the surface when in contact with it. But when the tool is off the surface no cutting takes place and the surface path consists of reproduction of the tool's motion in a contacting region taken by earlier rotations of the spindle.

For dynamic systems it is best to consider a profile taken axially along the surface. The spiral path of the cutting tool will intersect the profile once per revolution. The profile therefore contains a Poincaré section. (This is the section by which the dynamic system can be estimated without taking all of the data of the surface into account. The pattern of intersections on the Poincaré section gives a clue to the complete behaviour.) The actual turning mark which has been used so effectively in the spectrum analysis in section 6.8.2.1

(a)

E3

E2

E1

(b)

E3

E2

E1

An attractor found at A = 1, B = 10, C = 0.005

Z

15
10
5
0

200    210    220    230

Time

Z

16
14
12
10
8
6

500              550         Time    600

Motion of the tool associated with attractor

**Figure 6.130** (*a*), (*b*) Typical attractor at $A = 1$, $B = 10$, $C = 0.005$ (upper diagrams) and motion of tool associated with attractor (lower diagrams).

is no use here. It only contains the information about the tool shape. What is needed is the way in which the tool tip position develops as a function of time. How are the valley positions changing revolution to revolution? This contains the dynamic path information (figure 6.131).

Ideally, from these so-called Poincaré points in the Poincaré section the dynamics of the system can be estimated and from this the system parameters.

So, in principle, this is a technique for directly finding out what the tool path and hence the dynamics of the system are doing as a function of time. Furthermore they are found from the surface profile.

Practically, the situation is not so simple because there are not enough valley points in a typical profile to establish the system parameters if they are anything but the very simple cases described above (figure 6.131). Another problem is that if genuinely chaotic information in the form of noise is added to the system it can upset the parameter estimate. One way out of this is to consider the waviness profile (i.e. that in which the short wavelengths have been filtered out) as more suitable for the Poincaré section analysis than the roughness itself. This is compatible with earlier ideas on waviness except that the waviness here would have to be slightly redefined and would represent the smoothed normal behaviour of the machine tool rather than its unusual behaviour.

It has been pointed out [156] that despite this present lack of positive identification of the use of the technique, it could still be used as a monitor of manufacture because, if chaotic motion or high stochastic noise appears in the state space and Poincaré sections, this would indicate that something is wrong or about to go wrong in the system.

This seems to be a field which would benefit from using areal information from the surface rather than simple profiles. Then many more points would be available giving better estimates of attractors and thus model parameters.

For more information see references [6, 158–160].

**Figure 6.131** (*a*) Poincaré points on profile; (*b*) profile and waviness in turned surface.

## 6.9   Summary

The following trends have been observed in the generation of surfaces:

1. The availability of new surface parameters such as those derived from random process analysis gives the opportunity to examine the process and the machine tool in far greater detail than has hitherto been possible. In grinding, in particular, more use should be made of this. Some new possibilities have also been introduced. The concept of the surface as a fingerprint of the process has been established. New theory offers tremendous benefits.
2. There is a growing awareness of the benefits of linking manufacture directly to function. That the two are often indistinguishable is easy to demonstrate as shown in figure 6.132. Machining is an extreme example of wear.



**Figure 6.132** Linking of manufacture and function.

3. It is being recognized that to achieve surfaces on the atomic scale it is often advisable to match the 'unit' of the process to that of the detail required. This has led to an increase in atomic, ionic and similar new processes.

4. New materials such as ceramics have initiated the development of a new generation of machine tools and also brought into question some of the basic machining mechanisms.

5. There may be a possibility of making instruments mimic the machining process in order best to control the process by means of the surface data.

## References

[1] Donaldson R P 1974 Nomograms for theoretical finish. Calculations with round nosed tools *Paper* MR 937, SME

[2] Shaw M C 1984 *Metal Cutting Principles* (Oxford: Clarendon)

[3] Opitz H and Moll H 1941 *Herstellung Hochwertizwer brehflachen Ber Betriebswiss* (Berlin: VDI)

[4] Nakayama K and Shaw M C 1967/68 *Proc. IMechE* **182** 179

[5] Nakayama K, Shaw M C and Brewer R C 1966 *Ann. CIRP* **14** 211

[6] Sokolski P A 1955 *Prazision inter Metallbearbeitung* (Berlin: VEB Verlag Technik)

[7] Bramertz P H 1961 *Ind. ANZ* **83** 525

[8] Semmler D 1962 *Dissertation* TH Braunschweig

[9] Byrne G, Dornfield D, Inasaki I, Ketteler G, Koenig W and Teti R 1995 Tool condition monitoring—the status of research and industrial application *Ann. CIRP* **44** 541

[10] Crostack H A, Cronjager L, Hillman K, Muller P and Strunck T 1988 Rauhtiefe Beruhrungsios gemessen Erfassung er Oberflachenzustandes spanend bearbeiteter Bauteile mit Ultrashall TZ fur Melallbearbeitung **82** 10

[11] Whitehouse D J 1986 Some modern methods of evaluating surfaces *Proc. SM.E. Chicago Rackels J. H. J. Phys. E. Sci. Inst.* **19** 76

[12] Davies M A, Chou Y and Evans C J 1996 On Chip Morphology. Tool wear and cutting mechanism on finish hard turning. *Ann. CIRP* **45** 77

[13] Hingle H 1992 Control of precision diamond turned surfaces *Proc. VIII Inr. Oberflachen. Kolloq, Chemnitz*

[14] Wilks J 1980 Performance of diamonds as cutting tools for precision machines *Precis. Eng.* **2** 57–71

[15] Brammertz P 1962 Die entsehung der oberflachen raugheit beim Feinddrehen. *Industrieanzeiger* **25** no. 2 p32

[16] Weule H, Huntrup V and Tritschler H 2001 Micro cutting of steel to meet new requirements in miniaturization *Ann. CIRP* **50** 61

[17] Radford J D and Richardson D B 1974 *Production Engineering Technology* (London: Macmillan)

[18] Dickenson G R 1967/68 *Proc. IMechE* **182** 135

[19] Tlusty J, Zaton W and Ismail F 1983 *Ann. CIRP* **32** 309

[20] Martelloti M E 1941 *Trans. AME* **63** 677

[21] Klocke F 1997 Dry cutting *Ann. CIRP* **46** 519

[22] Weinert K 1997 Thamke D 1996 Report at VDI Conference Dusseldorf. *VDI-Verlag Dusseldorf* 111–124

[23] Kustas F M, Fehrehnbacher L L and Komanduri R 1997 Nanocoatings on cutting tools for dry machining *Ann. CIRP* **46**

[24] Abebe M and Appl F C 1988 *Wear* **126** 251–66

[25] Salje E 1988 Relations between abrasive processes *Ann. CIRP* **32** 641

[26] Jacobson S, Walter P and Hogmark S 1987 *Wear* **115** 81–93

[27] Torrence A A 1988 *Wear* **123** 87–96

[28] Childs T H 1970 *Sci. Int. J. Mech. Sci*. **12** 393–403

[29] Lindsay R P 1984 *Ann. CIRP* **33** 143–97

[30] Hoshimoto F, Kanai A and Miyashita M 1983 *Ann. CIRP* **32** 287

[31] Bhateja C R 1984 *Ann. CIRP* **33** 199

[32] Furijawa X, Miyashita M and Sliozaki S 1971 *Int. J. Mach. Tool Des. Res*. **11** 145–75

[33] Rowe B W, Miyashita M and Koenig W 1989 Centreless grinding research and its application in advanced manufacturing technology *Ann. CIRP* **38** 617

[34] Rowe B W 1973 An experimental investigation of grinding machine compliances in productivity *MTDR* Machine Tool Design and Research Conf. 479–88

[35] Whitehouse D J 1988 A revised philosophy of surface measuring systems *Proc. IMechE* **202** 169

[36] Salje E, Teiwes H and Heidenfelder H 1983 *Ann. CIRP* **32** 241

[37] Malkin S 1985 Current trends in CBN grinding technology *Ann. CIRP* **34** 557

[38] Vickerstaff T J 1976 *Int. J. Mach. Tool Res. Des*. **16** 145

[39] Saini D P, Wagner J G and Brown R H 1982 *Ann. CIRP* **31** 215

[40] Hahn R S and Lindsay R P 1966 *Ann. CIRP* **14** 47

[41] Hahn R S 1964 *Trans. ASME* **86** 287

[42] Konig W and Lrotz W 1975 *Ann. CIRP* **24** 231

[43] Pahlitzch G and Cuntze E O 1980 Private communication

[44] Hashimoto F *et al* 1965 *Proc. 6th Int. MTDR Conf*. p507; 1984 *Ann. CIRP* **33** 259

[45] Miyashita M, Hashmoto F and Kanai A 1982 *Ann. CIRP* **31** 221

[46] Trmal G and Kaliszer H 1975 Optimisation of grinding process and criteria for wheel life *Proc. 15th MTDR Conf*. 311

[47] Loladze T N *et al* 1982 *Ann. CIRP* **31** 205

[48] Yossifon S and Rubenstein C 1982 The surface roughness produced when steel is ground by alumina wheels *Ann. CIRP* **31** 225–8

[49] Brinksmeier E and Schneider C 1997 Grinding at very low speeds *Ann. CIRP* **46** 223

[50] Messer J 1983 *Dr Ing. Thesis* Aachen

[51] Bifano T, Blake P, Dow T and Scattergood R O 1987 Precision machining of ceramic materials *Proc. Ceramic Soc. ASME* (Pittsburgh, PA: AES) 99

[52] Miyashita M and Yoshioka J 1982 *Bull. Jpn Soc. Precis. Eng.* **16** 43

[53] Papoulis A 1965 *Probability, Random Variables and Stochastic Processes* (New York: McGraw-Hill)

[54] Whitehouse D J 1971 *PhD Thesis* Leicester University

[55] Terry A J and Brown C A comparison of topographic characterization parameters in grinding. *Ann. CIRP* **46** 497

[56] French J W 1917 *Trans. Opt. Soc.* **18** 8–48

[57] Preston F W 1983 *J. Soc. Gloss. Technol*. **17** 5

[58] Bielby 1921 *Aggregation and Flow of Solids* (London: Macmillan)

[59] Home D F 1972 *Optical Production Technology* (Bristol: Hilger)

[60] Grebenshchikov (1931, 1935); reference not available (referred to in [30])

[61] Touge M and Matsuo T 1996 Removal rate and surface roughness in high precision lapping of Mn-Zn ferrite. *Ann. CIRP* **45** 307

[62] Snoeys R, Staelens F and Dekeyser W 1986 *Ann. CIRP* **35** 467–80

[63] Taniguchi N 1983 *Ann. CIRP* **32** 573

[64] Electro chemical machining *PERA Rep*. 145

[65] Radhakrishnan V, Krishnaiah Chetty O V and Achyutha B T 1981 *Wear* **68**

[66] Masuzawa T and Tanaka K 1983 *Ann. CIRP* **32** 119

[67] Tani Y and Kawata K 1984 *Ann. CIRP* **33** 217

[68] Geiger M, Engel U and Pfestof M 1997 New developments for the qualification of technical surfaces in forming processes *Ann. CIRP* **46** 171

[69] Stout K J *et al The development of methods for the characterization of roughness in 3 dimensions.* Report EUR 15178 EN EC Brussels ISBN 07044 13132

[70] Lenard J G 2000 Tribology in rolling *Ann. CIRP* **49** 567–589

[71] Azushima A *et al* 1996 Direct observation of micro contact behaviour at the interface between tool and workpiece in lubricated upsetting *Ann. CIRP* **45** 206

[72] Sorensen C G, Bech J I, Andreasen J L, Bay N, Engel U and Neudecher T 1999 A basic study of the influence of surface topography on mechanisms of liquid lubrications in metal forming *Ann. CIRP* **48** 203

[73] Azushima A, Miyamoto J and Kudo H 1998 Effect of surface topography of workpiece on pressure dependence of coefficient of friction in sheet metal forming *Ann. CIRP* **47** 479

[74] Nee A Y C and Venkatesh V C 1983 *Ann. CIRP* **32** 201

[75] Miyamoto I and Taniguichi N 1983 *Bull. Jpn Soc. Precis. Eng.* **17** 195

[76] Bryan J B 1972 Private communication

[77] Mori Y and Sugiyama K 1979 OKUDAT *Bull. Jpn Soc. Precis. Eng*.

[79] Preston F 1926 The nature of the polishing operation *Trans. Opt. Soc. London*

[80] Venkatesh V C 1995 Observations on polishing and ultra precision machining of semiconductor substrate materials *Ann. CIRP* **44** 611

[82] Masuzawa T, Tsukamoto J and Fujino M 1995 *Drilling of deep holes by EDM Ann. CIRP* **38** 195

[83] Allu D M and Lechaheb A 1996 Micro EDM of ink jet nozzles *Journal of Materials Processing Technology* **58** 53

[84] Toenshoff H K, Hesse D, Kappel H and Mommsen J 1995 *Eximer laser systems manufacturing systems* **24** 395

[85] Miyamoto I, Taniguchi N and Kawata K 1984 *Proc. 8th Int. Conf. on Production Engineering (Tokyo)* 502

[86] Miyamoto I and Taniguchi N 1982 *Precis. Eng.* **4** 191

[87] Taniguchi N and Miyamoto I 1981 *Ann. CIRP* **30** 499

[88] Yoshioka J, Hasimoto F, Miyashita M, Kanai A, Abo T and Diato M 1985 *Proc. ASME (Miami Beach, November 1985)*

[89] Konig W *et al* 1985 *Ann. CIRP* **34** 537

[90] Furukawa and Kakuta A 1996 Molecular beam epitaxy as an ultraprecision machining process *Ann. CIRP* **45** 197

[91] Reason R E 1951 *Talyrond design* Taylor Hobson

[92]   Scheider N 1972 *Structured rollers* Private communication
[93]   Toenshoff H K, Von Alvensieben, Temme T and Willmann G 1999 Review of laser micro structuring *Proc. 1st Conference of Euspen Bremau Germany* **2** 16
[96]   Peklenik J and Zun I 1995 The energy quanta and the entropy – new parameters for identification of machining processes *Ann. CIRP* **44** 63
[97]   Malkin S and Anderson R B 1974 Thermal aspects of grinding part 1 Energy partition *ASME Journal of Engineering for Industry* **96** 1117–1183
[98]   Zhu B, Guo D, Sunderland J E and Malkin S 1995 Energy portion to the workpiece for grinding ceramics *Ann. CIRP* **44** 267
[99]   Bifano T G, Dow T G and Scattergood R O, 1991 Ductile regime grinding a new technology for machining brittle materials *Trans ASME Journal Engg. For Industry* **113** 184
[100]  Inamura T, Shimada S and Takezewa N 1997 Brittle/ductile transition computer simulations *Ann. CIRP* **46** 31
[101]  Shimda S, Ikawa N, Inamura T, Takezewa N, Ohmori H and Sata T 1995 Brittle-ductile transition phenomena in microindentation and micromachining *Ann. CIRP* **44** 523
[102]  Komanduri R, Chandrasekaran N and Raff L M 1999 Orientation effects in nanometric cutting of single crystals; a M. D. simulation approach *Ann. CIRP* **48** 67
[103]  Carlsson T and Stjernstoft T 2001 A model for calculation of the geometric shape of the cutting tool workpiece interface *Ann. CIRP* **50** 41
[104]  Duc E, Lartigue C, Tourner C and Bourdet P 1999 A new concept for the design and manufacture of free form surface: the machining surface *Ann. CIRP* **48** 103
[105]  Leu M C, Maiteh B Y and Blackmore D Fu L 2001 Creation of free form solid models in virtual reality *Ann. CIRP* **50** 73
[106]  Field H and Kahles J 1971 *Ann. CIRP* **20** 153
[107]  Von Turkovich B F 1981 *Ann. CIRP* **30** 533
[108]  Lucca D A, Brinksmeier E and Goch G 1993 Progress in assessing surface and subsurface integrity *Ann. CIRP* **47** 669
[109]  Roster W P, Gatto L R and Canunett J D 1981 Influence of shot peening on the surface integrity of some aerospace materials *Proc. Int. Conf. Shot Peening (Paris)* (Oxford: Pergamon)
[110]  Field M 1982 Metallurgical alterations and surface integrity produced by material removal techniques *Metal Res. Assoc. Intern. Rep.*
[111]  Brinksmeier E, Cammel J T, Konig W, Leskovar P, Peters J and Tonshaff H K 1982 Residual stress measurement and causes in machining processes *Ann. CIRP* **31** 491
[112]  Week M 1979 *VDI Ber*. no 354, 125
[113]  Schmohl H P 1973 *Dr Ing. Diss*. TU Hannover
[114]  Jutzler W I 1982 *Dr Ing. Diss*. RWTH Aachen
[115]  Syren B 1975 *Dr Ing. Diss*. University of Karlsruhe
[116]  Snoeys R, Maris M and Peters J 1978 *Ann. CIRP* **27** 571
[117]  Jaeger J C 1942 *J Proc Soc New South Wales* 76 pt 3
[118]  Metcut Research Associates 1978 *Publ no* MDC 78-103
[119]  Field M Private communication
[120]  Brinksmeier E 1982 *Dr Ing Diss* TU Hannover
[121]  Schreiber E 1973 Die Engenspannungsausbildung bein schliefen geharten Stahls *HTM* **28** 3 1979 186–9
[122]  Fanz H E 1979 *HTM* **34** 1 24–32
[123]  Kaczmarek I 1966 *Ann. CIRP* **1** 139
[124]  Matsumoto Y, Hashimoto F and Lahoti G Surface integrity generated by precision hard turning. *Ann. CIRP* **48** 59
[125]  Kiethe H 1973 *Dr Ing Diss* TH Karlsruhe
[126]  Machinabriity Data Centre 1980 *Machining Data Handbook* 3rd edn, vol 2 (Cincinnati, OH Metcut Research Associates Inc)
[127]  Gilley F H 1901 *Am. J. Sci*. **11** 269
[128]  Sach G 1927 *Z. Metallkd* **19** 352
[129]  Treuting R G and Read W T 1951 A mechanical determination of biaxial stress in sheet materials *J. Appl. Phys*. **22** 130
[130]  Aksenov G J 1979 Measurement of elastic stress in a fine grained material Z *Angew Phys USSR* **6** 3–16
[131]  Tonshoff H K 1974 *Ann. CIRP* **23** 187–8
[132]  1975 *Proc Workshop on Non-destructive Evaluation of Residual Stress NTIAC 72–2 (San Antonio Texas)*
[133]  1979 Eigen Spannungsmesungen an Kreissageblattern nut elektro magnetischen Verfahren *Forschungsber Laudes NRW* no 2817
[134]  Tonshoff H K and Brinksmeier E 1981 *Ann. CIRP* **30** 509–13
[135]  Macherauch E and Muller P 1961 *Z. Angew. Phys*. **13** 305
[136]  Bollenrth F, Hauk V and Muller E H 1967 *Z. Metallkd* **58** 76–82
[137]  Hauk V and Kockelmann H 1980 Eigenspannungen Tagungsband DGM 241—60
[137]  Kiezle O 1955 Biegen und Ruckferden nut Forsch-Ges. *Blechverarb* 17
[138]  Field M and Koster W 1978 Optimising grinding parameter to combine high productivity with high surface integrity *Ann. CIRP* **27** 523

[139] Alexander J M and Brewer R C 1963 *Manufacturing Properties of Materials* (London: van Nostrand)
[140] Timoshenko S 1934 *Theory of Elasticity* (New York: McGraw-Hill)
[141] Whitehouse D J 1978 *Proc. IMechE* **192** 179–88
[142] Tonshoff H K, Peters J, Inasaki I and Paul T 1992 *Ann. CIRP* **41** 677–88
[143] Sato K 1955 *Tech. Rep. Tohaka Univ.* **20** 59–70
[144] Yang C T and Shaw M C 1955 *Trans. ASME* **17** 645–60
[145] Reichenbach G S and Mayer 1956 *ASME* **18** 847–50
[146] Onoka T 1961 *Bull. Jpn Soc. Grinding Eng.* p27–9
[147] Salje E 1953 *Werkstattstech Betr*. **86** 177–82
[148] Brown R 1969 *Rep*. no 4 *American Grinding Association*
[149] Snoeys R and Peters J 1974 *Ann. CIRP* **23** 227–37
[150] Weinent K 1976 *Dr Ing. Diss*. TU Braunschweig
[151] Snoeys R 1975 Private communication
[152] Netterscheid T 1984 *Dr Ing. Diss*. RWTH Aachen
[153] Knop M 1989 *Dr Ing. Diss*. TH Aachen
[154] Whitehouse D J and Zheng K G 1992 The use of dual-space frequency functions in machine tool monitoring *J. Inst. Phys. Meas. Sci. Technol.* **3** 796–808
[155] Zheng K G and Whitehouse D J 1992 The application of the Wigner distribution function to machine tool diagnostics *Proc. IMechE* **206** 249–64
[156] Scott P J 1989 Non linear dynamic systems in surface metrology *Surf. Topogr*. **2** 345–66
[157] Tlusty J and Ismail F 1981 Basic non linearity in machining chatter *Ann. CIRP* **30** 299–304
[158] Fraser A M and Swinney H L 1986 Using mutual information to find independent coordinates for strange attractors *Phys. Rev*. A **33** 1134–40
[159] Broomhead D S and King G P 1986 On the qualitative analysis of experimental dynamic systems *Physica* **20D** 217–49
[160] Broomhead D S *et al* 1987 Topological dimension and local coordinates from times series data *J. Phys. A: Math. Gen.* **20** 563–9

# Chapter 7
# Surface geometry and its importance in function

## 7.1.   Introduction

This area of the subject is the most important because it is here that the significance of the surface is noticed, not from a theoretical point of view, but in practical terms. The term function is taken to mean the use of the surface. The word performance could also be used.

It has always been known that surface effects are important in friction, wear, reflection, etc; what has not been so obvious is the ability to quantify such characteristics. In the case of length standards, the early ways to measure required some sort of resort to the human body. The same is true in the case of surface metrology, but here it is not the length of the arm or the hand but the response of the fingernail or the eye to the surface which was used. In both cases this was better than nothing, but only just. The advent of modern instrumentation described in chapter 4 has enabled a better understanding of the effect of the surface to be attempted. Unfortunately, whereas in the control of manufacture it is reasonably easy to test out the theory, in the case of function it is not so. This is mainly because of the extreme variety of possible working environments in which the component may be expected to perform. This variety is so vast in terms of speeds and loads alone that the task is virtually impossible.

It is well to be reminded of the objective. It is required to test a surface in some way which will help in the prediction of whether that surface will behave functionally as it was intended by the designer. Ideally, the test should in some way mimic the actual function itself; that is, if the part has to have a certain coefficient of friction for a given load and speed then the test should be a friction test at that load and speed. Obviously this requires a different test instrument for every possible application, which is an economically unsound proposition. For this reason the 'direct' functional test can be thought of as being impractical. This is why the 'indirect' test has become so popular. It is relatively easy to design an instrument for the general purpose of measuring the geometry of the surface. But once the direct method has been dropped it becomes necessary to try to relate geometric parameters to functional ones; this means that there is another step required in the indirect method over the direct test. So although the indirect test method is possible and practical it has a weaker correlation with any particular function. This is why there has been such an enormous effort in theoretically trying to tie down geometry to function. This is in lieu of doing the experiments directly.

There has been progress in this field because of the introduction of random process analysis and of digital techniques into the instrumentation. As well as these a better understanding of the function has developed in the last forty years so that there has been an almost parallel development of instrumentation and functional knowledge. However, the danger is that academics who are good at the mathematics fail to understand either the instrumentation or the function adequately and as a result tend to cloud the investigation rather than clarify it. Perhaps this is an inevitable penalty when dealing with a multi-disciplinary subject!

### 7.1.1  The function map

Probably the biggest problem has been a lack of focus when discussing function. This has been due to two reasons. The first is the fact that the number of applications where surfaces are important is so large and diverse that a systematic approach has been impossible. The second is that, in many cases where a relationship between the surface and function has been found the link has been suppressed due to confidentiality. For these reasons a generic approach is given here to provide a guide. A major aspect of this approach is the classification of the function into a simple format. Obviously no classification can cover all functions but the approach here has been to classify tribological applications and in particular, contact, friction, wear, and lubrication and failure mechanisms. These make up the bulk of engineering functions.

The characterization of function, by using the function map has been achieved by breaking down the behaviour into two independent variables.

(1)  The normal separation of two surfaces
(2)  The lateral relative velocity of the two surfaces. Some of the tribological functions are indicated on figure 7.1.



**Figure 7.1**  Function map.

The function map treats the gap between the surfaces and the relative velocity as *variables*. This enables some *measure* of performance to be made. It is the first attempt to realize that characterization of the surfaces can only be effective if it follows from the functional classification. It is however, possible to treat them as *attributes*. An example is 1996 ISO 1208 Annexe B page 828 which attempts to link motif characterization with function. The branching is shown in figure 7.2.



**Figure 7.2**  Surface classification—by attributes.

The map in figure 7.1 is here called a 'function map'[1]. What is omitted from this diagram are the scales. In general terms the ordinate axis is in micrometres. Elastic contact for example could be 0.1 $\mu$m as shown in figure 7.7. Also the abscissa could have a maximum realistic value of 5 m/sec. Not all of these functions are subject to the same forces. A 'force map' is shown in figure 7.3. The designer should be aware of these different force regimes when specifying a function for surfaces.



**Figure 7.3** Force regime map.

It is interesting to note that the force map fits comfortably into a second order dynamic system—indicated by the curve. Notice that this force map is concerned with systems of surfaces, comprising usually two, and not with individual surfaces. The gap can be liquid-or air-filled and have laminar or turbulent flow (i.e. moving from inertial to viscous forces).

The objective is to use figures 7.1 and 7.3 to fit surface characteristics. Before this it should be remembered that the geometric properties of the surface can aid or inhibit the functional properties.

A simple breakdown of this aspect is shown in figure 7.4.

The general rule from this figure and figure 7.1 is that *average* geometrical properties of the surfaces (and hence the gap) improve performance whereas extreme values tend to ruin performance. An average statistic could well be the $R_a$ and an extreme could be $R_{max}$ or $R_y$.



**Figure 7.4** Dual roles of surface metrology and relationship to the statistics.

The 'defect' category does not mean an extreme of the usual distribution of values; it means the presence of unusual features on the surface not at all concerned with surface statistics. This category is called singular. Table 7.1 lists some examples.

Another point to note is that often different aspects of a surface affect different functions—the surface roughness need not have just one function. Plateau honing for cylinder liners in diesel engines have to be robust to cater for the high forces involved in a combustion cycle. Another part of the same surface has to retain oil. Also, the same finish can have two or more effects. For example in contact, to be discussed in the next section, the nature of a typical contact is determined to an extent by the radius of curvature, but chemical and physical properties also play an equal part. Geometry is necessarily the main constituent in a single contact. However, the distribution of contacts in space is determined almost completely by the height and peak distribution of the roughness and waviness.

**Table 7.1** Types of surface parameters.

| Averages | | Extremes | Singular |
|---|---|---|---|
| Heights | $R_a$ | $R_t$ | Point defects |
| | $R_q$ | $R_{max}$ | Scars |
| | | $R_{tm}$ | |
| MR(%) curve | | $R_z$ | |
| Spacings | HSC | | |
| | $S$ | | |
| | $S_m$ | | |
| | $\lambda_a$ | | |
| Hybrid | $\Delta_q$ | | |

Same statistical distribution  |  Different statistics



**Figure 7.5** Application of metrology chain.

To put the surface metrology in perspective relative to the function map of figure 7.1, consider Figure 7.5. This breaks the map more or less into 'layers of metrology' starting at 'dimension' and moving down the figure.

### 7.1.2 Nature of interaction

One point to remember when dealing with surfaces is that they are fundamentally different from the waveforms encountered in, say, communication signals. This is true even if the waveforms look the same. The reason for this concerns the function of the surface. In almost all cases the function is concerned with top-down or bottom-up considerations. Also, as the format of the function map suggests, a lateral movement is also common.



**Figure 7.6** Difference between surface and conventional signals: (a) spatial (parallel); (b) temporal (serial).

In figure 7.6 (a) the first diagram shows a contact situation and the second the reflection of light. In both cases the interaction starts from above the waveform representing the surface. Also there is no particular significance to be attached to interactions occurring on the left or right or in the middle of the waveform. The action (e.g. the contact) can be anywhere or everywhere. Also the action affects the peaks first rather than the valleys. These two cases are good examples of parallel operation, which in these cases are height sensitive (i.e. non linear). Figure 7.6 (b) shows an electrical waveform entering an analogue circuit. Here the interaction takes place at the right of the signal immediately in front of the circuit which is typical of serial operation. Also the peaks and valleys are equally important: a linear situation.

There is a fundamental problem that arises when considering the characterization of function. Even with the simple concept of the function map the mechanism of behaviour is complicated. This is because there are two independent aspects corresponding to the two axes. As mentioned above the normal approach axis involves parallel operation in which the $x$ (and $y$) co-ordinates of the surfaces are not important and involve top down behaviour. However, along the abscissa, time is involved; this is the dynamic axis. Because of this it is basically a serial activity. The function map can therefore be shown as in figure 7.7.

Even within this breakdown there are complications. For example the serial properties of the abscissa only apply to a profile. In reality, for the areal (3D) case there are many profiles in parallel. So the axis should refer to an averaged profile, where all the $z$ values for a given $y$ are averaged and moved in the $x$ direction—the dynamic direction. This is exactly equivalent to saying that the properties of the ensemble average are serial, which is the correct statistical term for the surface behaviour and which does not invoke the ergodic theorem. How these two axes are characterized mathematically will be considered in section 7.10 in the discussion.

**Figure 7.7** Function map—nature of interaction.

There is another point to remember when dealing with surfaces: the surface geometry represents rather more than just geometry. The way in which the surface has been generated produces this geometry but, in addition, it affects the subsurface in a number of physical ways including stress. The geometry is therefore not just geometry; it is a reminder of other changes that have taken place in the subsurface. It follows therefore that a geometric profile should be accompanied by other physical profiles such as stress, microhardness, etc, if the surface properties are to be considered comprehensive.

The way in which the relationship of surfaces to function is to be approached in this chapter is shown in table 7.6. The chapter is likely to be of most value to designers and researchers interested in finding out how the surface geometry influences performance. It is not an easy topic because of the complexity of the surfaces and the wide variety of possible uses to which they might be put.

Because most of the functionally important cases of surfaces are in tribology these will be considered first. For this reason section 7.2 has been split up into two solid body-effects, which are further split into static and dynamic characteristics as in the function map representation. Into the two-body category falls contact, friction, wear and lubrication. Characterizing two-body function in terms of surface geometry is extremely difficult because, for functions such as contact, wear or friction, it is the interaction between two bodies that is important. It needs both bodies to be specified in a way which reflects this interaction. It is not sufficient to specify the detail of either one of the surfaces—both have to be considered.

**Table 7.2** Breakdown of tribological function

Following the two-body study is a small section on three-body interactions. Then there are one-solid-body applications involving the interaction of one surface with, say, electromagnetic waves.

An important point is that what follows is not just true for machined surfaces. Thomas [2] rightly points out that the roughness of surfaces can have important effects in many situations. He describes instances including ships' hulls; airport runways and so on where the values of surface parameters are very large when compared with engineering surfaces. For example surface wavelengths of ten metres or so along a runway are important from the point of view of undercarriage vibration produced when an aircraft lands. Such examples have been difficult to assess because of the awkward sizes involved: not quite surveying but big for conventional surface metrology. Usually the engineering measurement problem is the opposite: the dimensions being too small for easy examination.

## 7.2 Two-body interaction—the static situation (See table 7.2)

Figure 7.1 can be expressed in two ways, the way shown here is in which both axes are expressed in terms of absolute units (i.e. in metres per second for the abscissa, and metres for the ordinate). This is viewing the function from the operational point of view (i.e. the macro situation). Another way is to make the units along the axes relative to the surface geometry. In this case the ordinate would be the ratio of the normal separation of the main planes to some function of the surface finish height. The abscissa would be calibrated in terms of the ratio of the relative lateral speed to the dominant wavelengths on the surfaces. This axis could then be described in frequency terms. In this format, the micro situation, it becomes difficult to relate to actual transfer energies. From the point of view of usage the first format is probably the easiest to picture because it deals with the macro situation in which the direct relevance of the surface geometry is masked. On the other hand, the second format gives a better representation of the relative scale of size of the surface characteristics to that of the function.

### 7.2.1 Contact

One problem immediately emerges in this subject, and this is verification. What exactly happens in a contact zone is most difficult to determine experimentally because of the inaccessibility of the places involved in the contact. By its very nature it is sealed off from the prying eyes of the investigator. Furthermore, it is on a very small scale, often only micrometres in size. This makes for an altogether investigation-unfriendly situation. Many attempts to see inside the zone have been made as have been described here, but one or both of the solid surfaces have had to be varied from their real state or considerably tampered with. In all cases the results become questionable. This is especially true in the dynamic functions.

Contact and other functional properties can best be described in terms of two phenomena: the unit event and the distribution of these unit events in space.

Take for example contact. The 'unit event' here is a single contact. How a typical contact is formed and what are the properties of the surface which determine it are fundamental issues. Equally important is their distribution laterally and normally. It is obviously important to know the number of contacts per unit area, their average size and their height above the mean plane of the common surfaces. It will be shown later on that the average 'unit event' contact properties are influenced to some extent by the geometry of the surfaces, but also their physical and chemical properties have to be considered. Regarding the distribution of events, however, the surface geometry is critical. Taken as a whole it is clear that the surface geometry is important in the two factors which dictate contact. Consequently, the surface geometry is indirectly important in all applications and situations in which contact is involved, such as friction and wear.

Table 7.3 shows the procedure to resolve the contact parameters of the real area of contact and the separation of the surfaces. This table illustrates that peak characteristics are not always the only feature of interest. In thermal conductivity, for example, the average gap between the two surfaces decides how much heat is exchanged. Also, in lubrication it is often the angular flanks of the peaks which initiate the fluid film and

not the peaks themselves. The peaks often only serve to establish the gap. Another misconception is the belief that the peaks on opposing surfaces touch.

Williamson [3] showed vividly that this is not so. Most contact is between the peaks on one surface and a shoulder or flank on the other, so to investigate a contact situation involves both peak distributions and surface height distributions at the same time.

**Table 7.3** Chain of considerations for surfaces under load.

Load
↓
General shapes of surfaces
↓
Nominal area—pressure distribution
↓
Number of contacts
↓
|              | → summit height    |
| Distribution | → summit curvature |
|              | → ordinate heights |
↓
Plastic/plastic index
↓
Real area of contact
↓
Separation of surfaces
↓
Function

The surfaces can be loaded vertically together with no lateral movement; they can be moving transversely relative to each other in solid-solid contact or they can be separated by a fluid film which can be liquid or gas. It is via the contact that the load is supported elastically or plastically, it is via the contact that electrical energy is transferred and it is also via the contact that a large part of the thermal energy is transferred. Unfortunately, the amount of contact is critical to what actually happens. The chance of two perfectly smooth and flat surfaces making a perfect contact is zero, notwithstanding the conceptual problem of the what is perfect smoothness and flatness. In practice what tends to happen is that the actual area of contact is considerably smaller than the apparent area of contact suggested by the macrodimensions of the parts. Contact is restricted in reality to a relatively few points. The number of points is determined by the fine detail of the surface texture as well as longer wavelengths and general shape, which also determine the possibility of interaction.

To determine the real contact area and hence the number of actual contacts for many geometrical shapes is often very difficult, especially if both surfaces are rough, so it is often expedient to consider one of the surfaces to be perfectly smooth and the other rough, the rough surface taking the combined roughness of the two. The validity of this is quite questionable except for smooth surfaces where the amplitude is low and the wavelengths are long. Any surface slopes above 0.1 rad are likely to produce problems in the analysis.

Solution of the two-body contact problem is often iterative. In most attempts the first consideration is that of determining the nominal area of contact. This has to be done whenever the geometrical shapes of the two bodies are non-conforming.

This area is generally arrived at using Hertzian theory. From this some idea of the size of the contact region is found. Knowing this it is usually possible to get an idea of the number of contacts within this area.

Then the balance of plasticity to elasticity in these local contacts is found. This is where the distribution of summit and ordinate heights and the distribution of even more detailed features such as summit curvature

become important. The iteration between these local aspects and the total load being supported depends a lot on the specific nature of the assumptions made in the contact model. There is no recognized answer yet to this problem. The usual flow is shown in table 7.3.

In what follows, this general pattern will be observed. Even so, evaluation of the real area of contact and the compliance from the distribution of pressure within the contact region is difficult and only rarely solvable analytically. It can usually only be done numerically by digitizing the surface and carrying out a computer simulation of contact. The crux of the problem is the establishment of the pressure distribution.

### 7.2.1.1 Point contact

The basic solution for a deflection at a point due to a vertical force elsewhere is shown in figure 7.8. From this the pressure distribution in more complex loading can be derived, for example for spheres and in the case of other macroshapes as shown in figure 7.9.



**Figure 7.8** Compliance due to point load.



**Figure 7.9** Compliance of two curved bodies in contact.

In figure 7.8, $W$ is a point load. This gives a compliance $\delta$ corresponding to the movement in the vertical direction at a distance of $r$ from the loading point. Thus

$$\delta = (1 - v^2)/\pi E \qquad (7.1)$$

Of course this is a simplification, for in practice there is no point load; the load is distributed in the sense that two finite geometrical shapes make contact. Typical configurations are spheres on flats, spheres on

spheres, cylinders on cylinders, etc. There are two basic issues, one macroscopic and the other microscopic, as already indicated. The former determines the possible area over which asperities contact and the latter corresponds to the individual behaviour at each asperity. In this latter case the asperity is often simplified in shape to make the calculations possible. Typical shapes used are spheres or paraboloids.

### 7.2.2 Macroscopic behaviour

The treatment which follows derives from that given by Love [4]. It highlights the basic problem, which is that the behaviour at any point is determined by the pressure distribution over the whole contact area.

It is very necessary to consider all possible shapes of component because in many frictional experiments with wear, friction and lubrication, all types of geometric configurations are used, ranging from four-ball machines to crossed cylinders. In each case the real area of contact has to be evaluated. This is even more general in practical situations. Also the contact of spheres, planes and cylinders describe what happens at the microscopic scale in which asperities of the surface roughness make contact.

When two bodies are pressed together displacement will occur in both. It is assumed that the forces operate parallel to the $z$ axis and that displacements in this direction only are being considered. The effect of this displacement is to provide an area in nominal contact (figure 7.9).

If the displacements at a point are $\omega_1$ and $\omega_2$ then, for points inside the area of contact, since the bodies touch over this area,

$$(z_1 + \omega_1) + (z_2 + \omega_2) = \delta \tag{7.2}$$

choosing the axes such that

$$z_1 + z_2 = Ax^2 + By^2$$
$$Ax^2 + By^2 + (\omega_1 + \omega_2) = \delta. \tag{7.3}$$

It follows from reference [5] that, assuming the surface to be plane, the deformation at a point $x$, $y$ due to this applied force is given by

$$\omega(x, y) = \frac{(1 - v^2)}{\pi E} \frac{p(x', y')}{r} dx' \, dy' \tag{7.4}$$

where $r$ is the distance from the point $x$, $y$ to the point $x'$, $y'$. Also this can be extended by using the theory of superposition. Hence the displacement at a point $x$, $y$ due to the distribution of pressure over an area $A$ is given by

$$\omega(x, y) = \frac{1 - v^2}{\pi E} \iint_A p \frac{(x', y')}{r} dx' \, dy'. \tag{7.5}$$

Putting equation (7.5) into (7.3) yields

$$\left( \frac{1 - v_1^2}{\pi E_1} + \frac{1 - v_2^2}{\pi E_2} \right) \iint_A p \frac{(x', y')}{r} dx' \, dy' = \delta - Ax^2 - By^2 \tag{7.6}$$

where the subscripts 1 and 2 represent the elastic constants for the two bodies. It also follows that

$$\frac{\omega_1}{\omega_2} = \left( \frac{1 - v_1^2}{E_1} \right) \Big/ \left( \frac{1 - v_2^2}{E_2} \right). \tag{7.7}$$

Equation (7.6) is the essential equation because the solution of this should give expressions for the area of contact, which is the parameter of principal interest, the pressure distribution over the area, and the compression [6].

Equation (7.6) can be evaluated using the method of Kellog [7], from which it emerges that the probability density

$$p(x, y) = k\left(1 - \frac{x^2}{a^2} - \frac{y^2}{b^2}\right)^{1/2}. \tag{7.8}$$

Equating the integral of $p(x, y)$ to the total applied force $P$ tending to

$$k = 3P/2\pi ab \tag{7.9}$$

from which

$$p(x, y) = \frac{3P}{2\pi ab}\left(1 - \frac{x^2}{a^2} - \frac{y^2}{b^2}\right)^{1/2} \tag{7.10}$$

hence

$$\tfrac{3}{4}P(V_1 + V_2) \int_0^\infty \left(1 - \frac{x^2}{a^2 + \psi} - \frac{y^2}{b^2 + \psi}\right)\frac{1}{[(a^2 + \psi)(b^2 + \psi)\psi]^{1/2}}\,\mathrm{d}\psi = \delta - Ax^2 - By^2 \tag{7.11}$$

yielding the three equations

$$\delta = \tfrac{3}{4}P(V_1 + V_2) \int_0^\infty \frac{\mathrm{d}\psi}{[(a^2 + \psi)(b^2 + \psi)\psi]^{1/2}} \tag{7.12}$$

$$A = \tfrac{3}{4}P(V_1 + V_2) \int_0^\infty \frac{\mathrm{d}\psi}{(a^2 + \psi)[(a^2 + \psi)(b^2 + \psi)\psi]^{1/2}} \tag{7.13}$$

$$B = \tfrac{3}{4}P(V_1 + V_2) \int_0^\infty \frac{\mathrm{d}\psi}{(b^2 + \psi)[(b^2 + \psi)(a^2 + \psi)\psi]^{1/2}}. \tag{7.14}$$

The variable $\psi$, is a dummy variable and

$$V_1 = \frac{1 - v_1^2}{\pi E_1} \qquad V_2 = \frac{1 - v_2^2}{E_2}.$$

The routine to get the axes of the contact ellipse, $a$ and $b$, is not straightforward. However, the general approach is to use estimated values of $A$ and $B$ to include in equations (7.13) and (7.14), from which $a$ and $b$ are determined. The way in which this works is best seen by an example.

### 7.2.2.1 Two spheres in contact

Let the spheres have diameters $D_1$ and $D_2$ respectively; then

$$z_1 = \frac{x^2}{D_1} + \frac{y^2}{D_1} \qquad z_2 = \frac{x^2}{D_2} + \frac{y^2}{D_2}$$

and the area of contact reduces to that of a circle radius $a$, from which it follows, making $\Psi = p^2$, that from which $a$ can be found thus:

$$A = B = \frac{3\pi}{8a^3}P(V_1 + V_2) \tag{7.15}$$

which gives the appropriate form of the equation which applies to most practical situations:

$$a \simeq \left(\frac{PR}{E}\right)^{1/3}. \tag{7.16}$$

The compliance can be worked out from equation (7.15) moving $a$. Thus

$$\delta = \tfrac{3}{4} P(V_1 + V_2) \int_0^\infty \frac{2\mathrm{d}\rho}{a^2 + \rho^2}$$

$$= \frac{3\pi}{4a} P(V_1 + V_2) \tag{7.17}$$

from which $\delta$ is given for the general case of two spheres of different size and material as:

$$\delta = \left(\frac{3\pi}{2}\right)^{2/3} P^{2/3} (V_1 + V_2)^{2/3} \left(\frac{1}{D_1} + \frac{1}{D_2}\right)^{1/3} \tag{7.18}$$

and for spheres of the same material and size as:

$$\delta = \left(\frac{9}{2}\right)^{1/3} \left(\frac{1 - v^2}{E}\right)^{2/3} P^{2/3} \left(\frac{2}{D_2}\right)^{1/3}. \tag{7.19}$$

This reduces by a factor of $21/3$ for the case of a sphere on a plane of the same material.
Following Puttock and Thwaite [8], if $P$ is in grams force $D$ is in millimetres:

$$\delta = 0.00002 P^{2/3} \left(\frac{2}{D}\right)^{1/3} \text{ mm}$$

and if $P$ is in pounds force and $D$ is in inches:

$$\delta = 0.000016 P^{2/3} \left(\frac{2}{D}\right)^{1/3} \text{ in.}$$

It is much simpler to estimate results rather than to work them out in detail, and formulae derived even for simple spherical contact are quite crucial in all branches of metrology. This is because the general case of contact between two quadratic forms is very difficult to evaluate. Invariably the single spherical model is used unless there is a justification for doing otherwise.



**Figure 7.10** Contact of sphere on flat.

An interesting result occurs in the case of a sphere on a plane (figure 7.10). The actual area of contact is

$$2\pi r_{\mathrm{d}}^2 = \pi r_{\mathrm{u}}^2. \tag{7.20}$$

In words, the actual contact area is half that which would be expected if the sphere (or parabola) were simply to intersect the undeformed plane. The case where a sphere is contacting inside a sphere can be catered for simply by making the substitution $1/D_1 + 1/D_2 = 1/D_1 - 1/D_2$, where $D_2$ is the larger radius.

### 7.2.2.2 Two cylinders in contact

For cases where cylinders are involved, the equations become somewhat complicated. However, for two equal cylinders at right angles, two of the curvatures are equal and two are infinite. This gives the obvious result that the situation is the same as that of a sphere on a plane:

$$a = \left[ \left( \frac{3\pi}{8} \right) P(V_1 + V_2) D^{-1} \right]^{1/3}. \tag{7.21}$$

Also the compliance $\delta$ is the same as for a sphere on a plane.

For unequal cylinders at right angles it is necessary to resort to elliptical integrals of the first kind.

Let $e$ be the ellipticity of the ellipse of contact. Multiplying both the numerator and the denominator of equation (7.13) by $(1/\delta^2)^{5/2}$, after some calculation it can be shown that

$$Aa^3 = 1.5P(V_1 + V_2) \int_0^{\pi/2} \frac{\sin^2 \theta \, d\theta}{(1 - e^2 \sin^2 \theta)^{1/2}} \tag{7.22}$$

$$Ba^3 = 1.5P(V_1 + V_2) \int_0^{\pi/2} \frac{\sin^2 \theta \, d\theta}{(1 - e^2 \sin^2 \theta)^{3/2}} \tag{7.23}$$

$$K = \int_0^{\pi/2} \frac{d\theta}{(1 - e^2 \sin^2 \theta)^{3/2}}$$

is an elliptical integral of the first kind from which

$$\frac{dK}{de} = e \int_0^{\pi/2} \frac{\sin^2 \theta \, d\theta}{(1 - e^2 \sin^2 \theta)^{3/2}}.$$

Also the complete elliptical integral of the second type is

$$E = \int_0^{\pi/2} (1 - e^2 \sin^2 \theta)^{1/2} \, d\theta$$

and

$$\frac{dE}{d\theta} = -e \int_0^{\pi/2} \frac{\sin^2 \theta \, d\theta}{(1 - e^2 \sin^2 \theta)^{1/2}}.$$

Hence

$$Aa^3 = -\frac{2QP}{e}\frac{dE}{de}$$

$$Ba^3 = -\frac{2QP}{e}\frac{dK}{de}$$

$$\delta = \frac{2QP}{a}K \tag{7.24}$$

where

$$Q = \tfrac{3}{4}(V_1 + V_2)$$

from which

$$\delta = 2K(PQ)^{2/3}\left[2D\left(-\frac{1}{e}\frac{dE}{de}\right)\right]^{-1/3}$$

given the relationships

$$\frac{dE}{de} = \frac{1}{e}(E - K)$$

$$\frac{dK}{de} = \frac{1}{e(1-e^2)}[E - (1-e^2)K]$$

$$\frac{A}{B} = \frac{-(1-e^2)(E-K)}{E-(1-e^2)K} \qquad \frac{1}{e}\frac{dE}{de} = \frac{1}{e^2}(E-K). \tag{7.25}$$

Therefore, for any value of $e$, the values of $A/B$, $K$ and $-(1/e)/(dE/de)$ can be found from tables (e.g. [8]). From these values $a$ and $\delta$ can be found for all configurations in this class.

### 7.2.2.3 Crossed cylinders at any angle

For crossed cylinders at any angle $\theta$ (acute)

$$\delta = 2K(PQ)^{2/3}\left(\frac{A}{-(2/e)/(dE/de)}\right)^{1/3}. \tag{7.26}$$

The ratio $A/B$ has to be determined prior to this.

### 7.2.2.4 Sphere on a cylinder

For this

$$z_1 = \frac{x^2}{D_1} + \frac{y^2}{D_1}$$

so

$$A = \frac{1}{D_1} \quad \text{and} \quad \frac{A}{B} = \frac{1/D_1}{1/D_1 + 1/D_2}$$

from which

$$a^3 = \frac{2QP}{A}\left(-\frac{1}{e}\frac{\mathrm{d}E}{\mathrm{d}e}\right)$$

$$\delta = \frac{2QP}{a}K. \tag{7.27}$$

### 7.2.2.5 *Sphere inside a cylinder*

For this

$$a^3 = \frac{2QP}{A}\left(-\frac{1}{e}\frac{\mathrm{d}E}{\mathrm{d}e}\right)$$

$$\delta = \frac{2QP}{a}K \tag{7.28}$$

where

$$A = \frac{1}{D_1} - \frac{1}{D_2} \qquad \text{and} \qquad \frac{A}{B} = \frac{1/D_1 - 1/D_2}{1/D_1}.$$

From these expressions it is obvious that the combinations with other geometries are endless, but nevertheless it is always possible to estimate the area to within a factor of 2 or 3 using the results above. This is adequate considering the other uncertainties in the treatment, which may produce considerably greater deviations.

Nevertheless, the nominal contact area has to be worked out as a preliminary to any estimate of real contact. The next step is to find the true area of contact from within the nominal contact region as a function of surface roughness.

### 7.2.3 *Microscopic behaviour*

#### 7.2.3.1 *General*

Remembering earlier discussions on the characteristics of surfaces, it is well to be reminded that a summit is defined as a peak within an area, that is an areal (or three-dimensional) peak. The term peak is restricted for the sake of argument to a maximum value on a profile graph taken by some sectioning process from a surface by means of an instrument.

There are a number of issues which need to be considered in contact situations. These include the density of peaks or summits and the distribution of their heights and shapes. These are important in deciding the balance between plastic and elastic deformation.

In a true contact situation, obviously two surfaces are involved and it could be argued that the behaviour in between two surfaces in contact should be experimentally determinable. Unfortunately this is rarely so. The contact region is almost totally inaccessible and very small relative to the component size. Many attempts have been made to make direct measurement of the number of contacts, but most have been relatively unsuccessful. The usual method has been to resort to estimating the contacts either from precontact measurement of the mating surfaces with some degree of extrapolation or by computer simulation. Both techniques are better than nothing but not ideal.

Perhaps the best example of contact prediction has been made by Williamson [3], although area-tracking methods have been used extensively (see e.g. [9]). His technique goes part way to carrying out the contact experiment with full instrumentation. The surface data is quite genuine, at least in the sense that the data is taken off a real surface. The actual contact between the surfaces is simulated directly on

the computer. Although this may be somewhat artificial in the sense that the contact is achieved in the computer, it demonstrates some interesting points. The most important is that actual contact does not always occur at the maximum value on either surface; it is invariably on the shoulder or flank of the respective summits. This feature of contact highlighted for the first time the fallacy of dealing altogether with the statistical properties of surfaces independently. In fact it is the point-to-point properties of the 'gap' between the surfaces which may well determine the function of the two surfaces in contact. It is a pity that practical instrumentation has not yet been able to achieve this result. The signal-to-noise ratio, and the mechanical correlation needed, seems even today to be formidable. Earlier on it has been demonstrated just how difficult it is to measure surfaces over an area using digital means. This is exacerbated when the problems of gap are considered. Also, quite often it is very difficult to procure a set of 'typical' surfaces in order to prove a particular tribological point. Invariably the available surfaces have extra features on them, such as waviness or a deep scratch, which can mask the demonstration of the desired effect.

One alternative is to carry out the experiment, at least in its pilot form, on the computer.

Methods of characterizing surfaces using autoregressive moving-average methods have already been discussed, although this characterization basis can be used as a method of generation in terms of a profile and as a two-dimensional method; it is in the 2D or area method where the technique comes into its own. Problems of mechanical correlation between the mean levels of successive profiles and synchronization of digital profiles are not an issue if a single profile is to represent the surface. They are very important for area representation. It is in these circumstances where 2D surface generation comes into its own. Such a method is described in section 3.11.2.

What happens when roughness is introduced to the situation where the nominal area has been found?

### 7.2.3.2  Types of parameter

Problems associated with contact parameters such as the number of asperities per unit area have attracted numerous investigators. Contact parameters have been evaluated for deterministic waveforms such as sinusoids or triangles, for random distributions of hemispherical bosses (representing the tips of asperities) [10] and for random processes in which the properties of the asperities are derived from the statistics of the process [11, 12]. More recently the surface parameters have been derived using fractal geometry [13, 14] which is somewhat surprising because, although fractals can justifiably be used for natural surfaces and growth type processes, they fit uncomfortably with machined surfaces.

Sometimes it is useful to consider the values of contact parameters as a function of height, (e.g. the density of peaks), or the perceived area of contact at each level. For example [15] the density of closed contours at any separation is $D_{cont}$ given below. It has to be remembered however, that such elementary analysis does not include peak correlation. Thus

$$D_{cont} = \frac{M_2}{M_0} \frac{1}{(2\pi)^{\frac{3}{2}}} \exp\left(-y^2/2S_q^2\right) \tag{7.29}$$

where $S_q$ is the root mean square value of the surface roughness. This is a composite roughness if two rough surfaces are involved. Sayles and Thomas [15]. also give a useful upper bound value of the average radius of the contacts $r_{av}$

$$r_{av} < (2\pi)^{\frac{1}{4}} \sqrt{\frac{M_0}{M_2 y}} \tag{7.30}$$

similar expressions can be derived using fractal analysis, e.g. $n_L$. is the number of contact points larger than $a$

$$n_L = A_L^{D/2} . A_L^{D/2} \tag{7.31}$$

given that area *A, D* is the fractal dimension. Majumdar [13] realized that there are theoretically an infinite number of small contacts but he maintains that they make a negligible contribution, although this treatment of the limiting behaviour can be questioned. It should be pointed out that this extraordinary behaviour is not confined to so called 'fractal' surfaces. Any function with an exponential autocorrelation such as a Markov process (often found in manufacturing processes such as grinding) suffers from the same problem.

Contact parameters can be evaluated at different heights using an assumed model for the surface, but these parameters sometimes lack conviction because they are a result of operations on a fixed model for the surface not dependent on the degree of contact: the model does not change. Intuitively the model has to change when load is applied. This will be addressed in the section on elastic/plastic models.

### 7.2.3.3 Peak correlation—possible mechanisms

Sometimes, however, the way in which the model should react to loading for example is anticipated in the mathematics.

Consider a Gaussian surface. According to random theory the number of maxima between the heights $y$ and $y + \delta y$ of a random profile is obtained from the probability density of a peak $f(\hat{y})$ where

$$f(\hat{y}) = \frac{1}{(2\pi)^{\frac{3}{2}}} \cdot \frac{1}{M_0 M_2} \left[ |M|^{\frac{1}{2}} \exp(-\bar{p}^2 k^2 y^2) + m_2^2 y \left[ \frac{\pi}{2M_0 M_2} \right]^{\frac{1}{2}} \exp\left( \frac{y^2}{2M_0} \right) \left[ 1 + erf(\bar{H}ky) \right] \right] \tag{7.32}$$

which is derived from the joint probability density function $p(y, \dot{y}, \ddot{y})$.

The probability of a maximum in the same interval is $p(\tilde{y})$ where

$$p(\hat{y}) = \frac{1}{4\pi\sqrt{M_0 M_2}} \cdot \exp\left( -\frac{y^2}{2M_0} \right) \left[ 1 + erf(\bar{H}ky) \right] \tag{7.33}$$

where

$$\bar{P} \text{ is } \sqrt{\frac{M_4}{M_2}} \text{ and } \bar{H} \text{ is } \sqrt{\frac{M_2}{M_0}}, |M| = M_2(M_0 M_4 - \tfrac{2}{2}) \tag{7.34}$$

Integrating equation (7.32)

$$F(\hat{y}) = \int_{y'}^{\infty} f(\hat{y}) \, dy = \frac{1}{4\pi} \left[ \underbrace{\bar{P}\left(1 - erf\left(\bar{P}ky'\right)\right)}_{A} + \underbrace{\bar{H} \, \exp\left( -\frac{y'^2}{2M_0} \right) \left[ 1 + erf\left(\bar{H}ky'\right) \right]}_{B} \right] \tag{7.35}$$

The interesting point about equation (7.35) is that the two components *A* and *B* which describe peaks at and above $y'$ show different types of peak! The two components above occur naturally from the mathematics.

A typical situation is shown in figure 7.11. In the case shown it is plateau honed.

**Figure 7.11** Process with different types of peak. In the case shown it is plateau honing.

In the figure peaks marked '*P*' are represented by the component *A* in equation 7.35 and peaks marked *H* are represented by *B*. The *A* component is dominated by $\overline{P}$ and *B* by $\overline{H}$. In the terminology here $\overline{P}$ type peaks are called 'dependent' peaks and $\overline{H}$ type peaks are called 'independent'. In fact the *A* and *B* components can be modified to allow the peak definitions to be more flexible and realistic, by making the definition dependent on only one of its flanks and not both.

The reason why equation (7.35) is important is that it makes the distinction between independent and correlated peaks. The same distinction is not possible in equation (7.32) because it only refers to one level. Take for example a type *H* peak shown in figure 7.12. When loaded it will simply crush in plastic deformation or elastically yield depending on its curvature, load and physical parameters; its behaviour does not affect surrounding peaks. On the other hand the '*P*' type peaks are so close that the deformation of one will affect surrounding ones as illustrated in figure 7.12.

Notice that the *P* type peaks are always part of a high spot which is not counted as a peak and yet is fundamental to the behaviour of the *P* peaks under load. The *P* peaks are correlated with each other by virtue of being included in one high spot at $y'$. This correlation is not simple as will be seen. Current thinking does not take the correlation into account adequately with the result that real contacts at a given height will be over estimated.

Loading peak $P_1$ causes an elastic compliance $\Delta y_1$ at $P_1$ but also a compliance $\Delta y_z$ at $P_2$ despite the fact that the loaded surface never actually touches $P_2$. This is a case of correlated peaks, which equation (7.35) indicates are the '*P*' type. Obviously if plastic deformation is the mode there can be no correlation.

The effect of this correlation is to reduce the effective number of peaks making contact at any level with respect to the number of peaks obtained simply by truncating the surface profile at the load level $y_L$.
In all cases of contact under load the most important issue is to determine the mode of deformation; whether peaks distort elastically, plastically or both.

There is, however, still a problem in the elastic mode which is due to the peak correlation mentioned above.



**Figure 7.12** Breakdown of peak types.

To illustrate the effect consider figure 7.12. This shows *physical* correlation which is shown in figures 7.13, 7.14 and 7.15. They show that physical correlation is obtained by compliance propagation, which exists only if there is an underlying high spot. This allows compliance to propagate via valleys to all adjacent peaks. This corresponds to impinging step by step via valleys if above the nominal load level to all the $P$ peaks but not the $\overline{H}$. In the latter case the valleys are below $y_L$—the load level, and also most likely below the meanline.



**Figure 7.13** Compliance propagation.



**Figure 7.14** Closed loop contact behaviour.



**Figure 7.15** Interaction of compliance loops.

There are two cases shown. One is when only one peak makes contact and the other is when both peaks make contact. If $K = \dfrac{\hat{y}_1}{\hat{y}_3}$ is the compliance ratio then, for just the two peaks making contact, the effective stiffness of the peaks is

$$\Delta S \sim \frac{1}{1-2K}. \tag{7.36}$$

For $M$ contacting peaks

$$\Delta S = \frac{1}{1-MK} \quad \text{where } MK < 1, K < 1. \tag{7.37}$$

The closed loop compliance case produces a situation in which the actual asperities appear to get stiffer and the underlying waviness (i.e. the base wave of the high spot count) appears to become more compliant. See figure 7.15 which shows the interdependence loops.

Actually the load is shared serially between the local peaks $a$ and the high spots. So there are two mechanisms: a serial effect at any one point of contact and a parallel mechanism between contacts.

In mechanical terms the Winkler model can be adopted (see [64] and figure 7.60) with the important difference that the elastic rods can be connected so that there is some shear which corresponds to the physical correlation. This is shown as a compliance diagram in figure 7.13.

This surprising result of local peak stiffness may explain why asperities seem to persist even in the presence of big loads such as may be found in a cylinder liner. By taking the compliance ratio K to be a height ratio the conclusions could be criticized, but it turns out that this stiffness result holds whatever model of deformation is used.

The compliance propagation model is a variation on the Winkler model, which describes the surface by a set of 'elastic rods' at the heights of the independent asperities (i.e. the $\overline{H}$ peaks as shown in figure 7.16). The baseline is shown at the position of the deepest valley or at the mean line depending on the model.



**Figure 7.16** Independent peaks.

This simple mechanism has two components, where $K_1\Delta y_1$ is elastic compliance of point $A$. Then $\Delta y_1$ is $\Delta V_1$ and $\Delta y_2$ (i.e. $\Delta y_2 = \Delta y_1$). $K_1$ is $\dfrac{V_1}{y_1}$ and $K_2$. The $P$ peaks are shown in figure 7.16; $P_1$ and $P_2$ are good examples.

Notice that, when the two correlated peaks make contact, each affects the other and closed loop compliance results, which is smaller than $\Delta z_1$ and $\Delta z_2$ of figure 7.17.

The question arises as to the test for independent peaks. There is no absolute way but it is easy to make an estimate of $x_{\text{p corr}}$.

(1) Determine the autocorrelation length of the surface ordinates $x_{\text{corr}}$. This determines the independence length of the surface and is a lower bound.

**Figure 7.17** Correlated peak (Winkler modified).

(2)  The peak independence is found by the formula

$$Nx_{\text{corr}} = x_{\text{p corr}}$$

$$\text{Hence } n + 1 = 2 \int_{\substack{1m \\ k \to \infty}} \left( 1 + 2 \sum_{i=1}^{k} \left( \frac{1}{2} \right)^{i+2} - 2(K+1)/2^{k+1} \right) = 4 \tag{7.38}$$

So n = 3.

Hence the Winkler 'rods' for independence have to be spaced about $3x_{\text{corr}}$ apart. Obviously $P$ peaks can be much closer.

The argument above is not about waviness. It is more concerned with the density of local peaks relative to the density of high spots (i.e. $S$ relative to $S_m$ in surface roughness parlance). If the average number of local peaks relative to high spots is high then there will be a large number of correlated peaks so that there would be fewer contacts at any level.

Waviness determines the number of high spots that would be in the contact zone. This in turn determines the number of local peaks and indirectly the number of contacts.

The credibility of the compliance propagation approach is based on the partitioning of the surface into $\overline{P}$ and $\overline{H}$ components. Looking at a particular profile can give the impression that the $\overline{P}$ type peaks are superimposed $H$ peaks. The $\overline{P}$ peaks are additive. When loaded the $\overline{P}$ type peaks can offload compliance onto the $H$ type, which in general have lower curvature, thereby effectively stiffening up the smaller local $\overline{P}$ peaks.

Another possible mechanism for transmitting compliance from the contact zone is the 'punch' method. In this the compliance is determined by a compliance weighting function $f(x,w)$ centred on the contacting peak $P_1$ and whose magnitude is a function of the distance $x$ from $P_1$. This compliance degrades the height of $P_2$. A typical weighting function could be obtained by considering the contacting peak $P_1$ to be a flat punch indenting the high spot containing $P_2$ and other peaks. Such a function is $f(x,w)$ where

$$f(x,w) = \frac{(1-v^2)}{\pi E} w \left\{ (x+a) \ln\left( \frac{x+a}{a} \right)^2 - (x-a) \ln\left( \frac{x-a}{a} \right)^2 \right\} + C \tag{7.39}$$

which is an extension of the point load in figure 7.2. $a$ is the width of the peak as seen from the base. More closely resembling an asperity is a triangular weighting function.

$$f(x,w) = \frac{(1-v^2)}{\pi E} \frac{w}{a} \left\{ (x+a)^2 \ln\left( \frac{x+a}{a} \right)^2 + (x-a)^2 \rho n\left( \frac{x-a}{a} \right)^2 - 2x^2 \rho n\left( \frac{x}{a} \right)^2 \right\} + C \tag{7.40}$$

**Figure 7.18** Archard model of elastic contact.

The constant $C$ depends on the arbitrary origin of the deflection axis. These functions are a means whereby the lateral spacings of peaks are taken into account. In the case where more than one local peak is contacting, the resultant compliance at any spot has to be evaluated. Using superposition is only allowable if the contacts approximate to cylinders rather than spheres (i.e. anisotropic contacting surfaces).

Both methods, the mattress (modified Winkler) or punch are approximations whereby a realistic correlation between peaks is determined by geometric conditions—the surface geometry, the physical parameter $E$, the elastic modulus and applied forces $W$.

The model above is different from other models dealing with elastic deformation. Instead of dealing just with the roughness and sometimes the waviness, a further element has been brought in—the high spot. This is embedded in the roughness signal and is defined relative to the profile mean line. It represents the basic structure of the roughness.

The sequence for contact is therefore:

Workpieces → form → waviness → high spots → local asperities → contacts.

In the above treatment two possible mechanisms have been suggested to explain correlated peaks. Both are plausible but they do not represent all possibilities.

In what follows plastic deformation of peaks as well as elastic deformation and the balance between the two will be considered.

### 7.2.3.4    Microcontact under load

Archard [16] considered the very simple case of a sphere on a flat. He then added asperities to the sphere (figure 7.18) and went on adding asperities to asperities in much the same way that fractal surfaces are built up. He assumed all asperities were spherical of various scales of size. From this model he attempted to simulate the calculation of the real area of contact $A$ as a function of the load. In figure 7.18($a$) the relationship is Hertzian, $A \propto W^{2/3}$. However, in ($b$) $A \propto W^{4/5}$ and as the surface becomes rougher and rougher $A \rightarrow W$. This is an astonishing result and is not at all obvious. It demonstrates that an apparently proportional relationship between $A$ and $W$ can be possible using a completely elastic regime and not invoking plasticity at all [16].

Archard used the relationship $\delta A \propto W^{2\psi}$ to characterize the increase in area with load. The load per unit contact is $W/N$ and $\sigma A \propto (W/N)^{2/3}$ because the deformation is elastic. Therefore, knowing that $A = N\delta A$ it follows that $\sigma A \propto (W1 - \psi)$ This expression demonstrates principles which are independent of the use of any model. Under multiple-contact conditions an increase in load can be utilized by increasing the size of existing areas, by bringing into play new areas or by both in different proportions, $\psi$ is a measure of these proportions. If an increase in the load is used mainly to enlarge existing areas, $\psi \rightarrow 1/3$ and the relationship between $A$ and $W$ approaches the single-contact condition $A \propto W^{2/3}$, whereas if the increased load is used mainly to create new areas, $\psi \rightarrow O$ and $A$ becomes nearly proportional to $W$. Dyson and Hirst [17] did demonstrate directly the division of the area into a large number of small areas. They showed that the number of

these contacts increases with increasing load but that their size increases very little. This fact alone is of considerable importance in electrical and thermal contact.

The general problem is to know how to proceed in the contact situation. There are many options. Many investigators such as Archard did not tackle the 3D or areal problem directly but concentrated on the profile problem. The reason for this is that it is simpler: the calculations are tractable. They do not, however, show the true behaviour. Perhaps it is best to investigate the contact properties using the profile method in order to follow the actual path in which understanding grew. (A flow diagram for the contact situation is shown in table 7.3.)

In what follows each of these paths will be looked at briefly after which moving surfaces and their associated tribology will be considered. It is hoped that this will provide a better insight into the overall contact problem. It might demonstrate the nature of the errors and misconceptions that can arise. It should also indicate what needs to be done in order to get some answers, remembering that the surface properties and their subsequent metrology represent only a small part of most investigations and into which not a great deal of effort can be directed.

The complicated flow diagram shows the difficulty of the problem. Ideally, real surfaces should be used and the physical behaviour as a load is impressed should be observed and quantified. Instrumentation to do this directly is still not available. The options still seem to revolve about degrees of fidelity to the real situation. Real surfaces can be mapped and then forced together in the computer. This is what Williamson initiated and other people such as Sayles carried further [15]. It does allow behaviour to be observed, albeit at second hand. Measurements can be taken of the degree of compliance of the two surfaces under load with real surfaces, but what happens on a point-to-point basis is still not known; it can only be inferred by relocation techniques or seeing the changes in electrical conductivity or some other indirect property.

Because of the instrumental failure much effort has been made to approach the problem theoretically. This has been done on two levels, one partial and the other complete. In the former the surfaces have been generated by computer, sometimes as profiles and sometimes areally, and then investigated on the computer in exactly the same way as if they had been real surfaces. This method has certain advantages, one being that extraneous topographic features such as flaws and waviness can be eliminated, so giving a more targeted experiment. It has to be remembered, however, that they are contrived surfaces and that they may not be convincing. The other method is to assume a certain statistical model for the surfaces and then to carry out the loading experiment completely analytically. This may be somewhat unrealistic but it does enable a feel for the problem to be obtained. It is also very difficult to incorporate the peak correlation phenomena. Working point by point in a computer with real surfaces is a closer option.

### 7.2.3.5   Elastic/plastic balance-plasticity index

Greenwood [23] and Greenwood and Trip [24] considerably simplified an understanding of the rough contact problem by showing that one of the contacting surfaces under certain conditions could always be considered flat and smooth. Although the reasoning behind their assertion is questionable, it has allowed progress to be made. The basic theory subject to reservations allows expansion to non-complex bodies to be made where local curvatures can be considered.

It does not matter which surface is assumed to be rough. Greenwood and Trip assumed the configuration shown in figure 7.19. One important fact which emerged from the analysis was that the behaviour of rough contacts is determined primarily by the statistical distribution of the roughness asperity heights and only to a second degree by the mode of the deformation. The idea of such a model is to derive the relationship between the separation of the mean surfaces and the pressure. If the sphere and plane were both rigid and the asperities elastic the separation would be a known function of the radial distance from the centre of the contact and the pressure distribution could be obtained.

**Figure 7.19** Contact sphere and rough surface.

Recently the original work of Archard has been re-examined [18] in the light of fractal analysis. This analysis by Claverella *et al* for elastic contact has shown clearly [19] that following on from Greenwood and Tripp's work [20] fractal description of contact behaviour is unrealistic as the scale dimension gets small. Real area of contact approaches zero and pressures become infinite—much the same as Markov process models of the surface. Archard's simple model assumed that each scale of size could be decoupled from the next, which is not realistic, as the compliance propagation concept developed in equation (7.37) shows. Attempts to deduce a multiscale model of the surface from the spectrum seems to inhibit the credibility of fractal models. Obviously for high values of frequency the power law can be deduced from the fall-off rate. Where any scale of size begins is more difficult. This highlights one of the basic problems of fractals at long wavelengths which is that there is no break point. This in itself precludes the serious use of fractals derived from typical engineering surfaces. The mechanism of production e.g. grinding, automatically has the break point corresponding to the density of graining. Failure to be able to encompass this fundamental fact is a worry. On the other hand the ease with which a Markov approach deals with this break problem leads to the question as to whether the fractal description is really valid anyway [22]. Questions such as how a Weierstrass profile can be produced in engineering are unanswered and yet it is still regarded as useful model!

However, since they are elastic they are deformed and the separation is changed. This leads to a difficult and analytical situation which can only be solved iteratively.

Solutions of the equations of contact for low loads and high loads are given in figure 7.20.

The concept of radius is not straightforward with roughness present—unlike the Hertzian situation. It is clear that roughness becomes less important as the load increases.

The maximum surface pressure usually quoted from Hertzian theory is equal to one and a half times the mean pressure assuming an elliptical pressure distribution. It therefore increases as one-third of the power of the load. When roughness is present at high loads the same is true, but at small loads the pressures are much lower because they are affected by the roughness. For low loads, if the effective radius replaces



**Figure 7.20** Effective radius variation with load: (a) low loads; (b) high loads (after Greenwood and Trip).

the smooth Hertzian radius there is a similar relationship between the mean and maximum pressure as for the Hertzian case.

Maximum shear stress when roughness is present lies on average deeper than for the Hertzian case and is of a much smaller value. An interesting result of this and similar analyses is that, as the load is increased, the effect is to increase the number of microcontacts and not to increase the average size or load on a microcontact. Individual microcontacts obviously have to grow as the load increases to ensure that the *average* contact area is about the same. Also, despite the increase in the size of the contact region due to contact outside the Hertzian zone, owing to the roughness, the total real area of contact is less.

Archard demonstrated that, if it is assumed that the finest asperities support the load, then as the scale is reduced the area becomes effectively independent of load, leading to the statement that the number of contacts is proportional to the load.

The whole issue of plastic and elastic deformation and its implications is not well understood.

Greenwood [25] as usual sums up the difference in a neat diagrammatical way as shown in figure 7.21.

In figure 7.21(a) the contact is plastic and the area of contact is greater than the geometric area. In figure 7.21(b) the contact is elastic and the opposite is true and the area is one-half, as shown in equation (7.20). The reason for this difference derives from the fact that in the plastic case material volume is conserved, whereas for the elastic case it is not a spring and can displace different volumes whereas a plastic article cannot. Greenwood concludes that, even in truly plastic deformation, elastic effects still have an influence. For the case of a contact between a sphere and a plane, taking the contact area to be equal to the geometrical area will rarely be in error by 20%.



**Figure 7.21** Greenwood's plastic (*a*) and elastic (*b*) contact models.

Since Archard a number of contact models have been suggested of varying complexity. All of them have in common the use of the statistics of random processes in one form or another. The realization that most practical contact situations involve random rather than deterministic surfaces is true at least at present—whether or not this state of affairs will continue in future depends on the new manufacturing methods. As previously said, these models give a good indication of the number of contacts and the contact area but little information about the chemico-physical properties within each contact. Both aspects are important in the overall function. (Unfortunately the latter properties are beyond the scope of this book.) Here lies an important message. Surface geometry such as roughness and its metrology plays a vital role in the average statistical behaviour of the contact region. *It does not play such a dominant role in the functional properties of the local microcontacts.* The chemical and physical properties play more of a part in this aspect of contact.

Also the somewhat complex nature of what constitutes a summit or valley in the two-dimensional (areal) cases can add confusion. Nayak, following on from Longuet–Higgins, talks about the way in which the number and nature of contacts changes as a function of height. Simply thinking of a summit count does presuppose that the nature of a summit is understood. Consider an arbitrary plane cutting the surface; a number of possibilities of summits and valleys occur as the level is lowered. According to ideas originally ascribed to Maxwell by Greenwood there is a relationship between the number of maxima, minima and closed contours at any level.

This is considered later (as in figure 7.25) for the cases when plastic flow occurs. Before this some quantitative expressions will be obtained for rough surfaces under load to determine elastic and/or plastic behaviour.

For the effect of curved contact in the presence of roughness there are two situations to be considered. One is when the roughness is considered to deform elastically and the other plastically. Greenwood and Trip [24] and Greenwood *et al* [26] consider the first, and Mikic and Roca [27] the latter. In both cases the maximum central pressure $p(0)$ can be compared with the Hertzian maximum pressure $P_0$.

The reason why such calculations are necessary is because of the fact that it is often necessary to know the maximum stress at a contact involving spheres or circles, for example railway wheels, ball and roller bearings, cams, etc. The presence of roughness can seriously change the estimate. Usually in such considerations the roughness is considered to be of a wavelength which is very small compared with the curved surface. The assumption here is that if long-wavelength components of roughness exist they will modify the estimate of the local curvature at the contact and so be taken into account that way. As it is, such effects are usually ignored. This wavelength effect appears to have been considered only by Whitehouse (reference [55] in chapter 4) and then not for loaded situations. It should be mentioned that the transition from plastic to elastic deformation is part of the 'shakedown' which occurs in repeated loading (7.4.4).

In some earlier work Greenwood and Trip [24] did an elastic analysis as seen in figure 7.22. They assumed a smooth sphere and rough flat. The roughness was considered to be made up of spherical asperities radius $\beta$ and $\eta$ per unit area. The peaks were assumed to be Gaussian and having a standard deviation of $\sigma$. They determined two independent non-dimensional variables $T$ and $\mu$ where

$$T = \frac{2P}{\sigma E (2R\sigma)^{1/2}} \qquad \mu = \tfrac{8}{3}\sigma\eta(2R\beta)^{1/2}$$

(7.41)

where the variable $T$ is meant to be governed by load considerations and $\mu$ by the geometry of the surfaces.

Later a surface parameter $\alpha$ was introduced. This was also non-dimensional

$$\alpha = \frac{\sigma R}{a_0^2}$$

(7.42)

where $a_o$ is the contact radius given for smooth surfaces as

$$a_0 = \left(\frac{3PR}{4E^*}\right)^{1/3}$$

(7.43)

and where $1/R = 1/R_1 + 1/R_2$ and

$$\frac{1}{E^*} = \frac{1 - v_1^2}{E_1} + \frac{1 - v_2^2}{E_2}.$$

(7.44)



**Figure 7.22** Contact of rough surface with sphere.

The two surface are flattened by $d$ within the contact area where

$$d = \frac{\alpha_0^2}{2R_1} + \frac{\alpha_0^2}{2R_2} + \frac{\alpha_0^2}{2R}.$$

(7.45)

Here in equation 7.42 $\sigma$ is given by

$$\sigma = \left(\sigma_1^2 + \sigma_2^2\right)^{1/2}$$

(7.46)

where it is assumed that, in the general case, both surfaces are rough and curved.

The smoothed-out pressure distribution calculated for two different values of $\alpha$ (large $T$) where the asperity deformation (assumed elastic) is small compared with the bulk deformation, which in this case is very nearly Hertzian, and the case where $\alpha$ is large ($T$ small) and the effect of the asperities becomes highly significant, the contact pressure is smaller and spaced over a wider area (figure 7.21(a)). In figure 7.23 the central pressure $p(0)$ is plotted versus $\alpha$ as a fraction of the maximum Hertz pressure $P_0 = (6PE^{*2}/\pi^3 R^2)^{1/3}$ for the two values of $\mu$ which probably encompass most practical surfaces. From these figures it is clear that the effect of surface roughness on the maximum pressure is governed more by $\alpha$ than $\mu$. The detailed geometry of the surface seems to have a secondary effect. If the peak distribution is taken to be exponential rather than Gaussian then it can be shown that the pressure distribution is a function of $\alpha$.



**Figure 7.23** Greenwood's contact parameters $\mu$ and $\alpha$.

For rough surfaces it should also be noted [9] that the first yield would be expected at a depth of about $0.4a^*$ under a contact pressure of $p(0) = 3.5k$, where $k$ is the yield strength in shear.

Here $a^*$ is the effective contact radius which has to be defined because a Gaussian height distribution allows infinite peaks. Greenwood defines it as

$$a^* = \frac{3}{8} \int_0^\infty rp(r)\mathrm{d}r \bigg/ \int_0^\infty p(r)\mathrm{d}r$$

(7.47)

so chosen to make $a^* = a_0$ where the pressure distribution is Hertzian.

For a treatment in the case when the asperities deform plastically rather than elastically [27] the surface height is taken to have a Gaussian distribution rather than the peaks. Also the real contact pressure is taken to be constant and is assumed to be equal to the hardness $H$ of the softer of the two materials. Then the pressure distribution is found to be primarily governed by the non-dimensional parameter where $\bar{p}_0 = 2p_0/3$ and is the mean Hertz pressure. To a lesser extent $p(r)$ is also influenced by $H/\bar{p}_0$.

From Hertz theory

$$\bar{\sigma} = \frac{\pi \sigma E^*}{a_0 \bar{p}_0} \tag{7.48}$$

The variation of $p(0)/p_0$ from the above plastic expressions is plotted in figure 7.23.

This graph is fundamental. It shows that under these circumstances the mode of deformation of the asperities is unimportant! The actual contact pressure is dictated by $\alpha$. From this, if $\alpha < 0.05$ the prediction of the contact stress by Hertz theory is not likely to be out by more than 5%, which is more than adequate for most purposes.

In the case when the surfaces are non-spherical the contact zone can be elliptical rather than circular and an equivalent spherical case is defined such that $R_e = (R'R'')^{1/2}$ *to* enable the same formula for $\alpha$ to be used. Other shapes such as rough cylinders and annuli have been investigated assuming plastic deformation of the asperities [28] and they have shown similar spreading of the contact area, as above, relative to Hertz theory. Some of the assumptions are, however, somewhat suspect because of misconstruing the relationship between displacement and contact ratio via the material ratio curve.

One big advantage of using non-dimensional parameters as the variables in this case is that the problem can easily be scaled in size to cater for large wheels or miniature ball bearings.

From this type of analysis it has emerged that for practical cases involving a large radius of curvature and small roughnesses the texture is not too important from the point of view of maximum stress. It is, however, of fundamental importance in the case of contact points and individual contact areas which are important in electrical and heat conduction. It is the presence or absence of roughness that tends to be more important.

One other point that has been generally ignored is the fact that because contact is often not at the top of asperities the contact forces are not perpendicular to the mean planes of other surfaces but normal, locally to the point of contact.

Rossmanith [29] has considered this aspect for elastic contacts using Hertz–Midlin theory. He considered the cases of no slip and slip-afflicted cases on the vertical force to see how it changed as a result of these lateral forces. Not surprisingly he found that the mean bending moment, torsion moment, shear and vertical force per unit area vanished if the contact zone is circular; however, they turn out to be also zero for arbitrarily shaped asperities. In general, the effect of this is to reduce the effective normal load somewhat.

The starting point for the introduction of roughness is the case of a rough flat surface contacting a smooth flat surface [23]. Assume that the rough surface has asperities on it whose probability density as a function of height is $p(z)$. This means that the probability density of a peak being of height $z$ above the reference plane of the smooth surface is $p(z)$. The number of peaks having this height will be $Np(z)$ where $N$ is the number of asperities on the rough surface. Consider the case when the other surface is lowered onto it assuming no peak correlation. When the separation is $d$ there will be a contact at all asperities with heights greater than $d$ (figure 7.24). There will be distortion of the asperities and



**Figure 7.24** Approach of rough and smooth surfaces.

the plane to accommodate the excess height *(z − d)*. Therefore the relationship governing the deformation of a single contact in terms of the extra height absorbed in the deformation is needed, that is the local compliance *w*. This will depend on the asperity shape and the mode of deformation. For example, as previously explained for purely elastic conditions, the Hertzian relations can be used. Thus

$$P \propto w^{3/2} \quad \text{and} \quad A \propto w \tag{7.49}$$

or more generally

$$P = h(w) \quad \text{and} \quad A = g(w) \cdot \tag{7.50}$$

To find the total load and the total area therefore the excess height $z − d = t$ is substituted for *w* and the condition summed for all load-carrying asperities. Giving the geometrical results

$$P = \int_0^\infty h(z-d)Np(z)\mathrm{d}z \qquad A = \int_0^\infty g(z-d)Np(z)\mathrm{d}z \cdot \tag{7.51}$$

So, in principle, once the distribution of heights of the asperities and the mode of deformation of a single asperity are known then the total load and the total area can be found as a function of the separation of the surfaces. Also the number of apparent contacting asperities will be given by

$$n = \int Np(z)\mathrm{d}z \cdot \tag{7.52}$$

Here the separation can be eliminated to get the relationship between the number of contacts and the load. Greenwood points out [23] that this is where the problems start even ignoring correlation.

Workers have made many assumptions of the distribution of height, ranging from uniform, linear exponential, Gaussian and log normal, and the asperities have been wedges, cones, spheres or rods, whilst the deformation modes considered are elastic and ideal plastic and sometimes elastic plastic. Obviously these represent an enormous number of combinations, all of which cannot be correct! There is at least one redeeming factor which can be derived from working with this simple model of a flat smooth surface against a rough one, and this is that issues about the properties of the gap between surfaces and the properties of the peaks are one and the same.

When the earliest theories were put forward concerning contact theories the distribution of heights was chosen mainly for ease of calculation (e.g. [30]). Recently Greenwood and Williamson [31], Bickel [32], Tallian *et al* [33], Pesante [34] and Ling [35] found that the peaks were often Gaussian or at least close. What is more often the case is that the actual surface heights (or ordinates) are Gaussian and that the peaks have a distribution which is slightly different, for example the gamma distribution.

One special case can be used to highlight some of the effects as a load is imposed on the surface, and this is the exponential

$$p(z) = \frac{1}{\sigma}\exp\left(-\frac{z}{\sigma}\right)\cdot \tag{7.53}$$

Substituting this into (7.52) and (7.51) gives

$$P = N\int_0^\infty \frac{1}{\sigma}\exp\left(-\frac{z}{\sigma}\right)h(z-d)\mathrm{d}z = N\exp\left(-\frac{d}{\sigma}\right)\int_0^\infty \frac{1}{\sigma}\exp\left(-\frac{t}{\sigma}\right)h(t)\mathrm{d}t \tag{7.54}$$

$$A = N\int_0^\infty \frac{1}{\sigma}\exp\left(-\frac{z}{\sigma}\right)g(z-d)\mathrm{d}z = N\exp\left(-\frac{d}{\sigma}\right)\int_0^\infty \frac{1}{\sigma}\exp\left(-\frac{t}{\sigma}\right)g(t)\mathrm{d}t \tag{7.55}$$

and

$$n = N \exp\left(-\frac{d}{\sigma}\right) \tag{7.56}$$

where the separation $d$ appears only in the exponential term. Eliminating it gives the interesting result that

$$A \propto P \text{ and } n \propto P \tag{7.57}$$

Hence for this special distribution the mode of deformation and the shape of the asperities *do not* influence the load, area or number of contact relationships with the result that there will be exact proportionality between area and load. This is what would be expected for a purely plastic regime.

So, according to this, there are two possible ways of looking at Amonton's laws. Either there is plastic flow which is independent of surface asperity shape or there is an exponential distribution of peaks and any mode of deformation.

What could happen in the latter case is that as the flat, smooth surface approaches the rough one existing contact spots will tend to grow, but at the same time new, small contacts form. When the process of growth and formation balance the average spot size is constant and independent of load. All that changes is the number of contacts.

In practice there is not likely to be either a completely plastic regime or an exponential distribution of peaks; rather the situation will be somewhere in between. However, the exercise of considering the exponential distribution is useful in that it does illustrate one extreme of behaviour and can therefore usefully be considered to be a reference. Greenwood [23] uses this idea extensively when attempting to get a feel for the usual situation where the contacts between a rough surface and a flat are a mixture of elastic and plastic.

Consider a simple sphere on a flat. For any compliance $w_p$ the deformation is elastic; for $w > w_p$ it becomes plastic.

For a sphere on a plane

$$w_p = \beta\left(\frac{H}{E^*}\right)^2 \tag{7.58}$$

where $H$ is the hardness $\beta$ is the radius of the sphere and $E^*$ is the effective modulus:

$$\frac{1}{E^*} = \frac{1 - v_2^1}{E_2} + \frac{1 - v_1^2}{E_1}. \tag{7.59}$$

Two surfaces in contact will, because of the wide range of asperity heights, always have a certain percentage of plastic deformation so that it is meaningless to consider that the condition for plasticity is achieved when all peaks are plastically deformed. The nearest it is possible to get is to consider a certain proportion of peaks to be strained past the elastic limit. In the special case of the exponential distribution this leads to the result that the ratio is $\exp(-w_p/\sigma)$; that is, it is a constant independent of load. It therefore depends only on the topography which, as has been pointed out, is contrary to the intuitive belief that an increase in load results in an increase in the proportion of plastically deformed peaks. This result, however, does match up to the concept that only the number of peaks is increased.

The criterion for whether a surface gives elastic contact, plastic contact or lies in between—the load-dependent range—can be expressed as a function of $w_p/\sigma$ where $\sigma$ is the standard deviation of the peak distribution. Greenwood expressed a 'plasticity index' $\psi$, in terms of $E'$, $H$, $\sigma$ and $\beta$ as

$$\psi = \frac{E}{H}\sqrt{\frac{\sigma}{\beta}}. \tag{7.60}$$

This concept of a transition from mainly elastic to plastic as a function of the surface topography was profound because it began to explain the nature of running in. Although Archard postulated the idea that running in was complete when all the contacting asperities became elastic rather than plastic, he never quantified this idea. The basic idea of Archard that run-in surfaces were elastic put the study of surface topography firmly back in favour in tribology from the rather limited importance given to it previously by Bowden and Tabor [36].

Another route was to consider the criterion for elastically compressing an asperity down flat. One such criterion related the mean slope $m$ to the hardness and the elastic modulus. Thus if

$$m < kH/E^*$$

(7.61)

then no plastic flow would occur [37,38], where $k$ is a constant taking values between 0.8 and 2.0 depending on asperity shape.

Other indices include that by Mikic and Roca [27] who introduced $\psi_3$ where

$$\psi_3 = \frac{E^*}{H} \tan\theta$$

(7.62)

where $\tan\theta$ was the mean slope of the asperity. $\psi_3$ is such that at least 2% of the asperities deform plastically when 4, 3 >~ 0.25. This model has been resurrected by Greenwood [39] in a tentative look at pressure variations within the contact zone.

A more complete theory for the plasticity index was given by Whitehouse and Archard [40] who considered a distribution of peak curvatures as well as heights. Equation (7.51) then becomes

$$\psi = \frac{E}{H} \frac{\sigma}{\beta^*}$$

(7.63)

where $\sigma$ is the RMS height of the surface and $\beta^*$ is the correlation length or independence distance, not the peak radius. Another refinement by Onions and Archard [41] put

$$\psi = 0.69 \frac{E}{H} \frac{\sigma}{\beta^*}$$

(7.64)

directly, indicating that there were more plastic contacts than for the Greenwood index for the same numerical value of $\psi$. Upon reflection [40] this became obvious when it was realized that the higher peaks are necessarily sharper on average than the lower ones. This means that they crush more easily.

The importance of this discussion on the behaviour of the asperities under load cannot be over-emphasized because it relates to tribological function in wear, friction and lubrication as will be seen.

It is not too important which plasticity index is used. The best one will be dependent on the particular application and will probably be found out empirically. What is important is the realization that there is always a possibility of regimes in which plastic and elastic conditions are present simultaneously. In terms of the Greenwood index this is in the region where $\psi$ takes values between 0.6 and 1.0.

Other topographic indices have been suggested. One by Gupta and Cook [42] defines a topographic index $\zeta$ essentially as a measure of the ratio of the nominal to real contact area of two surfaces in contact. Thus

$$\zeta = \frac{1}{\sqrt{2D_2}} \frac{1}{R_m} \left( \frac{E^*}{H} \right)$$

(7.65)

where $D$ is the average number of peaks per unit area (or the smaller of the two densities when the surfaces are different).

$$R_m = \frac{R_1 R_2}{R_1 + R_2}.$$
$$(7.66)$$

In their analysis the distribution of radii was log normal [42].

Nayak [43] moves away from isotropic surfaces to those in three dimensions in which a lay was present. Thus

$$\psi = \frac{E}{H}\sigma^*$$
$$(7.67)$$

where $\sigma^*$ is the RMS value of the differential of that component of a surface in contact which has a narrow spectrum. He shows that if $\psi > 5$ then most asperities will deform plastically.

Other workers have approached the problem somewhat differently using the spectral moments of a profile, $m_0$, $m_2$, $m_4$, described earlier. They determine the probability density of peak curvature as a function of height. Then, assuming that a profile comprises randomly distributed asperities having spherical tips, they estimate the proportion of elastic to plastic contact of the profile on a flat as a function of separation. Thus if $R_p$ is the maximum peak height and $h$ is the separation of the flat from the mean plane, the probability of elastic deformation becomes $P_{el}$ where

$$P_{el} = \text{Prob}[c < (H/E)^2/(R_p - h)] \quad \text{given that } z^* \geqslant h$$
$$(7.68)$$

so that in this context the important parameter is that which describes how curvature C changes with height, that is $f(c/z^*>/h)$.

Unfortunately such height-dependent definitions are impractical because $R_p$, the peak height, is unbounded; the larger the sample the bigger is $R_p$.

One comprehensive index has been derived by Francis [44] who in effect defines two indices $\psi_1$ and $\psi_2$

$$\psi_1 = \frac{\sigma'}{k^{1/4}} \frac{E_1^*}{P_m}$$
$$(7.69)$$

$$\psi_2 = \frac{\sigma}{k^{1/2}} \frac{E_1^*}{P_m}$$
$$(7.70)$$

where
$$k = \sqrt{1.5} \frac{\sigma'^2}{\sigma \sigma^*}$$
$$(7.71)$$

which is similar to that of Nayak in that $\sigma''$, the RMS curvature, is included ($\sigma'$ is the RMS slope and $\sigma$ is the RMS of the gap between the surfaces). In general $\psi_1$ has been found to be the most sensitive of the two indices. He finds that for $\psi_2 < 4$ there should be elastic conditions and plastic when $\psi_2 > 14$.

More recently generalized indices have been proposed [45] which try to take isotropy into account. These authors arrive at a plasticity index in which the isotropy $\gamma$ is included. When written in terms of the Greenwood index this becomes

$$\psi = \left(\frac{1}{w}\right)^{1/2} = \frac{\pi}{2\sqrt{2k}}\left(\frac{1+\gamma}{K(e)H(e)}\right)^{1/4}\frac{E^*}{H}\left(\frac{\sigma^*}{R_x}\right)^{1/2}$$
$$(7.72)$$

where $K(e)$ and $H(e)$ are elliptical integrals for the pressure distribution, $\gamma = R_x R_\gamma$, the ratio of the radii of curvature in the two directions, is taken to be orthogonal and $k$ is a function of the yield stress and approximates to unity over a range of $\gamma$. The interesting result from this work is that for surfaces which have $\gamma < 10$ the plas-

ticity index can reasonably be approximated to

$$\psi = (\psi_x \psi_y)^{1/2} \tag{7.73}$$

where $\psi_x$, $\psi_y$ are the corresponding conventional plasticity indices in the $x$ and $y$ directions. It must be presupposed that the $x$ and y directions are taken to give $y$ a maximum value.

A number of points need to be made here about the concept of plasticity indices. In principle the idea is good—the idea is that of providing a non-dimensional number which gives a measure of predictability to a functional requirement, in this case contact and wear. The fact that this number is arrived at by reference to both topographic and material properties has to be applauded because functional behaviour cannot be divorced from either. However, once having accepted the validity of the idea the practical implementation of it is much more questionable. Putting aside the material properties until later, consider the topographic features. There seems to be a variety of parameters basically involving either first differentials $\overline{m}$ tan $\theta$, $\theta'$, $\sigma/\beta^*$, etc, or second differentials $\sigma''$ $1/R$, and sometimes a height. This state of affairs is theoretically unsatisfactory.

An approach to the question of isotropy has been made by Nayak [43] using random process analysis. Similar, more extensive coverage has also been made by Bush *et al* [46, 47] Nayak makes an interesting comparison of the onset of plastic flow as shown in figure 7.25. The diagrams show that for loads of $0.10 - 0.4$ kg mm$^2$ the ratio $A_p/A_c$ *is* less than 0.02 for $\psi_1 < 0.6$ and greater than 0.02 for $\psi_1 > 1$. Here $A_p$ is plastic area and $A_c$ elastic area $C = 1 - 0.9/\alpha$ ($\alpha$ is Nayak band width parameter).

Bush *et al* eventually derives a plasticity index in terms of the moments of the spectrum in the case of isotropic surfaces:



**Figure 7.25** Onset of plastic flow for different models (after Bush). (*a*) Greenwood and Williamson, (*b*) Nayak, (*c*) elliptical, (*d*) Onions and Archard, (*e*) anisotropic.

$$\psi = \left(\frac{8}{3}\right)^{1/2} \left(\frac{E^*}{H}\right) \left(\frac{m_0 m_4 C}{\pi}\right)^{1/4}. \tag{7.74}$$

For anisotropic surfaces the result requires five parameters: the variance of height $m_{00}$, the two principal mean square slopes $m_{20}$ and $m_{02}$, and the two principal mean square curvatures $m_{40}$ and $m_{04}$. From two profiles, one along the lay and one across it, the moments can be found by using peak densities and zero crossings (see chapter 2).

Using the Van Mises yield criterion and surface parameters mentioned above gives the result for $A_p / A_c$ in figure 7.25(e). In fact it is clear that whatever the form of analysis for a given load and numerical value of $\psi$, the degree of plastic deformation is approximately half that of the corresponding isotropic surface, because although the number of peaks is smaller along the lay their curvature is proportionately larger and hence they are more elastic (see figures 7.26 – 7.28).



**Figure 7.26** Variation of $A_c$ with $d$.



**Figure 7.27** Variation of $P$ with $d$.

**Figure 7.28** Variation of $A_c$ with $d$.

It is true to say that often the idea of a plasticity index makes more sense when the separation of the surfaces is relatively large (large means that there is little contact). The reason for this is at least twofold. One concerns the assumption which is almost always made that the asperities are independent of each other, the other is that the forces are normal at the contact. Neither is strictly true. Nayak has relied heavily on the observations of Pullen and Williamson [48] that when an incremental deformation of the asperity tips occurs the incremental plastic volume reappears where there is no contact in the case where the surface is constrained this effect may inhibit asperities collapsing with the application of the load. Pullen and Williamson call this asperity persistence. It has the effect of making the surface appear to be harder than it is.

What emerges from the analysis of plasticity under this regime is a picture of contact which exhibits more larger contacts than would be expected—the extreme condition predicted by Archard. The density of contacts is therefore considerably lower, except at lower pressures than would be predicted.

So the situation is that the surface peaks appear to have a higher stiffness—in the case of elasticity—according to the Whitehouse model described early in this chapter. It also seems to have greater resistance to



**Figure 7.29** Variation of $A_c$ with $P$.

the movement of material—the plastic case—according to the model of Nayak/Pullen [33]. It may be that both mechanisms occur thereby enhancing asperity persistence [48].

Contact patches are perforated with holes, the density of which is about the same as that of the contact patches. Another feature likely to arise is that the average patch is likely to be non-circular except at light loads. This fact is of considerable importance in thermal and electrical contact.

Nayak investigated the nature of plastic contacts in the same way as Longuet–Higgins examined areal contours and comes to some potentially interesting conclusions.

If closed contours bounding finite contact patches are considered then the situation could be as seen in figure 7.30(a).

The presence of holes is guaranteed by having maxima below the mean plane of the surface. Nayak goes on to show that the presence of such holes is likely if the spectrum of the surfaces is wideband and Gaussian. Because many surfaces fall into this category holes must be expected in plastic conditions rather than be considered a rarity.

It should be remembered that it is not possible to consider how the parameters change as a function of separation as if they were independent of the behaviour under load because the loading distorts the model! The contacts are spatially correlated.

The real problem is that the $H$ values supporting the $P$ peaks are not recognized by just counting the number of asperities yet these play an important role in the peak correlation.

In what follows it is inherently assumed that correlation effects are negligible.



(a)

No (max) + No (min) − No (saddle) = 1

(b)

No (max) + No (min) − No (saddle) > 2 − No (holes= 1)

(c)

$$\text{Dens (max)} + \text{Dens (min)} - \text{Den (saddle)} + \text{Den (holes)}$$
$$>z \qquad >z \qquad >z \qquad >z$$
$$=\underset{z}{\text{Den}}\ \text{(patches)}$$

**Figure 7.30** Nayak behaviour of areal contour patches (No means number).

### 7.2.4 Effect of waviness on contact.

Seabra and Berthe [49] have investigated the effect of waviness on the pressure distribution in the contact zone. They use a rigorous tensor analysis of the stress fields to assess the effect of surface geometry. Because of the complexity only a few results are obtained. Nevertheless, despite the fact that the rough surface is computer-generated and takes a raised cosine form, some useful results can be obtained as shown

in figure 7.31. Not surprisingly, the pressure points above the Hertzian pressure distribution correspond with the peaks of the waviness. A measure of the pressure concentration factor PCF is



**Figure 7.31** Effect of waviness on contact.

$$\text{PCF} = C_1 \left( \frac{\lambda}{A_w} \right)^{\alpha} \left( \frac{a}{R} \right)^{\beta} \left( \frac{\lambda}{a} \right)^{\gamma} \tag{7.75}$$

with $C_1 = 4.3884$, $\delta = -0.4234$, $\beta = -0.4204$ and $\lambda = -0.0155$. $\lambda$ is the waviness wavelength, $A_w$ the amplitude, $R$ the radius of smooth body and $a$ the Hertzian radius

The reduction in area of contact, *RA,* is similarly

$$RA = \ln \left[ C_2 \left( \frac{\lambda}{A_w} \right)^{\theta} \left( \frac{a}{R} \right)^{\eta} \left( \frac{\lambda}{a} \right)^{\psi} \right] \tag{7.76}$$

with $C_2 = 0.9404$, $\theta = 0.3802$, $\eta = 0.3349$ and $\psi = 0.0395$.

In both cases the correlation between the observed results and theory in (7.75) and (7.76) is 0.96.

Similar expressions result in the *y* direction. Roughness in the presence of waviness is qualitatively the same as the waviness, although here care has to be exercised to identify any plastic deformation.

Thus, the effect of waviness in combination with roughness is:

(1) modification of the pressure distribution in the contact zone;
(2) existence of several points of maximum pressure instead of one (usually corresponding to waviness peaks);
(3) several contact areas instead of one;
(4) maximum pressures higher than the Hertzian;
(5) overall contact area equal to or smaller than the Hertzian.

The most important parameters in the case of waviness are, again not surprisingly, wavelength, amplitude and load. Both the pressure concentration and the reduction of area can be expressed in terms of these parameters.

### 7.2.5 Non-Gaussian surfaces and contact

In a lot of analysis it is assumed that the contacting surfaces are Gaussian in character. For many abrasive processes this is true. Also it should be remembered that even if the contacting surfaces are not Gaussian the gap between them is more likely to be. However, there are occasions where asymmetrical surfaces are involved. These are often produced in multiprocessed surfaces such as in plateau honing.

The starting point is to generate a sequence of random numbers having either a uniform distribution or a Gaussian distribution. It is easy in practice to generate any required distribution. Take for example use of the uniform distribution.

Two random numbers $n_1$ and $n_2$ are generated which serve as an address on the distribution curve (figure 7.32). On this curve lies the desired distribution. Simply only accepting points which lie inside the curve (shaded) gives the desired distribution. Points which lie outside the curve are rejected. This, after many points $P(n_1n_2)$, gives the amplitude distribution. Obviously the initial random numbers are independent. This corresponds to white noise. To generate a desired autocorrelation function the accepted sequence $P(n_in_j)$, simply Fourier transform it multiply it by the transform of the desired autocorrelation and then retransform back to the spatial domain.



**Figure 7.32** Complete axis

For example if an exponential autocorrelation $\exp(-\alpha x)$ is required. The transform of $P(n_in_j)$ (i.e. $F(w_1w_2)$ is multiplied by $\dfrac{\alpha}{\alpha^2 + w^2}$ and then invert-transformed back.

For an areal rather than profile surface three random numbers per point are used $n_1$ $n_2$ $n_3$. To get a specific skewness $SK_i$ and kurtosis $K_i$ on the input, the Johnson translator method is used [2] [50].

This relates to the required output $SK_0$, $K_0$ by

$$SK_0 = SK_i \left( \sum_{i=0}^{i_{in}} t_i^3 \right)^2 \Bigg/ \left( \sum_{i=0}^{i_{in}} t_i^2 \right)^3 \tag{7.77}$$

$$K_0^{\frac{1}{2}} = \left( \sum t^4 K_i + \sum_{i=0}^{i_{in}-1} \sum_{j=i+1}^{i=in} t_i^2 t_j^2 \right) \Bigg/ \left( \sum_{i=0}^{i=in} t_i^2 \right)^2 \tag{7.78}$$

Given that it is possible to generate the required surface it becomes necessary to investigate its contact properties. Bhushan [51] produced some contact properties. A useful physical constraint was employed to minimize the potential energy which is allowable for elastic effects. How mixes of elastic and plastics contacts are dealt with is not clear. The results indicated that surfaces with slight positive skew and kurtosis of over 4 yielded an optimum surface with minimum contact area. It was concluded that negative skew coupled with low kurtosis ($< 3$) resulted in severe friction, whereas very high kurtosis and $R_q$ tend to plastically deform. In this work the presence of surface liquid masked the generality of the results.

There have been a number of attempts to model the non Gaussian surfaces, sometimes incorporating load parameters. Many investigators are reluctant to veer from the Gaussian model of the surface. One reason is that it is easy to superimpose Gaussian characteristics, (e.g. to find the gap properties) as a function of the two contacting surfaces. For bimodal surfaces and multiprocesses where the skewness of the surface is not zero, other models have to be used [50]. Probably the first in 1978 was due to Whitehouse [52], who used beta functions, McCool who used the Weibull function [53] Adier and Firman [54] who used a model based on the $x^2$ (Chi squared) distribution.

McCool tried to compare the Weibull surface with the Gaussian surface. He concluded that the Weibull function very nearly had the additive properties of the Gaussian distribution but could deal with non-zero skewed surfaces well.

He found that the mean asperity pressure can be higher than the Gaussian surface, having the same RMS value, by 1.4 if the skew is +1 and lower by 1.7 if the skew is −1.0. As the range of skew here is not very high the differences in pressure seem to be very high.

### 7.2.6 *Fractal and contact*

Fractal theory is described in chapter 2. In principle fractal dimensions can be used wherever spectral analysis is used. Fractal dimensioning is jut another characterization of the surface. Usually to be useful fractal characterization has to extend over a wide range of size.

Evidence of fractal properties involve a power law of some phenomena

$$P = Kr^{f(D)} \tag{7.79}$$

$K$ is constant, $r$ is a characteristic length scale of the system and $D$ is the fractal dimension. Strictly for the system to be fractal the power law above should be valid at all length scales. However, there is evidence of fractal behaviour over typically less than two decades. The power law could well be useful over the range but its limited range does not necessarily indicate fractal behaviour. Claiming fractal behaviour over a few octaves is not significant. Neither is claiming that a corner or break in the spectrum is evidence of 'multifractal' behaviour.

To recapitulate, fractals have some strange properties, for example, the differentials do not exist, then are often tending towards infinity. A true fractal looks the same and has the same properties at any scale thereby tempting investigators to work at a convenient size (e.g. micrometres) and to interpolate their results down to nanometric or even atomic levels where experimentation is very difficult. The problem here is that the actual properties (e.g. friction or contact), change with size! What is friction at the atomic level? Fractals having the scale invariant properties are said to be 'self similar'. These fractals do occur in nature but are rarely man made. Most fractals require another factor which ensures the self similarity. This is called the topothesy and is added when the scale changes to keep the property of self similarity.

Fractals which require the scale are said to be 'self affine'.

In a review of fractal applications Thomas [2] recalls the efforts of a number of investigators.

Normally engineering surfaces do not have fractal properties. The property of a cutting process together with the tool path are not fractal. In fact with single and multiple tooth cutting like turning and milling, and abrasive processes such as grinding the process is definitely not fractal. Fractal effects could occur in the material cracking produced by the tool and in fact are Markov processes. However, there are cases, such as in magnetic tapes which are coated, where fractal behaviour is seen. This is mainly due to the fact that a process which is deposited onto a surface (i.e. in the form of a coating) is similar to the characteristics of growth mechanisms and therefore is likely to have fractal properties [56].

In some cases [15] it has been possible to show that the area of contact and the relative proportions of elastic and plastic deformations are sensitive to fractal parameters. What has been found, however, is that fractal mechanics tend to underestimate the contact area [56] It may well be because of the influence of asperity correlation mentioned earlier.

On an intuitive level it seems unlikely that manufactured surfaces have fractal properties. What has been reported so far is unimpressive [57] and [58]. It could be that phenomena such as wear or damage [59] which grow in time could develop fractal properties but contact situations seem unlikely.

### 7.2.7 *Coated contact*

The problem of elastic and sometimes also plastic contact between non-conforming coated surfaces has received much attention in recent years. One good example of their application is in wear protection.

The coated analogue of Hertz equations can be solved using layered body stress functions and the Fourier transform method [60]. A neat method for finding that contact dimension and the maximum pressure has been developed by Oliver [61] but perhaps the most practical method has been due to Sayles [62] who manages to include, albeit with a numerical method, the effect of roughness.

When one or both surfaces in contact are coated a somewhat different situation arises because this is now a three body situation rather than the conventional two body. So instead of transitions from elastic to plastic being straightforward the introduction of the third body—the film—presents problems [63] the basic result is that for light loads there is a considerable difference between coated and uncoated surfaces.

For those asperities making contact only with the film it is the coating, being soft, which deforms plastically. The substrate has the effect of stiffening the coating.

## 7.3 Functional properties of contact

### 7.3.1 General

Many functions devolve from contact mechanisms. Some involve static positioning and the transfer of one sort of energy normal to the surface, as in thermal and electrical conductivity. Some involve a small normal displacement, as in stiffness where the transfer is normal force, yet other cases involve small movements over a long period, as in creep, and others fast or slow transverse movement, as in friction, lubrication and wear.

In practically all cases of functional prediction the basic approach has been to split the problem into two. First the behaviour of a typical contact is considered and then the behaviour of an ensemble of such contacts taken over the whole surface.

The starting point has therefore usually been to define each surface in terms of some sort of asperity model as shown in figure 7.33, and then to assign values to its features such as radius of curvature, presence or absence of surface films such as oxides, the elastic model and hardness.



**Figure 7.33** Contact of the asperities.

The behaviour is then considered with respect to movement and load. Such contact behaviour is complex and beyond the scope of this book. However, there are texts which deal in great detail with such matters: (see for example the definitive book by K L Johnson [64]). It suffices to say that the properties of asperity contact are investigated classically in much the same way as in section 7.2, which deals with perfectly smooth, simple geometric bodies.

In the case of simple loading the mode of deformation was considered paramount—whether the contact would behave elastically, plastically or as a mixture of the two. When sliding or lateral movement is considered other issues are present, for example whether or not stick-slip and other frictional properties exist.

This philosophy of using such asperity models has been useful on many occasions but it does suffer from one drawback. This is that real surfaces are not comprised of lots of small asperities scattered about, each independent of the other; the surface is a boundary and can be considered to be continuous. Random process methods (see e.g. [65]) tend to allow the simplistic model to be replaced by a more meaningful picture-sometimes! Even then the philosophy has to be considered carefully because surface contact phenomena are principally concerned with normal load behaviour and not with time, as in electrical signals,

so that top-down characteristics tend to be more important than time or general space parameters most often used in random process analysis.

The first and most important assumption using simple asperity model is the nature of the geometry of the typical asperity.

Having decided on the asperity model involved in the contact the next step involves an investigation of how these contacts will be distributed in space. This involves giving the surfaces asperity distributions. From these the density of contacts is found, the mean height of the contact and the real area of contact, to mention just a few parameters. This brings in the second main assumption. What is this distribution? Various shapes have been tried such as the linear, exponential, log normal, beta function and so on. More often than not the Gaussian distribution is used because of its simple properties and likelihood of occurrence. It can be applied to the asperity distribution or the surface height distribution. Often it is applied to the gap distribution. In any event the Gaussian distribution can be considered to be a useful starting point, especially where two surfaces are concerned. Even if the two mating surfaces are badly non-Gaussian their gap often is not.

Another technique is emerging which does not attempt to force asperity shapes or distributions into analytical treatments of contact. Instead actual surface profiles are used (see [66]). With this technique digital numerical records of real surfaces are obtained. The normal movements of the opposing bodies are made on a computer, which ultimately results in overlap between them. This overlap is then equated to the sum of the strains created by the compliances of a series of adjacent pressure 'elements' within the regions of apparent contact. Usually the pressure elements are taken to be the digitized ordinates of the surfaces. This is in effect utilizing the discrete method of characterizing surfaces advocated by Whitehouse, rather that the continuous method. The sum of all these pressurized elements, none of which is allowed to be negative, is equated to the total load in an iterative way until the accumulated displacements at each element equate to the rigid-body movement imposed and the sum of the corresponding pressure steps needed to achieve this displacement equals the nominal pressure. This method, usually called the *numerical contact method* (see [67]), does not rely on a model of the topography. Also the contact geometry is fully defined across the whole interface. The technique does imply that the behaviour of each of the 'elements' is independent, which is questionable at the sample spacing of digitized surfaces. However, if the element width is taken as $\beta^*$, the autocorrelation length, then this criticism could be relaxed.

If the force-compliance characteristic for an 'element' is given by $W(z)$ and the probability density function of ordinates is $p(z)$, the effective load-carrying capacity $L$ of a surface $f(z)$ being pressed by a flat is given simply by equation (7.80) where $A$ is the height at first contact:

$$L(z) = \int_z^A p(z')W(A-z')\mathrm{d}z' \qquad (7.80)$$

which is very easily expressed as the characteristic function of the height distribution multiplied by the characteristic function of the load-compliance curve!

It may be that this sort of approach will become used more often because of its simplicity and freedom from ensnaring assumptions. To resort to this method does not say much for the theoretical development of the subject over the past few years!

*7.3.2 Stiffness*

Early work concerned with the dynamic vibrational behaviour of machine tools [68, 69], realized the importance of the stiffness of bolted joints in the construction of the machines. In this work it was often assumed that the asperities were distributed according to a parabolic law and that the surfaces were made up of cusplike asperities. The average behaviour of the asperities was then based upon the material ratio (bearing area) curve for the surface. Such approaches were reasonable for surfaces especially manufactured to have the cusp-like characteristics but not for surfaces having random components of any significance.

Thomas and Sayles [65] seemed to be the first to attempt to consider the stiffness characteristics of random surfaces. They replaced the deterministic model for the surface, and a process mark and waviness, by a continuous spectrum of wavelengths. In this way they transformed a simple deterministic model of the surface into a random one in which overlapping bands of wavelength from a continuous spectrum had to be considered.

Nowadays other features of the joint which are concerned more with friction are also considered to be important in the stiffness of machine tools and any other structural joints because damping is produced.

Conventionally, stiffness can be defined as $S$ where

$$S = -\, \mathrm{d}W/\mathrm{d}h \cdot \tag{7.81}$$

It is here regarded as a purely elastic mode. This can be defined in terms of Hertzian contact theory and Greenwood and Williamson's surface model ($W$ is the load and $h$ the compliance).

Hence, if it is assumed that two rough surfaces can be approximated to one rough surface on a plane, the following relationships can be used:

$$W = C \int_t^\infty (s - t)^{3/2} \varphi(s)\mathrm{d}s \tag{7.82}$$

where $t = h/\sigma_s$ and $\varphi(s)\mathrm{d}s$ is the probability of finding an asperity tip at the dimensionless height $s = z/\sigma s$ from the mean plane of the surface. By differentiating and reorganizing slightly the stiffness becomes

$$S = \left(\frac{3}{2\sigma_s}\right)\frac{F_{1/2}(t)}{F_{3/2}(t)} W \tag{7.83}$$

where, using Thomas and Sayles's terminology where $s$ is a dummy height variable and

$$F_n = \int_t^\infty (s - t)^n \varphi(s)\mathrm{d}s. \tag{7.84}$$

In equation (7.83) the ratio of $F_{1/2}(t)$ to $F_{3/2}(t)$ changes very slowly with respect to mean plane separation and mean plane separation changes little with load; it being assumed that the movement is small and centred on a much higher preload value.

Hence, from equation (7.83)

$$S \propto W \tag{7.85}$$



**Figure 7.34** Stiffness as a function of load and $\sigma$.

This implies that within a practical range around the preload value the stiffness is proportional to the load but not independent of it. (See the first part of graphs in figure 7.34).

Examination of equation (7.84) shows that this formula is the same as that for the $n$th central moment of the distribution. It has already been shown from the work of Longuet–Higgins and Nayak that the moments of the power spectrum are themselves conveniently related to the curvatures of signals, density of peaks, etc (chapter 2). Use of the moment relationship can enable some estimation of the stiffness as a function of surface parameters to be obtained.

Although tractable there are a number of assumptions which have to be made to get practical results. Thus, let the spectrum of the surface be $P(\omega) = K/(a^2 + \omega^2)$, where $K$ is constant and $a$ is time constant. This assumes that the autocorrelation function is exponential which, as has already been shown, is rather questionable over a very wide range but likely over a limited range. If this frequency range is $\omega_u$ which is where plastic flow occurs, and letting $\Omega = \omega_u/a$, then

$$D = \left(\frac{1}{6\pi\sqrt{3}}\right)(m_4/m_2)^{1/2} = \frac{a^2}{6\pi\sqrt{3}}\left(\frac{\Omega^3/3 - \Omega + \tan^{-1}\Omega}{\Omega - \tan\Omega}\right) \tag{7.86}$$

and

$$R = \left(\frac{1}{3a^2\sigma}\right)(\Omega^3/3 - \Omega + \tan^{-1}\Omega)^{-1/2} \tag{7.87}$$

where $m_4$ and $m_3$ are the moments of the spectrum which correspond to the even differentials of the autocorrelation function at the origin.

The constant in (7.82) is

$$C = \tfrac{4}{3}DAE^*R^{1/2}\sigma^{3/2} \tag{7.89}$$

where

$$\frac{1}{E^*} = \frac{1 - v_2^1}{E_1} + \frac{1 - v_2^2}{E_2}$$

and $D$ is the density.

Putting equation (7.89) in (7.83) and integrating gives $S$ in terms of $R$, $D$ and $\sigma$; the average curvature, the density and the root mean square value of peaks. A simple way to do this from profile graphs has been explained by Sherif [70] using

$$\sigma = \sqrt{R_{q1}^2 + R_q^2} \quad \text{and} \quad D = \left(\frac{H_{sc_1} + H_{sc_2}}{2}\right)^2 \frac{1}{L^2}$$

(making the wrong assumptions that the area of density is the square of the profile). However, it is close enough for this approximation. He approximates the curvature $1/R$ to be

$$\frac{S_{m1}}{32R_{pi}} + \frac{R_{pi}}{2}. \tag{7.90}$$

Assuming that this model is acceptable, the way in which $R$, $D$ and $\sigma$ affect the stiffness can be readily obtained. The graphs are shown in figures (7.35 – 7.39).

Approaching the stiffness from the Whitehouse and Archard theory [40], Onions and Archard [41] derived the important surface parameters from a knowledge of the RMS of the ordinate heights, $\sigma^1$, and the

correlation length, $\beta^*$. Hence

$$W = \frac{4}{3}DAE^*(2.3\beta^*)\sigma F(h)$$

Where

$$F(h) = \int_h^\infty (s-h)^{3/2} \int_0^\infty \frac{p^*(s,C)}{NC^{1/2}}\,\mathrm{d}C\mathrm{d}s \qquad (7.91)$$



**Figure 7.35** Stiffness as a function of load and *R*.



**Figure 7.36** Stiffness as a function of load and *D*.



**Figure 7.37** Stiffness as a function of load and correlation length.

**Figure 7.38** Stiffness as a function of load and σ.



**Figure 7.39** Stiffness as a function of load for *GW* and *OA* models.

where $p^*(s, C)$ is the probability density of an asperity having a height and curvature $C$ given by

$$p^*(s,C) \simeq \frac{1}{2^{3/2}\pi} \exp(-s^2/2)\exp[-(s-C/2)^2]\text{erf}(C/2) \cdot \qquad (7.92)$$

For independence $N$ $1/3$ for profiles and $1/5$ for areal mapping, given the simplest discrete patterns. Hence

$$R = \frac{2}{\pi^{1/2}}\left(\frac{2.3\varphi^*}{9\sigma}\right)^2 \qquad (7.93)$$

and

$$D = \frac{1}{5}\frac{1}{(2.3\beta^*)^2}$$

giving

$$S = -\frac{1}{\sigma}\frac{\mathrm{d}W}{\mathrm{d}h} = \frac{6AE^*}{5(2.3\beta^*)}\int_h^\infty (s-h)^{1/2}\int_0^\infty \frac{p^*(s,C)}{C^{1/2}}\mathrm{d}C\mathrm{d}s \cdot \qquad (7.94)$$

It can be seen that σ, $R$ and $D$ affect the stiffness in Greenwood's model and $\alpha$ and $\beta^*$ in Onions and Archard's. Furthermore, the stiffness for a given load is about 50% higher for the latter than the former presumably because Greenwood's model underestimates the contact pressure because it uses asperities of

constant curvature. The simple conclusion to be reached from this exercise is that it is possible to increase the normal stiffness at a joint by either having the fine surfaces or having a large preload for $W$. The former is expensive; the latter imposes a high state of stress on the components making up the structure. A word of caution here concerns the independence or otherwise of the parameters. Although $\beta^*$ and $\sigma$ are given for the Onions and Archard theory [26], in fact the peak curvatures and densities are obtained from it using the discrete properties of the random surface explained earlier. Also, in the Greenwood model $\sigma$, $R$ and $D$ are not independent. It has already been shown by Whitehouse that the product of these is about 0.05 for any random surface:

$$\sigma RD = 0.05 \cdot \tag{7.95}$$

This gives the important basic result that normal stiffness relies on *two* independent geometrical parameters in this type of treatment. These two parameters have to include one in amplitude and one with horizontal importance. Furthermore it is the instantaneous effective curvature of the two surfaces at a contact which is important, not the peak asperities which rarely touch the other surface directly.

What does happen when two surfaces are bolted together? The result depends on whether tangential or normal relative movement takes place in the contact zone. This determines whether 'static' friction occurs which depends on the relative elastic moduli of the two surfaces. If they are the same, no slip takes place at the contact—the displacement has to be equal on both surfaces. If not there is a component of force tangentially which results in slip and a certain amount of tangential compliance, which fortunately is not very dependent on surface roughness [64]. In the event of tangential movement the effective contact area is changed considerably whereas the pressure distribution is not [71]. There have been attempts to take into account stiffness effects for both elastic and plastic regimes. For example, Nagaraj [71] defines an equivalent stiffness as

$$S_e = m(S_c) + (1 - m)S_p \tag{7.96}$$

where the plastic $S_p$ corresponds to slip material either in the normal or tangential directions when the asperities are deforming plastically. Because plastic and elastic regimes are usually involved at the same time both types of movement must take place. Ironically, the presence of friction inhibits tangential movement which has the same effect therefore as asperity interaction, yet the tangential compliance is to a large extent controlled by the elastic region immediately under the contact zone. The way in which the plastic stiffness term changes the apparent stiffness is shown in figure 7.40.



**Figure 7.40** Deformation as a function of load-plastic and elastic.

This graph from Nagaraj [71] seems to suggest that stiffness is independent of load over a considerable range.

Although normal and not tangential stiffness is under discussion, it should be mentioned that the latter is of fundamental importance in wear. Tangential movement of small magnitude due to mainly elastic compliance always precedes gross sliding and hence wear and friction.

Friction can be used to indicate what is happening when two surfaces are brought together with a normal load. Static friction or non-returnable resistance to the imposition of load has an effective coefficient of friction [54–56].

A typical coefficient of friction curve might look somewhat as shown in figure 7.41. It is split into four distinct regions:

I In this region the surface films of oxides are being slid over. The friction is usually small because the surface films have a lower coefficient of friction than metal-metal.

II As contact pressure is increased surface films are progressively broken down, metal-metal contact occurs and the friction begins to rise.



**Figure 7.41** Different regimes in friction.

III Here substantial metal-metal contact occurs, this time with no interceding films; hence the coefficient of friction is high.

IV Extensive plastic surface deformation occurs as the material fails in compression, eventually resorting to the steady condition under that load when $\mu = 0$ and elastic conditions on the whole prevail and there is no plastic movement.

This sort of graph occurs when the load is applied uniformly in time. Should there be any vibration at the same time it has been noticed that the coefficient of friction drops considerably. This effect is great if the variation in load is tangential rather than normal.

### 7.3.3 Mechanical seals

A subject which is close to contact and stiffness is that of mechanical seals. The parameter of importance is related to the maximum gap which exists between the two surfaces. Treatment of this subject is fragmentary. Much emphasis is put on profiles, from which estimates of gap can be made, but in reality it is the lay of the surface which dominates. Early attempts found that the peak curvature was very important. This was due to the fact that the curvatures determine the elastic compliance of the contacting points which in turn determines the mean gap and hence the leakage. The nearest to providing a suitable model for this were Tsukizoe and Hisakado [72], although Mitchell and Rowe [73,74] provide a simple idea in terms of the ratio of the mean separation of the centre places of both surfaces divided by their mutual surface texture, that is

$$\text{leakage} = d \big/ \sqrt{\sigma_1^2 + \sigma_1^2} \cdot \tag{7.97}$$

A particularly simple way to determine the leakage capability was provided by George [75] who showed that the amount of leakage could be estimated by the mismatch of the power spectra of the metal surface and the polymer ring seal pressed on it (figure 7.42). To measure the polymer it had to be frozen *in situ* so as to get a profilometer reading from it representing the as-sealed state.

**Figure 7.42** Difference between metal and polymer.

### 7.3.4 Adhesion

So far the discussion has been about forces, either elastic or plastic or both, which oppose an imposed load between two surfaces these are not the only type of force which is operative. Another kind of force is due to the atomic nature of the surface. A direct result of the competing forces due to attraction and repulsion between atoms or molecules in the two surfaces is that there should exist a separation between them $z_0$ where these molecular forces are equal; the surfaces are in equilibrium. Distances less than $z_0$ produce a repulsion while separations greater than $z_0$ produce attraction. This is shown in figure 7.43 as the force-separation curve and surface energy for ideal surfaces [64].



**Figure 7.43** Type of force as a function of separation.

In the figure the tensile force—that of adhesion—has to be exerted in order to separate the surfaces. The work done to separate the surfaces is called the surface energy and is the area under the tension curve from $z = z_0$ to infinity. This is $2\lambda$, where $\lambda$ is the surface energy of each of the two new surfaces created. If the surfaces are from dissimilar solids the work done in order to separate the surfaces is $\lambda_1 + \lambda_2 - 2\lambda_{12}$, where $\lambda_1$ and $\lambda_2$ are the surface energies of the two solids respectively and $\lambda_{12}$ is the energy of the interface.

The actual formula for the force between the surfaces $F(z)$, is given by

$$F(z) = Az^{-n} + Bz^{-m} \qquad m > n \tag{7.98}$$

where $A$, $B$, $n$ and $m$ depend on the solids. Sometimes it is argued that this adhesive force does not exist since it has been found to be difficult to observe.

The only real evidence for such a force, apart from the purely theoretical one, was provided experimentally by Johnson *et al* [76] who investigated the contact between a smooth glass sphere and flats of rubber and gelatine and found a good agreement between the theory which predicted a force of $3\pi R\gamma/2$ for a radius $R$ of the sphere and surface energy $\gamma$ out. That is, the pull-off force is given by

$$F_0 = 1.5\pi R\gamma \cdot \tag{7.99}$$

The reason why this is so is a classical example of why surface texture is important functionally. Fuller and Tabor [77] resolved the paradox of the non-existence of the adhesive force by attributing it to surface roughness.

Since the surface energies do not differ greatly for different surfaces and the relationship is apparently independent of any other physical property (and in particular the elastic modulus), it seems incredible that adhesion is not commonly seen. Figure 7.44 shows the pressure distribution for a sphere on a flat surface for the elastic Hertzian case and the case where adhesion is present.



**Figure 7.44** Pressure distribution for sphere on flat, for elastic and adhesion cases.

Fuller and Tabor used the Greenwood and Williamson model of surface roughness to show that, as the two surfaces are separated, the contacts at the lower peaks will be in tension (attraction) and broken while the higher ones are still in compression, with the result that there will be negligible adhesion unless the adhesion index

$$\theta = \frac{E^*\sigma^{3/2}}{R^{1/2}\Delta\gamma} < 10 \tag{7.100}$$

where $R$ is the characteristic radius of curvature of the asperity and $\sigma$ is the standard deviation of the asperities.

To test this, very elastic substances were used. Gelatine and rubber have elastic moduli orders of magnitude less than other solids so that the adhesion index-equation 7.100 may be small for moderately rough surfaces (~$1\mu m$). Unfortunately, as will be seen, this model only seems to work for low-modulus solids.

Before looking at other theories, note that the very important parameters involved are curvature and RMS height. In this application they appear to be the direct relevant parameters. This does not mean to say that some correlation could not be found between other parameters. For surfaces having uniform statistics parameters are often interdependent as has been seen.

It is perhaps relevant to note that this property of adhesion between rough surfaces is not the same property as that of 'wringing', first mentioned by Whitworth in 1840 [78] and Tyndall in 1857 [79]. In wringing a thin film of liquid has to be present between the two surfaces before the very substantial normal attraction can appear. Although not directly adhesion, this wringing property, which does not depend on atmospheric pressure, has been made great use of in stacking precision-made stainless-steel blocks to make convenient

metrology length standards. (This use has been attributed to C E Johansson in Sweden at about 1900.) The current thinking [36] suggests that the force is largely due to the surface tension of the liquid.

Reverting back to adhesion, the reason why it is important is because it influences the friction between surfaces in relative motion. Consequently much attention has been given to its theoretical basis, especially with regard to engineering surfaces, a second theory has been proposed by Derjaguin *et al* (DMT) [80] and Muller *et al* [81] which differs in approach from that of the Johnson *et al* (JKR) [76] model. Whereas the JKR model is based on the assumption that attractive intermolecular surface forces result in elastic deformation of the sphere and thus increase the area of contact beyond the Hertzian prediction and only allow attractive forces inside the contact zone; the DMT model assumes that the surface forces do not change the deformed profile from the expected Hertz theory. It assumes that the attractive forces all lie outside the contact area and are balanced by the compression in the contact region, which it takes as having the Hertzian stress distribution. This alternative idea produces a different pull-off force from the JKR model (equation (7.93)). Instead of the constant 1.5 the constraint becomes 2.0 for the DMT model. Muller *et al* in 1980 [81] eventually reconciled the two by producing a parameter $\gamma$ given by

$$\lambda = \left( \frac{R\gamma^2}{E^2 \varepsilon^2} \right) \tag{7.101}$$

Here $E$ is the elastic modulus and $\varepsilon$ the interatomic spacing. In this revised model for $\lambda < 1$ the DMT model applies and for high values of $\lambda \gg 1$ the JKR theory applies. Extra attempts have been made to refine the theories for the elastic deformation case [9], but effects which include plasticity have also been tried [82, 83, 84]. It has been shown that the surface forces alone could be sufficient to produce local plastic deformation, so it appears that theories exist for both very elastic and totally inelastic materials. What about in between? One idea has been presented by Chang *et al* [85, 86] who retain the basic concept of the DMT model, which is that it is where the gap is formed between the two surfaces (just outside the contact zone) that supplies the attractive force. Attraction and repulsion seemingly cancel out for reasons already suggested within the contact. This 'attractive pressure' $\tilde{p}$ according to the 'Lennard–Jones interaction potential' [87], is given by

$$\tilde{p} = \frac{d\varphi}{dz} = \frac{8}{3} \frac{\Delta\gamma}{\varepsilon} \left[ \left( \frac{\varepsilon}{\zeta} \right)^3 - \left( \frac{\varepsilon}{\zeta} \right)^9 \right] \tag{7.102}$$

where $\varphi$ is the interatomic potential, $\Delta\gamma$ is the change in surface energy specified by $\Delta\gamma = \gamma_1 + \gamma_2 - 2\gamma_{12}$, $\xi$ is the separation of the two surfaces outside the contact and 8 is the intermolecular distance. The adhesive force therefore is given by

$$F_9 = 2\pi \int^{\infty} \tilde{p}(z) r \, dr \tag{7.103}$$

where $r$ is the distance from the centre of the contact and $a$ is the contact radius.

When the sphere and flat are just not in contact equation (7.102) becomes

$$F_8 = \frac{8}{3} \pi R \Delta\gamma \left[ \left( \frac{\varepsilon}{\zeta_0} \right)^2 - \frac{1}{4} \left( \frac{\varepsilon}{\zeta_0} \right)^8 \right] \tag{7.104}$$

where $\zeta_0$ is the closest distance. When there is a point contact, $\zeta_0 = \varepsilon$ and the adhesion at the point contact becomes

$$F_0 = 2\pi R \Delta\gamma \cdot \tag{7.105}$$

Compare this equation with (7.99) for the JKR model which is $1.5\pi R \Delta\gamma$.

Using the basic idea of the Greenwood and Williamson model together with the plasticity index, Chang *et al* were able to get a value for the adhesion index $\theta$ for rough surfaces. They found that if $\theta > 100$ the pull-off force due to adhesion becomes negligible for hard steel as opposed to 10 for rubber. They also showed that the adhesion force is negligible compared with the contact load when the plasticity index $\psi \geq 2.5$ or when the surface energy $\Delta\gamma \geq 0.5$ J m$^{-2}$. One interesting point to emerge was that for smooth, clean surfaces the adhesion can be well over 20% of the contact load and therefore it is not negligible as was often thought (figure 7.45).



**Figure 7.45** Adhesive force as a function of adhesive index.

What is strange and disconcerting about this aspect of tribology is that two attempts to explain a phenomenon, one within the contact zone and the other outside, produced practically identical results. When taken with convincing practical evidence this indicates that there is an element of truth in both. In fact this is probably what happens; there is no doubt that the JKR model works for light loads. The main thing to remember is that the parameters of functional importance of the surface roughness are the radius and the standard deviation and the degree to which these affect the adhesion!

The predictions of adhesion theory are such that what can definitely be said is that adhesion is most important for light loads and, taking this to its conclusion, there must be a finite contact area even for zero load!

Figure 7.46 shows the contact between an elastic sphere and a rigid flat in the presence of surface forces: *(a)* the DMT model; *(b)* the JKR model; *(c)* the change in radius of contact as a function of applied load. The pull-off force required is $z_1 = 1.5\pi R\Delta\gamma$ for the JKR and $z_2 = 2\pi R\Delta\gamma$ for the DMT model.

For ease of picturing what is going on in the contact region the adhesion index given in equation (7.100) can be rewritten as

$$-\frac{E^*\sigma^{3/2}}{R^{1/2}\Delta\gamma} = \frac{\text{Hertzian elastic forces of pushing apart the higher asperities}}{\text{pull together (adhesion) forces of the lower asperities}}. \qquad (7.106)$$

Recent work [88] on the adhesion of elastic spheres confirms the applicability of the pull off force being $2\pi R\gamma$ for the Bradley solution [324] and $1.5\,\pi R\gamma$ for the Johnson result in terms of the Tabor parameter [325] $\left(\frac{\sqrt{R\Delta\gamma}}{E^\bullet}\right)E^{\frac{2}{3}} \equiv \mu_T$. For $\mu_T > 3$ the Johnson formula is true. For $\mu_T < 1$ the Bradley formula works best.

One surprising result is that for $\mu_T > 1$. The load approach curves become S shaped leading to jumps in and out of contact!

**Figure 7.46** Contact between elastic sphere and rigid flat in the presence of surface force.

This phenomenon presumably would lead to energy dissipation and could be interpreted as elastic hysteresis. The practical significance of this interesting result depends on the value of $\mu_T$ for real surface asperities.

### 7.3.5 Thermal conductivity

Heat transfer between two solids in contact can take the form of conduction, radiation and convection, usually all taking place in parallel as shown in figure 7.47.



**Figure 7.47** Transfer of heat between two surfaces.

Conduction can be via the actual points of contact of the two surfaces or through the interstitial gap if a conductor is present. Convection can also be the mechanism of transfer if there is a fluid present. In the case when there is a fluid present it is usual to consider that the area involved is the nominal area of the contact rather than the area related to the size of the spots. This is because the gap does not vary appreciably with load [90]. Heat transfer does, however, have its difficulties in the sense that the macrodistortion often complicates the transfer of heat through the microcontact regions.

Radiation can occur across the gap as a further means of heat transfer but is not usually significant within the temperature range encountered under normal conditions.

The biggest problem in heat transfer is concerned with the conductance through the points of contact and this is why the surface texture of the surfaces is so important. Of principal importance is the number of contacts and their size. These two factors dominate thermal conductivity, although the quality of the contact also has to be taken into account. Following Sayles and Thomas [91], the conductance of an isotropic surface in elastic contact with a smooth plane can be evaluated using expressions already given for the real area of contact, so

$$A_r/A = \varphi(t)/2 \tag{7.107}$$

and the force per unit nominal area is

$$F/A = [E^* m_2/(2\pi)^{1/2}]\varphi(t)/t \tag{7.108}$$

where $t$ is the normalized separation of mean planes to the roughness $m_0^{1/2}$ using Nayak's terminology—valid here for $t > 2$, the highest spots—$E^*$ is the composite modulus and $\varphi(t)$ is the Gaussian probability density function.

It is obvious that the total conductance of the surface is intimately related to the number and nature of the contact points. Following for the present the argument using elastic deformation.

$$n = (2\pi)^{-3/2}(m_2/m_0)\varphi(t)\cdot \tag{7.109}$$

The total conductance under this regime is therefore

$$C_0 = A\sum_{i=1}^{n} 2a_i k \tag{7.110}$$

where $k$ is the harmonic mean thermal conductivity or

$$C = 2kn\bar{a}A \tag{7.111}$$

where $\bar{a}$ is the mean contact radius. However,

$$A_i = A\sum_{\alpha=i}^{n} \pi a_1^2 \cdot \tag{7.112}$$

Thomas and Probert [72] here consider an upper limit solution because of the lack of precise knowledge of the distribution of areas. This results in the inequality

$$\bar{a} < (2\pi)^{1/4}(m_0/m_2)^{1/2}t^{-1/2} \tag{7.113}$$

from which

$$C < (2\pi^5)^{-1/4}kA(m_2/m_0)^{1/2}t^{1/2}\varphi(t) \tag{7.114}$$

which is an estimate of conductance given elastic conditions.

As $t$ cannot be isolated from the above equations it is possible to plot the parameters in non-dimensional form (figure 7.48)

$$\begin{aligned}
C_0^* &= C_0(m_0/m_2)^{1/2}(2\pi^5)^{1/4}/kA \\
F^* &= F(2\pi)^{1/2}/E^* m_2 A \\
n^* &= n(m_0/m_2)(2\pi)^{3/2} \\
\bar{a}^* &= \bar{a}(m_2/m_0)/(2\pi)^{1/4} \cdot
\end{aligned} \tag{7.115}$$

Equation (7.114) according to Sayles and Thomas [91] is the upper bound conductance for an isotropic surface on a smooth one. This brings out a few points. The first is that the conductance is nearly proportional to the load, which is also true of the number of contacts in accordance with the general elastic theory and backed by Tsukizoe and Hisakado [92], Thomas and Probert [93] and Kraghelsky and Demkin [94].

Relative to the load, the relationships for $C_0$, $n$ and $a$ can be obtained [91]:

$$\begin{aligned}
d(\ln C)/d(\ln F) &= (1 + 2t^2)/2(1 + t^2) \\
d(\ln n)/d(\ln F) &= t^2/(1 + t^2) \\
d(\ln \bar{a})/d(\ln F) &= 1/2(1 + t^2) \cdot
\end{aligned} \tag{7.116}$$

For values of $t > 2$ where this treatment is valid

$$C \propto F^{0.9} \qquad n \propto F^{0.8} \qquad \bar{a} \propto F^{0.1} \tag{7.117}$$

**Figure 7.48** Estimate of conductance given elastic conditions.

As $t \rightarrow \infty$ the exponents of $F$ for $C$ and $n \rightarrow 1$ and for $\bar{a} \rightarrow 0$ which, judging from equation (7.113), makes sense physically.

These results for the elastic contacts are obtained in terms of bulk parameters of the material and identifiable surface parameters $a$ and $n$. Also the average spot size is almost independent of the load. This means that $C_0 \propto n$. For plastic conditions the conductance is different (from the work of Yovanovich and co-workers [95, 96]):

$$C = 1.5(A_r/A)^{0.95} \tag{7.118}$$

reported by Babus' Haq and co-workers [76, 77].

It is also fair to say that contact conductance has not yet been ideally predicted because of the messy problems that go with it and result directly from the passage of heat. Similar effects do not surround electrical conductance to the same extent. One such problem which is often found is that as the interface and structure heat up there is usually an associated thermal distortion if two dissimilar metals are used in the interface. This can obviously affect the pressures and even the nominal contact area. The direction of heat flow is also another variable.

The thermal distortion problem is shown in figure 7.49. The radius of curvature of the two surfaces as shown is $1/p = q\alpha/k$ where $1/p$ is the radius. Thus for the two surfaces the difference in curvature is

$$q(\alpha_1/k_1 - \alpha_2/k_2) \tag{7.119}$$



**Figure 7.49** Effect of distortion on heat flow in two solids.

This factor (7.119) indicates whether or not the macroscopic contact region so formed is a disc or an annulus. When the heat flows from material 1 to material 2 and $\alpha_1/k_1 < \alpha_2/k_2$ a disc contact ensues as seen in figure 7.49. If the heat flows across the interface in the *reverse* direction then an annular contact occurs.

In either case the heat is constrained to flow via the microscopic asperity bridges formed within the real contact zone. The overall thermal effect is therefore due to a combination of the effects of the macroscopic and the microscopic constrictions—the former corresponding to what is in effect a form error dependent on temperature.

Waviness, here meaning the long wavelength component of the spectrum, does have an effect in that it can determine the number of high spots which come into contact. It is sometimes said that the conductance of waviness and roughness is additive [15]. It is true that the waviness and roughness are both important but it is more likely that the waviness modulates the roughness (i.e. there is a multiplication rather than an addition when simple contact is being considered).

Waviness has a slightly different effect for thermal conductivity as for electrical contact because of the possibility of heat convection or radiation across the gap between the two contacting surfaces. Waviness determines the property of the gap as well as the number of local contacts. The gap as such is not so critical in electrical conductivity [97-100].

In plastic terms the ratio of real to apparent contact is given simply by the ratio of the applied pressure to the yield flow. The relative contact pressure [102] also influences somewhat the effective thickness of the layers of gas in the interstitial voids, which can affect the rate of conduction. Usually the rate of radiation exchange is ignored [101]. It is small compared with the heat transferred by conduction.

The microcontact conductance is proportional to $\bar{a}$ and $n$. Contacts involving surfaces having lower asperity flank slopes in general present a higher resistance to interfacial heat flows. This is attributed to the fact that geometrically this would mean fewer true contact spots, which in effect acts as a sort of constriction when compared with high macroslopes, that is

$$C = (\dot{Q}\sigma)(\Delta\tau Amk) \tag{7.120}$$

where $\Delta\tau$ is the temperature difference between the surfaces, $\sigma$ is the joint surface roughness and $\dot{Q}$ is the power dissipation.

It seems from both these approaches that the number of contacts is very important for light loading (elastic) and heavier loading (plastic).

If there is a fluid then it may transmit heat, but it is true to say that natural convection is suppressed and heat transfer occurs through the fluid by conduction if the gap is small (that is less than about 0.1 mm).

The effective conduction across the fluid is given by

$$C_f = k_f A/\delta \tag{7.121}$$

where $C_f$ is the fluid conductance [101] and $\delta$ is the mean interfacial gap produced by thermal distortion, although the conductance through the fluid (which is usually comparatively small) may become significant when the joint distorts and could under extreme conditions be comparable with the macroscopic conductance.

The total thermal conductance is computed as the reciprocal of the harmonic sum of the macro- and microconductances for each of the two contacting members. If a fluid or filler is present the conductance of the medium in the gap due to distortion must be taken into account [102-104].

Using the same nomenclature as for the elastic case the *microscopic* thermal contact resistance is

$$C = 2ank/g(a/b) \tag{7.122}$$

where $a$ is the mean radius of microcontacts; $g(a/b)$ is called a 'constriction alleviation factor' and accounts for the interference between heat flows passing through neighbouring microcontact bridges. The

*macroscopic* thermal constriction resistance $R$ is obtained by equipotential flow theory. If the contact plane is isothermal

$$R = \tan^{-1}(r/c - 1)\pi ck \tag{7.123}$$

where $r$ is the radial position on the interface from the centre and $c$ is the total radius always greater than or equal to $r$ [105].

This $R$ corresponds to the $g(a/b)$ of (7.122) with a small variation, that is

$$g(a/b) = (2/\pi)\tan^{-1}(b/a - 1)\cdot \tag{7.124}$$

Heat conductance is usually encouraged to be as high as possible so as to increase heat dissipation away from sensitive structures such as in propulsion units, space vehicles, [106, 107], nuclear reactors, and so on. Sometimes, however, there has been an attempt to maximize thermal contact in order to improve the thermal insulation of a variety of systems [108, 109].

In general the heat conductance is considered to be the summation of the macro- and microeffects and also within each the contributions due to convection and radiation, that is, the conductance is given by

$$C_0 = C_{\text{solid}} + C_{\text{fluid}} + C \tag{7.125}$$

where the radiation transfer per unit area between two infinite plates is approximately

$$C_{\text{gap}} = \sigma' \frac{\varepsilon_1 \varepsilon_2}{\varepsilon_1 + \varepsilon_2 - \varepsilon_1 \varepsilon_2} \left( \frac{T_1^4 - T_2^4}{T_1 - T_2} \right) \tag{7.126}$$

From this equation it can be seen that the radiation exchange becomes more important the higher the temperature, as in Wien's law where $\varepsilon$ is surface emissivity. This radiation contribution is negligible at room temperature. The transfer across contacts for temperatures less than 1000°C rarely exceeds 2% of the overall conductance.

Equation (7.125) will be considered again with respect to electrical contact.

### 7.3.6  Relationship between electrical and thermal conductivity

The prediction and measurement of electrical contact resistance are easier than for their thermal counterparts. For example, the relaxation times required for the contact to obtain the steady state is much shorter for an electrical contact. The contact system is much easier to insulate electrically than thermally because of the fewer modes of transfer.

Electrical conduction is more tractable than thermal conductivity because the properties of the gap between the surface is not so critical. There is only one mode of transfer, i.e. at contact points. This means that the distribution of contact points spatially is more significant than for thermal conductivity.

However, what happens at a contact point is more complicated than the thermal case because of the effect of thin films on the surface.

Electrical contact can be used as an analogue to thermal contact (with certain reservations) since both are functions of the same mechanical and geometrical parameters. However, because of the presence of surface films considerable departures from the Weideman–Franz–Lorenz law for electrical conductors have been found for contacts between metals according to O'Callaghan and Probert [103]. Therefore the analogy should only be used for perfectly clean contacts in high vacuum and when radiation can be neglected ($<$300K).

In the presence of the restrictions above, Holm's expression for electrical contact at a circular point is given by

$$R_{ce} = p_e/2a \tag{7.127}$$

In this equation PC can be replaced by $1/ks$ to describe the thermal conduction through the single-asperity contact.

**Figure 7.50** Constriction zone due to limited contact size.

Holm [110] considered the flux lines approaching the contact zone at the asperity point to be hyperbolic curves orthogonal to semi-ellipsoidal equipotential surfaces (figure 7.50), and thereby obtained

$$R_s = \frac{1}{2\pi k_s} \int \frac{d\mu}{(\beta_1^2 + \mu)(\beta_2^2 + \mu)\mu} \qquad (7.128)$$

where
$\beta_1$, $\beta_2$ and $\mu$ are respectively the major and minor semi-axes of an elliptical contact region and a parameter defining an equipotential surface in the vicinity of an interfacial fluid. When the contact is circular, equation (7.128) reduces to the simple normalized equation

$$R_s = \frac{1}{2ak_s} \qquad (7.129)$$

where $a$ is the radius of the contact and where, because the resistances of both parts of the contact are in series

$$\frac{2}{k_s} = \frac{1}{k_1} + \frac{1}{k_2}. \qquad (7.130)$$

For $N$ contacts

$$R_s = \frac{1}{2\bar{a}Nk_s} \qquad (7.131)$$

where $a$ is the mean radius.
According to Tsukizoe and Hisakado [71]

$$\frac{N}{A_0} = \frac{\pi m^2 \varphi(t)}{8} \qquad \bar{a} = \frac{2}{\pi mt} \qquad (7.132)$$

where $m$ is the normalized mean slope, $t$ is the normalized separation ($= u/\sigma$) and

$$b = \sqrt{\frac{A_n}{\pi N}} \qquad (7.133)$$

where $b$ is mean heat flow radius making the assumption that the sum of the cross-sectional areas of the flow channels equals the nominal area. Thermal and electrical contact theories are always functions of the load.

Thomas and Probert [87] collated hundreds of results to get the general relationship between the dimensionless conductance $C^*$ and the dimensionless load $W^*$ of

$$\ln C^* = 0.74 \ln W^* + 2.26 \qquad \text{for steel}$$
$$\ln C^* = 0.72 \ln W^* + 0.66 \qquad \text{for aluminium} \cdot \tag{7.134}$$

They arrived at the dimensionless variables

$$C^* = \frac{C}{\sigma k} \qquad\qquad W^* = \frac{W}{\sigma^2 H} \tag{7.135}$$

which differed from Holm's variables by $\sigma = \sqrt{A}$, where $\sigma$ is the RMS surface height and $A$ the nominal area. They derived the two dimensionless variables from the five possible ones using the pi theorem of dimensional analysis, giving $5 - 3 = 2$ dimensionless possibilities.

From equation (7.125)

$$\frac{1}{R_t} = \frac{1}{R_s} + \frac{1}{R_f} + \frac{1}{R_r} \cdot \tag{7.136}$$

This summation can explain detailed heat transfer. It explains the variation of the thermal resistance of a contact with load at a given temperature. It does not explain the variation of resistance with temperature at a given load.

Note that although in electrically insulating surface films free electrons are inhibited, they are thermally conducting.

Electrical resistance at a contact can be measured much more quickly than its thermal counterpart so that the former is often used to give qualitative indications of, for example, the thermal resistance upon the surface texture.

This electrical-thermal analogue applies more or less adequately when the contact is in high vacuum of better than $5 \times 10^{-6}$ Torr. Under these circumstances the simple electrical constriction form with $1/k$ replacing $p$ can be used (equation (7.127)). Hence,

$$C_c = 2ak \qquad \text{for thermal}$$
$$C_c = 2ap \qquad \text{for electrical} \tag{7.137}$$

In fact, because of the relative ease of measurement of the electrical resistance it has often been used as a way of estimating the real area of contact; that is, if there are $N$ asperities in contact of radius $a$ then $A_r = N\pi a^2$ and $R_c = \rho/2aN$, from which

$$N \simeq \frac{\pi H \rho^2}{4 W R_c^2} \qquad a \simeq \frac{2 W R_c}{\pi H \rho}. \tag{7.138}$$

One assumption often made is that the contact is purely metallic. However, metal surfaces are nearly always covered by thin films (oxides or sulphides usually) which affect the conduction. If these films are thin enough or even semiconducting, they permit a transmission of electricity. Then the resultant resistance of the contact could be represented by

$$R_T = R_c + R_F \tag{7.139}$$

for which Holm [110] defines $R_F$ as $\varphi/\pi r a^2 \cdot \varphi$ is film resistance per unit area, so combining equation (7.139) with the Hertz equation for $a$, that is $a = 1.1(WR/E)^{1/3}$, it can be seen that for elastic deformation a plot of $R_T W^{1/3}$ versus $W^{-1/3}$ should result in a straight line, the gradient and intercept for which enabling the constriction and film resistance contributions to be distinguished. Hence for purely metallic contacts the relation simplifies to

$$R_T \propto W^{-1/3} \tag{7.140}$$

or for chemical film

$$R_{\mathrm{T}} \propto \mathrm{W}^{-2/3} \tag{7.141}$$

Hence as a rule of thumb experimental points should have a basic relationship

$$R_{\mathrm{T}} \propto \mathrm{W}^{-q} \tag{7.142}$$

where $q$ for metallic contacts varies from between $1/3$ and $1/2$ and for film contacts between $2/3$ and unity. Electrical resistance measurements suggest that although surface films remain intact under quite high normal loads, they fracture under values of tangential loads which are less than those necessary to produce macrosliding. Often there have been sharp decreases in resistance although damage has not been apparent. The resulting increase of metal-metal contact has been used to explain the greater ease with which metals 'cold weld' when subjected to additional tangential loads. This is obviously of importance in friction and wear.

It has generally been assumed in the analysis of electrical and thermal conductance that there is a uniform distribution of contacts within the contact zone. This may be the case under controlled conditions but in practice, owing to large-scale macroscopic constrictions caused by the non-conformity of the mating surfaces, it may be necessary to add a macroscopic resistance term.

Note that because the thermal and electrical conductance are both intimately dependent on the real area of contact the *variation* found in both is about the same.

If there is any large-scale variation it is usually in the electrical rather than the thermal conductance because of the greater dependence on surface films [112].

Another phenomenon with many possible applications is a rectification effect. Although this is easily understood with electrical contact it is not quite so obvious in the thermal case. Cases have been reported where the contact thermal resistance for the heat flow from steel to aluminium has been found to be five times higher than in the opposite direction. This has been attributed to macroscopic deflections which would be asymmetric [113].

In both electrical and thermal contact there is a noticeable hysteresis effect. The conductance when first imposing pressure on a contact is generally lower than the value at the same pressure when the load is being reduced. It can, according to Howells and Probert [112], be attributed to two mechanisms:

1. The initial deformation of the asperities which is then succeeded by elastic relaxation.
2. The formation of cold welds which maintain a real area of contact larger than that which was made initially for a given load.

Also important in thermal and electrical contacts is temperature [114]. The thermal conductance of a contact measured at liquid nitrogen temperatures is several times lower than at room temperature. The same is seen in electrical conductance only more markedly, the change being of orders of magnitude. It is suggested that this is due to the presence of a contaminant film on the surfaces which breaks more easily at high temperature.

As the temperature is lowered the hardness and the yield pressure of metals increase with a consequent decrease in the real area of contact. This hardness change affects both the thermal conductance and the electrical conductance.

In general terms the electrical and thermal conductance of contacts are influenced greatly by the real area of contact between the surfaces, the degree of cleanliness and, to a lesser extent, the bulk properties of the components.

The simple relationship between conductance and load is

$$R \propto W^{-q} \tag{7.143}$$

$q$ is usually lower as the temperature reduces and is less for rough surfaces than for smooth. Again, $q$ is usually lower for the thermal case than for the electrical contact under the same conditions.

Because in electrical contact the pressure on a film can affect the conductance it is relevant to note that a film on the asperity does not always behave in a straightforward way. There are four basic possibilities:

1. It can be insulating and so must require some form of removal, disruption or fracture before current can pass.
2. The film can conduct with a certain resistivity.
3. It can rectify giving asymmetric conduction behaviour.
4. It can be very thin, in which case conduction can occur by the tunnelling effect.

It has been found [115] that, even when the contact is sliding, the asperity model is still approximately valid and can be used as a basis for explaining the contact behaviour. The presence of current does tend to accelerate the wear at the contact, the positive contact usually wearing far greater than the negative contact. Also, the relationship between wear and load is different if current is present. The current density adds to the load effect, mainly affecting wear by virtue of changing adhesion properties.

### 7.3.7 Summary

In the interaction between two bodies under load the significant advances have been twofold. One is the 'more realistic representation of the geometry by a random field rather than a set of model asperities. The other is the combining of geometric and physical properties to provide a numerical index to functional behaviour under load. What is missing is a further index to take account of surface film properties in functions like conductivity.

## 7.4 Two-body interactions—dynamic effects

### 7.4.1 General

In sections 7.2 and 7.3 contact has been considered. Although not embodied in the formal definition of tribology, contact phenomena are a key to tribological performance. Historically the study of the subject of lubrication and an understanding of its importance preceded that of wear [116] because of its effect on friction. It will be considered after friction and wear in this section because it is less directly influenced by contact where the surface roughness is a critical factor. In the dynamic effects one surface moves relative to the other. The movement considered is usually tangential rather than normal. There are three ways in which this can happen: sliding, rolling or spinning, or a mixture of them all. In this section sliding will be considered first because rolling is just a special case when the slide is zero. This leads directly to the subject of friction. Friction is at the very heart of tribology; wear results from it and lubrication aims to reduce it.

### 7.4.2 Friction

#### 7.4.2.1 Mechanisms—general

An elegant introduction to friction has been given by Archard (1975) in a Royal Institution Lecture [117] who quotes Dowson [118] that the first recorded tribologist is the man in the carving of a statue of Ti at Saqqara pouring lubricant in front of a sledge. This action is dated about 2400 BC, so it can be assumed that a great deal of thought has been given in the past to the reduction of friction. The study of friction initially was no more than the application of empirical knowledge. At a latter stage more general questions were asked and attempts were made to deduce general laws of behaviour. The laws of friction as set down as a result of this rationalization of knowledge are called Amonton's laws, although they may be equally claimed for Coulomb or Leonardo da Vinci.

The laws are straightforward:

1. The tangential force is proportional to the vertical load between the surfaces.
2. The frictional force is independent of the nominal area of the body in contact.

Coulomb thought that friction was involved with surface texture. He postulated that the effort was needed to make the sliding body lift over the hills of the surface being contacted. This required that the asperities interlock. Obviously this idea cannot work because any work done is retrieved when the body slides into the valleys. What was missing was a mechanism for losing energy.

This energy loss in practice appears as heat and has important consequences for the theory of tribology. As is now known the force of friction for solid contact is the force required to shear the actual junctions at points of contacts between the two surfaces. The actual area of contact is proportional to the load. Friction is also proportional to the load because this real area of contact is proportional to the load. This relationship has always been difficult to verify because of the difficulty of measuring the real area of contact.

In the earliest experiments Amonton came to the conclusion that the force required to overcome friction was one-third of the load—for all cases! Even Euler (1750) [95] agreed. This view expounded by great men took a lot of overturning. In fact it was much later, in the early 1900s, that rigorous experiments were made.

Coulomb's contention [120] that the friction was caused by the raising of the asperities of one surface over those of another was believed wholeheartedly, yet no one explained where the dissipation came from. It had to be inelastic behaviour in the bulk of the material or in some way dependent on the shear of the surface layers.

Hardy's [121] and Bowden and Tabor's [90] work gave a more realistic basis to the property. They made the whole approach fundamentally simple to start with by considering what Archard calls the unit event of friction, which is the deformation and rubbing of a single asperity. Because of this simple approach it was relatively easy to do experiments with the result that a mass of data was soon produced and some simple relationships derived, namely

$$\text{normal load } W = P_m A$$
$$\text{friction force } F = SA \tag{7.144}$$
$$\text{coefficient of friction } = W = S/rn$$

where $P_m$ is the mean contact pressure and $S$ is the shear strength of the contact. Also, initially $P_m$ was taken to be the indentation hardness of the softer of the two contacting materials. It was only when Archard developed certain aspects of his multiple contact theory that this assumption was seriously questioned.

Incidentally, Archard extended the so-called adhesion theory of friction to wear, in which similar very simple expressions to equation (7.144) were produced, namely

$$V/L = KA \cdot \tag{7.145}$$

The volume worn per unit distance of sliding will be proportional to the real area of contact. Here $K$ is a wear coefficient which reflects the probability that an asperity clash will produce a removal or movement of metal.

The important parameters here are the density of asperities and the probability that a contact produces a wear particle. Thus again the count of asperities (and associated contacts) is the important parameter, as it is in elastic and thermal contact and in friction, because each contact has an associated shear, and all contacts contribute to the friction. This is fundamentally different from wear in which a further random or probabilistic element enters into the equation. This is the probability of a summit breaking off in contact. The wear is much more of a fatigue problem than is friction. This latter statement refers to the mild wear regime rather than severe wear.

Returning to friction, this is affected by all contacts at the point where the contact occurs and it is the properties of the asperities which initiate the friction [122]. At the contact of the asperities junctions appear. These can stick or adhere, in which case they have to be sheared or broken when relative movement between the surfaces commences. Alternatively the asperities ride over each other as Coulomb suggested. A third possibility is that the asperities of the harder material simply plough through those of the softer

material. There are therefore at least two and possibly three mechanisms associated with friction: adhesion, ploughing and work done, which involves energy dissipation in some way, perhaps in the form of sound waves or compression waves.

Archard's demonstration that multiple asperities could obey Amonton's laws is elegant and simple. It is not meant to be rigorous because some assumptions have to be made which are not completely true. However, his argument predates the concept and usefulness of the fractal idea of surfaces mentioned often in this book. For this reason it is worth recapitulation.

Take an electric contact between a smooth sphere of radius $R_1$ and a smooth flat under load $W$. There is a circular contact area $A_1$ of radius $b$ where

$$A_1 = \pi b^2 = K_1 W^{2/3} \tag{7.146}$$

Using the Hertz theory [123] the load $\delta W$ supported by an annulus of radius $r$ width $dr$ is

$$\delta W = \frac{2}{3} \frac{W^{1/3}}{K_1} \left( 1 - \frac{r^2}{b^2} \right)^{1/2} 2\pi r \, dr. \tag{7.147}$$

If now the sphere is covered with small spherical caps of radius of curvature $R_2 < R_1$ evenly distributed at $m$ per unit area, the load $\delta W$ is now supported by $q$ contacts ($q = m2\pi r \, dr$) and the load $W$ on each is

$$\frac{\delta W}{q} = \frac{2}{3} \frac{W^{1/3}}{mK_2} \left( 1 - \frac{r^2}{b^2} \right)^{1/2}.$$

The area $a_2$ of each contact is

$$a_2 = K_2 W^{2/3} = K_2 \left( \frac{2}{3mK_1} \right)^{2/3} W^{2/9} \left( 1 - \frac{r^2}{b^2} \right)^{1/2} \tag{7.148}$$

where $K_2$ for $R_2$ corresponds to $K_1$ for $R_1$.

The area $a_2$ of each contact is

$$
\begin{aligned}
A_2 &= \int_0^{r-b} q a_2 = k_2 \left( \frac{2}{3mk_1} \right)^{2/3} W^{2/9} \int \left( 1 - \frac{r^2}{b^2} \right)^{1/2} m2\pi r \, dr \\
&= K_2 \left( \frac{2}{3mK_1} \right)^{2/3} W^{2/9} \tfrac{3}{4} K_1 W^{2/3} \\
&= \overline{K}_2 W^{8/9}
\end{aligned} \tag{7.149}
$$

where

$$\overline{K}_2 = \tfrac{3}{4} K_2 (\tfrac{2}{3})^{2/3} (K_1 m)^{1/3}. \tag{7.150}$$

So with one smaller scale of spheres added to the original sphere the power of $A$ against $W$ has changed from $2/3$ to $8/9$. This process can be repeated adding farther sets of smaller and smaller spheres giving the power $26/27$ etc.

Although the models are an imperfect representation of real surfaces, they nevertheless show that as the model approaches that of a real surface the power law approaches unity, which is the condition hitherto reserved only for plastic contact. Then as the complexity of the model increases the number of individual areas becomes more nearly proportional to the load and their size less dependent on it.

An interesting simulation to reconstruct the profile from the diffraction pattern of a nominally periodic workpiece is reported by Mansfield [124].

The exercise was carried out to validate the signal obtained in the back focal plane of a laws collecting coherent light diffracted from diamond turned cylinders. The instrument under test was the Talyfine of Taylor Hobson.

Cylindrical workpieces of copper are used extensively in copying machines. It is important to monitor the quality of the turning during manufacture as this largely determines the fidelity of the reproduction. Hence the use of a diffractometer which is potentially the most suitable technique for in-process surface measurement.

In order to be compatible with the output of other types of surface finish instruments it was deemed necessary to attempt to produce a profile from the diffraction pattern. As intensities only are measured, the task is that of returning phase into the signal in such a way as to enable a convincing periodic profile to be produced. This Mansfield achieved using some ingenious algorithms. He used a set of estimates of the possible spectra and then proceeded to optimize the estimates by means of an arbitrary but reasonable merit function. This was based on minimizing the sum of squared residuals $S$ where $E$ is the known intensity

$$S = \sum_{k=-N/2}^{k=N/2} \left( \left| E_k \right| - \left| E_k^{ref} \right| \right)^2$$

and $|E_R^{ref}|$ the measured simulation.

The conditional minimization based on the partial derivative of each order was used. He used a novel conjugate gradient technique optimizing technique and was able to produce some useful results.

Obviously, from a practical point of view, it is worthless producing the diffraction pattern at the speed of light only to waste time analysing it. Mansfield gets around this potential obstacle by investing heavily in nested FFT routines. So he seems to have matched the output inverse transform as near as possible to the input transform.

The results of a number of such simulations in which various levels of noise were introduced are encouraging enough to test the technique in practice.

The issue is whether it is possible to get back to the profile from the diffraction pattern. It seems that it is possible given some severe constraints.

The model used here is that of a periodic signal added to a random component—to simulate diamond turning. The question is whether such an iterative method could produce a profile when the periodic component is modulated by the random component. Even more difficult is the case for random surfaces as in grinding.

The model also demonstrated the important point that the surface asperities could not be ignored as they could be if the purely plastic argument prevailed. In practice it has been shown that both occur, but often the mode of material behaviour is less important than the presence, number and distribution of the asperities.



**Figure 7.51** Profile and reconstructed profile.

Nevertheless, this being said, in support for the case of the importance of surface geometry it cannot be denied that once relative motion between surfaces is involved the purely geometrical aspects of the surface are somewhat secondary to the material properties and the chemical properties.

The presence of a thin layer of oxide on the surface, for example, plays an important if relatively undefined role in friction (figure 7.52(*b*)). At first sight the contact is metallic (figure 7.52(*a*)).



**Figure 7.52** (*a*) Apparent contact metal to metal; (*b*) oxide layers at surface boundary.

The presence of the oxide layer may be one of the reasons why it is difficult to assert that friction is independent of the load. According to the characterization of the load versus average area of asperity contact it should simply be because the *average* asperity contact does not change much with load. Unfortunately this is where the idea of average behaviour falls down. The fact that the average asperity remains more or less constant does not stop the biggest contact getting bigger. This means that the stress regions around it will also be much larger than average. This increase could certainly tend to break down any oxide film and allow metal-metal contact, a condition which would not happen to the average contact. This oxide breakdown can cause a large and unpredictable increase in the coefficient of friction and as a result the wear rate. The presence or absence of oxide film is very important to the rubbing behaviour of surfaces [125].

There are two basic distance-related mechanisms according to Archard: one is junction growth and the other is called 'weld growth' which is comparable with or greater than the contact width. It is suggested that the protective action of films is of some significance in non-conforming mechanisms running under conditions of near rolling contact. As the slide/roll ratio is increased the distance of sliding experienced by any given asperity as it passes through the contact region will increase and this could be a significant factor in the damage and welding friction.

Quin [126-128] and Childs [129] have studied the effect of oxide films in sliding contact. They find that the wear rate is very dependent on temperature, which points to a correlation between the growth of oxide film and the temperature. Although the study of oxide films and their behaviour is not really surface geometry, the nature of the geometry can decide how the films behave under load. Also the study of the oxide debris produced in a wear experiment can give clues as to how and in what mode the contact is formed. For example, an analysis of debris from iron and steel in friction and wear situations is significantly different from the analysis of debris obtained in static experiments. For instance, there is under certain conditions a much higher proportion of $Fe_3O_4$ and $FeO$ than $Fe_2O_3$ which is the normally expected film material. It has been suggested that the presence of $Fe_3O_4$ coincides with low wear and high protection against metal-metal contact.

In brief the oxides are significant. Quite how they work in practice is not known for certain. Archard said that 'an analytical approach which may provide a theory for the mechanism and magnitude of frictional energy dissipation in the presence of oxide films is sorely needed' [130].

Mechanistic models for friction have been attempted (e.g. [131]). The idea is that contact can be considered to be between two bodies of cantilever-like asperities whose stiffness is less than the *bulk*. This stiffness is in the tangential rather than the normal direction (figures 7.53 and 7.54).

**Figure 7.53** Cantilever model for asperities in friction.



**Figure 7.54** Cantilever spring model for friction.

The asperities take the form of springs shown in figures 7.53 and 7.54. This model considers friction from the cold junction and asperity climbing modes. There are four stages in the development of the sliding characteristic of the model:

1. The stable junction (cold junction), together with the gradual increase in slope, help asperities to carry the totally applied normal load just at the downslip threshold.
2. An increase in the tangential force brings the full normal load from the downslip threshold to the upslip threshold.
3. The cold junctions fail, allowing the upper asperity to slip upwards; this ends when the asperities meet at their tips.
4. Gross sliding occurs.

This model and others like it cannot be expected to mirror the truly practical case but it does enable discussion to take place of the possible mechanisms to produce the various known phases of surface contact behaviour. Such phases are elastic displacement in the tangential direction (quasistatic behaviour) and the microslips along asperities and gross sliding. To be more relevant, areal behaviour encompassing summit distributions and a general random field characterization of these modelled asperities would have to be investigated.

*Note on computer simulation:*

This whole concept of carrying out speculative experiments on a computer is still in its infancy, not so much in the ideas as in the application and, more importantly still, in the verification of results. This point has been somewhat neglected, but in the field of surface metrology and the behaviour of surfaces it is critical. At one time very simple parameters such as $R_a$ were used to give an idea of the surface, whether or not it was likely to be useful in a functional situation. Later on wavelength parameters were added and later still the more elaborate random process analysis. The idea was to get from a control position, as in the control of manufacture, to

a predictive situation in which the functional capability could be guaranteed. In order to do this, rather than do the experiments fully the link between performance and specification is being increasingly filled by the computer simulation of the surfaces in their working environment.

However, the rapid move to more and more simulation should not slow down the practical verification of results. It may be cheaper and quicker to do simulations but this does not necessarily produce more valid results.

Other more recent models of friction have been put forward. An unusual approach has been made by Ganghoffer and Schultz [132]. They start by artificially inserting a third material in between the contacting surfaces. The intermediary has to be thin.

They in effect, exchange the two body situation with a three body which apparently is more tractable and which enables very general laws of friction and contact to be produced. In some respects this type of approach seems more useful for describing the properties of the thin layer between the surface, rather than the contact and friction between the two principal surfaces. They also get into some mathematical problems as the thickness of the layer is taken to the limit of zero.

Rice and Moslehy [133] make up their model of friction with surface asperities at the interface as expected but also with debris in the model. They include an interface layer as in the Ganghoffer approach but it comprises a much more straightforward substance which is debris. They then proceed to treat the two surfaces and the interface as a mechanical system. They then investigate the dynamics of the system and, in particular, parameters such as displacement and velocity at any point. Using this interface, fluctuations in the frictional force and the normal forces are obtained which apparently is the object of the exercise. In effect the coefficient of friction and the normal force are made time dependent and both sides of the interface are regarded as independent.

The asperities and debris are regarded as exciters to the system. What is not clear is whether they found the obvious connection between the frictional deviations and the differential of the gap profile, or did the presence of debris spoil the correlation?

Some examples of the importance of the roughness of surfaces and friction in everyday life have been reported by Thomas. Most are concerned with slipping on floors. That floors and shoes can be ranked by roughness measurement is not disputed but quantitative values are elusive [134].

The most commonly reported studies on friction are concerned with roads [135], o-ring seals and various modes of lubrication [136] [137].

Recent work on friction in engineering seems to have been concentrated in a few areas. These are:

(i) Frictional effects of manufacturing processes
(ii) Temperature effects of friction
(iii) Mechanisms of ploughing
(iv) Stick-slip.

In the hot rolling of aluminium the detail about the surface contact is the most difficult to quantify [138]. It is now possible to show the dependence of the coefficient of friction on the reduction and on the speed. As the reduction is increased the coefficient first drops, reaches a minimum then increases. At the highest speeds it is fairly constant at about $0.2 \pm 0.1$.

The magnitudes and variations have been attributed to the adhesion theory of friction. It is suggested that as there is a mixed regime of lubrication it is entirely possible that bonds form between the asperities of the roll and rolled strip. The level of reduction and the roll speed vary the number of asperities in the contact zone and the time in which they are under pressure.

It is somewhat strange that attempts to correlate the surface finish to friction in steel rolling [139] concludes that simple parameters $R_a$, $S_a$ do not have the expected correlation. This is to some extent understandable. What is more worrying is that some of the 'areal' peak parameters such as $S_{pk}$ were also ineffective.

Attempts to relate the surface texture directly to the frictional characteristics of aluminium sheet and formed steel have been made. Work with aluminium is becoming more important since car bodies and engines are being made with it. There has been an increase in experiments to control friction, drawability, paintability. Consistent results [140] indicate a friction model which depends on plastic deformation and the real area of contact under pressure. Coefficients of friction across the rolling direction were about 40% lower than with the rolling direction. As plastic flow seemed to be the dominant mechanism the actual size and shape parameters of the roughness turned out to be less important than is usually the case when controlling the process.

Some very careful experimentation has been carried out with sheet steel by Sharpelos and Morris [141]. They quite convincingly tie down frictional coefficients to $R_a$ and the average wavelength. In particular they determined that friction (hence formability) became lower as $R_a$ was increased with the proviso that the wavelength were short. They suspected that the trapping of lubricant was helping to keep the surfaces separated. What they did not appear to do was to realize that their results pointed directly to the use of a single slope parameter for control.

*Friction heating and general comment about surface texture*

One key factor in friction is what happens to the energy generated. Obviously, at least in macroengineering, it emerges as heat which is either dissipated or which raises the temperature in the immediate surroundings. The key to the thermal effects is the time scale involve in a dynamic situation where the lateral velocities are significant.

The impact time between contacting asperities is very short (e.g. at 1m/sec). An impact between asperities of $50\mu$m wavelength takes $10^{-4}$ seconds which is very large when compared with about $10^{-5}$ seconds for the apparent time of contact when considering nominal dimensions of bearings for example. The role of texture temporally is the same as it was spatially—the apparent area of contact got replaced with real area of contact based on the texture. It could be that the *real time of contact is similarly dependent on texture.*

The complication of this is that solving heat equations near to the surface should be tackled impulsively rather than by using the Fourier transform.— In fact, to cater for texture, the Laplace transform with its ability to work with Dirac impulses is the best choice. This is because asperity contacts can be best formulated into Dirac trains of impulses sometimes randomly spaced and in parallel for each and every asperity contact. The conduction equation has to be solved.

$$\frac{1}{K}\frac{\partial T}{\partial t} = \nabla^2 T \tag{7.151}$$

the two factors which are dominated by the texture are the initial and boundary conditions. Fortunately, by changing to Laplace transforms the initial time conditions are inherently included in the transformed equation. The left hand part of the equation in Laplace form can include almost infinitesimal initial time to be taken into account. The right hand side allows the small dimensions of the texture to be incorporated into the solution. Higher temperatures than expected are the outcome of typical calculations involving texture for the simple reason that the conductivity is not large enough to allow the heat to dissipate from the restricted size and shape of the asperity.

The implication is that the plastic mode of deformation of asperities is more likely to occur than elastic and therefore the working point of the function map moves to the right.

The temperature at single asperities has been investigated by Yevtuchenko and Ivany K [142] in the presence of a fluid film. They make some assumptions about the asperity height distribution and asperity shape and work out the 'steady state' maximum asperity temperature! For a practical surface it would be interesting for them to employ the Laplace approach.

Even failing to use the impulsive approach exemplified by Laplace very high temperatures are predicted, of the order of 1200°K, which is very near to important phase change temperatures in steel and titanium alloys.

The point is that failure to consider the asperities and their distribution as the initiator of such temperature leads to artificially low temperatures. Also it is necessary to take into account the shear properties of any lubricant between the surfaces [143, 144].

As the 'unit event' of function in this case is closely related to the unit contact, which in turn is influenced by the texture, it is important to consider all possible unit or basic mechanisms especially if friction prediction is required. One such is the ploughing effect of an asperity on a soft surface. This may sound irrelevant but unfortunately it is needed.

The ultimate requirement is to predict the frictional forces that one surface has when rubbing against another and in particular the ploughing component of that force. Even with today's computing power this aim has not been achieved so the problem is being approached step by step. The first of these steps being to consider the action of one asperity.

Azarkhin [145] developed an upper bound for an indentor ploughing a soft surface. In this and other papers he assumes a perfect plastic deformation and partially uses finite element analysis to check some of the results.

Torrance [146] calculates the frictional force of a wedge moving over an elastoplastic solid. Conditions for boundary frictional coefficients have been inferred from the results. He concludes that, for typical engineering surfaces, friction should be relatively insensitive to roughness but that it could rise dramatically with slopes > 5° and < 1°. The way in which boundary friction rises as the surfaces become smoother is heavily dependent on the boundary film as well as the mechanical properties of the surfaces in contact. In fact the presence of the boundary film makes the interpretation of the practical results very difficult because, as the temperature changes, the conditions also change. It is difficult to make truly predictive statements about real surfaces from this type of analysis. However, it is useful in its own right to help understand the nature and extent of possible damage caused by the stylus of tactile instruments when measuring soft surfaces. This aspect is commented on further in chapter 4 on instrumentation.

Analysis of the ploughing or pushing force due to an asperity also forms the basis of some of Kregelskii's classifications of wear, namely the ploughing class II and the chip forming class III. He uses these terms to describe the mechanisms of breakage of friction bonds.

One aspect of friction, namely stick-slip is important in many applications especially those involved in slow and precise movements. Many surface instruments fall into this category. So far it seems that there is no completely satisfactory solution.

One point which seems to be common [147] is that the instantaneous normal separation of the sliding bodies during stick-slip is an essential element in determining stick-slip transitions. The view is that friction should be regarded as a geometric effect which is related to normal motions and contact geometry and hence texture rather than a relationship between forces or stresses. However, use of the function map can bring the forces and the geometry together.

Early papers presented the stick-slip force as discontinuous changes between the static and a dynamic coefficient of friction e.g. [148]. The values chosen were treated as constants for fixed combinations of materials surface roughness and lubrication, or as a function of the time length of the stick and slip phases as a function of the driving velocity. More recently the friction force is taken to be a function of the relative velocity between the rubbing surfaces. One interesting observation [149] has been made that stick-slip can be induced by retardation between the surfaces. The possible influence of the roughness is in the static friction. The static friction force grows during the time of the 'stick' phase. This has been attributed to the increase in the real contact area due to creep in contacting asperities or possibly the slow breaking through of oxide fields on one or both surfaces by the contacting asperities. Another alternative when oil is present is the squeezing out of the oil in the contact region which itself is dependent on local geometries. One relevant conclusion from this work which indirectly involves the surfaces is the importance of the mechanical damping term. Unfortunately no attempt was made to quantify this, but it seems plausible that surface tension between contacting asperities could well be a factor. This point should be investigated because it should become relatively more important when dealing with miniature parts common in microdynamics.

Recent attempts to explain friction using fractal models have been made but because of the pathological behaviour of areas of contact and pressures at small scales has not been resolved.

### 7.4.2.2    Friction—wear, dry conditions

In the section devoted to friction it was virtually impossible to proceed without having a parallel discussion on wear. Contact friction and wear are integral to the moving scenario in the dry condition. It has to be said, however, that the probability of real situations where the surfaces are dry, loaded and sliding is very low. This condition only usually occurs when there has been a failure in lubrication or there is a very special use as in a reference sliding bearing which takes no load. (Such bearings are sometimes used in the design of metrology equipment for measuring straightness etc.)

Heavily loaded surfaces sliding over each other under dry circumstances usually only occur in cases when the lubricant between them has failed.

Perhaps one of the best experiments to determine the surface parameters to predict the damage when heavily loaded surfaces are rubbed together is that by Hirst and Hollander [150]. Although they used elements of random process analysis on profiles and not a complete mapping of the surfaces, they produced some very interesting results. They chose $\sigma$ and $\beta^*$, the RMS and correlation length of the surfaces, as advocated by Whitehouse.

### 7.4.3    Wear general

In the section devoted to friction it was virtually impossible to proceed without having a parallel discussion on wear. Contact friction and wear are integral to the moving scenario in the dry condition. It has to be said, however, that the probability of real situations where the surfaces are dry, loaded and sliding is very low. This condition only usually occurs when there has been a failure in lubrication or there is a very special use as in a reference sliding bearing which takes no load. (Such bearings are sometimes used in the design of metrology equipment for measuring straightness etc.)

Heavily loaded surfaces sliding over each other under dry circumstances usually only occur in cases when the lubricant between them has failed.

### 7.4.3.1    Wear classification

Basically Hirst and Hollander showed that the gross structure as typified by $\beta^*$ was the scale of size which controlled the failure. Their results are well known and shown in figures 7.55 and 7.56.

For the upper boundary between safe and unsafe they suggested that the values of $\beta^*$ (i.e. the long slopes) were important and not the individual values of $\sigma$ and $\beta^*$. The latter suggested that it was the level of stress that was important and determined the point of failure. They changed the plasticity index $(E/H)/(\sigma/\beta^*)$ by a factor of 0.6.

Wear is basically damage caused to the surface directly or indirectly caused by some contact mechanism, usually with components of normal and tangential forces. Ideas differ as to how best to describe it. Is it removal of material, movement of material or simply changes in the property of material?

Because wear is often associated with the failure of lubrication the various types of wear will be addressed in their relevant lubrication regimes.

In general terms wear and friction go together the main difference being that all contacts contribute to friction whereas only a proportion contribute to wear.

There is one wear mechanism due to Suh which he calls the delamination theory of wear [326]. In this situation considerable normal forces act on the surface with attendant subsurface shear forces. The combination causes flaking to occur: the wear debris is in the form of platelets of thickness often decided by the depth

of the maximum Hertzian stresses. This aspect of wear is not dealt with in this book because the very severity of the forces destroys the surface roughness. So delamination theory is one aspect of wear which is not affected by the surface finish.



**Figure 7.55** Wear damage as a function of surface parameters.



**Figure 7.56** Hirst and Hollander's results for severe wear damage.

It was found that the $\psi$ = constant curves ran almost parallel to the boundary between safe and unsafe surfaces. Furthermore the actual position of the boundary corresponded to a value of $\psi = 1.1$ for the highest load, 0.7 for the intermediate load and 0.45 for the lightest load. These values are in good agreement with the Greenwood and Williamson value of $\psi = 0.6$, the limiting value for elastic deformation when the surface asperities have been modelled on a single constant curvature. Notice the very small range of $\psi$ (0.45 to 1.1) corresponding to a 20-fold change in applied load. This result is important because it is a carefully carried out practical experiment which demonstrates that failure is much more likely to be determined by the nature of the surface texture than the load.

In the same experiment estimates of the size of the contact (ball on flat) corresponded to those wavelengths determined by the $\beta^*$ value at failure, namely $1/\beta^*$, thereby linking the function of the surface to the means of characterizing the surface.

The experimental work showed that it was possible to provide a useful characterization of the potential ability of a random surface to withstand wear. This is encouraging but has to be qualified because the experimentation was very carefully carried out and targeted to fit the specific type of model the theory of Whitehouse and Archard proposed. Later work by Poon and Sayles [152] suggests that this qualification is not necessary for a wide range of conditions.

Nevertheless, even the earlier work showed that the critical deformation is of the same order of magnitude as the main structure and in some ways removes an important restriction since it suggests that the relation between the main structure and the smaller-scale structure is not of great significance. The same is not true in other applications such as in optics.

The mechanisms of friction and wear involved in dry unlubricated situations have been classified in a number of ways. It is not within the scope of this book to deal in depth with the various aspects of the subject, only those in which the surface asperities are involved. It is true, however, that because the transfer of energy takes place at the surface and the surfaces are, so to speak, always modulated by the machining process, the asperities and contacts between them play an important role. Burwell [153] has classified wear. This classification is still used so it will be given here for completeness:

(1) adhesion
(2) corrosion
(3) loose abrasive particle
(4) cutting or ploughing
(5) erosion and/or fatigue.

Unfortunately these mechanisms are not independent of each other. More pertinent from the point of surface geometry is the classification by the great Russian tribologist Kragelskii [154] who defines the destruction of the frictional bond in terms of what can occur when surface asperities are in contact. These different types are shown in figure 7.57



**Figure 7.57** Kragelskii's breakdown of friction bonds.

A spherical asperity is assumed to be rubbing the softer surface. What decides the nature of the sliding is $h/R$ where $h$ is the asperity height and $R$ its radius.

It has been demonstrated [154] that if the situation regarding an asperity could simplistically be regarded as that of an indentor being moved then plastic flow occurs when

$$\frac{h}{R} \geqslant C\left(\frac{H}{E}\right)^2 \tag{7.152}$$

where C is a constant. Alternatively, the mean slope when the asperity is flattened elastically is given by

$$\tan\theta \leqslant C\left(\frac{H}{E}\right) \tag{7.153}$$

where $\theta$ is the slope of the asperity—nominally $\sigma\beta^*$. With plastic flow the friction band can be type II or III in figure 7.57. The important distinction is that in type II material is moved or ploughed whereas in type III material is removed as in chip formation with abrasives. If the asperity can move forward and the displaced material can slip round it a track is ploughed as in figure 7.58. If it cannot slip round it, pile-up occurs and is eventually removed. The important parameter here is the limiting friction between the hard asperity or abrasive and the softer material. This is determined by $\mu$, the true molecular adhesion between the surfaces. $\mu$ is therefore important in that it determines the mode of movement of material whereas $h/R$, the other factor, determines the amount of it.



**Figure 7.58** Ploughed track.

In Kregelskii's simple terms, when $h/R$ is small, welding, bonding and scuffing occur. He distinguishes between type IV and V bonding by means of the critical variation showing how $r$, the shear strength, changes with depth. If the shear strength of the surface film is less than that of the substrate the frictional band is type IV, but if the film or surface layer without the film is stronger, the band will be of type V. The implication of this classification is that the removal or destruction of the surface film by the mating asperities is an adhesive phenomenon.

Archard's approach is shown in table 7.4.

**Table 7.4** Classification of mechanisms contributing to friction and wear.

| Mode | Surface effects | Subsurface or bulk effects |
| --- | --- | --- |
| Elastic deformation | Amonton's laws of friction obeyed with deformation multiple contacts, otherwise wear likely to be involved with protective surface films, i.e. mild wear | Elastic hysteresis component of friction; fatigue wear which includes rolling contact fatigue (pitting) and perhaps some smaller-scale mechanism |
| Plastic deformation | $U = S/H$ [90]; adhesive wear can be mild deformation or severe dependent upon protective films; if severe wear surface and subsurface effects interrelated | Ploughing term of friction, abrasive wear and direct removal of material; if no material removal fatigue mechanisms possible |

It should be noted here that the mode of deformation is based on the static contact, rather than upon the sliding mode. In the case of the elastic contact the maximum shear stress is below the surface. Similarly in a static contact with plastic mode the plastic deformation is below the surface. This is also to a large extent true when moderate tangential forces operate.

In this classification a link between friction and wear becomes possible. This is that the subsurface effects can only play a part in wear when these become linked with surface effects.

### 7.4.3.2 Wear prediction and surface profilometry

Although the surface texture as such is not directly involved in the severe wear classification it is usually indirectly involved whenever wear volume is measured. In these cases the actual volume can be assessed by using relocation methods. Furthermore, specific point by point changes in the surface in the case of mild wear can be identified and quantified during the wear run by superimposing the worn profile back onto the original as in figure 7.59. Measuring weight loss is not so sensitive.



**Figure 7.59** Measurement of wear.

Surface parameters such as skewness can be a good measure of mild wear as well as in fatigue wear [155].

Stout *et al* have been very active in surface mapping to demonstrate that areal (which he calls 3D) information is very useful for characterizing surface wear [156 for example]. Torrance also advocates the use of profilometry [339]. Thomas et al [158] have used profilometry with an atomic force microscope to reveal wear in cylinder bores. They indicate that the bigger slopes and larger numbers of peaks and valleys significantly change ideas about the topographic influx.

In mild wear the surface texture itself can influence the wear rate. Chandrasekaran [159] found an almost exact proportionality between the reciprocal of the initial surface roughness to the reciprocal of the wear rate for steel sliding against steel. In this case he found, not unexpectedly, that fine surfaces produced less wear.

Summarizing, although the texture may not often influence severe wear it can certainly affect the assessment of wear damage.

### 7.4.3.3 Wear prediction models

Because of the economic importance of wear—there has been interest in developing wear models. One of the most comprehensive attempts is due to Meng and Ludema [160] who identify literally hundreds of models and variables. Even then one of the most useful due to Holm [161] and Archard [162] is also one of the simplest.

$$W = K \frac{FN}{H}$$

(7.154)

*W* is the wear volume per unit sliding distance, *K* is the dimensionless wear coefficient, *H* is the hardness of the softer material.

Archard refined the above equation to include qualifications that the contact was between local interacting asperities. The real contact is proportional to load the asperities are circular, the deformation for metals is plastic and the system is isothermal. A wear process is a dynamic process and the prediction of it can be seen as an initial value problem: the initial value is sometimes the roughness. In fact Barwell [163] put wear into three types of wear rate.

$$
\left.\begin{aligned}
&\frac{\beta}{\alpha}(1 - \exp(-\alpha t)) \\
&W = \alpha t \\
&W = \beta \exp(+\alpha t)
\end{aligned}\right\}
$$

<div align="right">(7.155)</div>

Here $\beta$ was described as some characteristic of the initial surface (texture?).

Wear has been simulated using the Winkler surface model [164].



**Figure 7.60** Winkler mattress model of elastic deformation.

In this model, called the elastic foundation model, the contacting surface is depicted as shown in figure 7.60. This is not the same model as used earlier for modelling peak correlation under load. In this wear model there is no connection between each spring element so each spring or elastic rod is regarded as independent. On just what basis the independence is determined is not clear. The contacting surface is in effect a spring mattress with no lateral connections. Although this method is simple it seems to agree well with the finite element results, at least in terms of pressure distribution under the asperity. Friction and plastic deformation in the contact were ignored. It seems unlikely that a method which does not take into account lateral subsurface stress could ever be useful in a dynamic situation.

The idea of simplifying the surface by means of discrete independent units in order to predict average properties of two bodies in rubbing contact is not new. Instead of having mechanical properties such as elasticity it is possible to give each independent unit a geometrical property, for example the probability of height (Fig. 7.61).



**Figure 7.61** Independent event model of surface.

In figure 7.61 the probability density of the height movement of the upper member probability density function can be found in terms of the densities chain of the surface $pd_1. \ldots \ldots \ldots pd_n$. The difference here is that the independence distance $\beta*$ is obtained directly from the autocorrelation function of the surface.

Also, correlation between units can be introduced. This approach for predicting normal displacements in terms of lateral movements between bodies seems to be much tidier than Winkler albeit for a different objective.

In fact the above is much nearer to a real situation (geometrically) of an asperity of one (upper) body moving against a completely rough lower body rather than the flat usually inferred in wear modelling.

In most cases of wear modelling it appears that the condensation of the two real bodies in contact into a sphere on a flat follows almost as a matter of course, presumably because the problem is more tractable. It is surprising that the technique adopted above for stylus integration and similar problems has not been adopted in wear thereby making the model more realistic. Some simplification is allowed by making the autocorrelation function exponential as befits abrasive processes and allowing the surface to be elements in first order Markov chain [11].

One analysis for wear on a coated contact uses an out of date asperity model to investigate a relatively modern problem. Chang [296] is able to differentiate between wear mechanisms for the coating and for the substrate. The mechanism for the coating turns out to be microcutting and that for the substrate is adhesion. One conclusion is that thin coatings under light load conditions can effectively stiffen the asperities.

Various models for the surface have been examined. Two body abrasive wear tests have been carried out using a special set of steel inserts of known shape. Each shape represented a different asperity shape—various angles 60°, 120° of ridges and one radius—all of the same height. These were run against chalk! This choice had some advantages in ensuring pure abrasion. The only variables in the tests were the surface topographies. It was found that the wear rates depended heavily on the surface profile but the investigators chose to measure fractal dimensions of the surface rather than measuring the surface slope directly. $R_a$ was measured but showed little correlation with wear rate. It is a pity that fractal parameters are continually being measured instead of common sense parameters. Wear in these cases has a steady state dynamically and is not a growth mechanism suited for fractal description. Surface gradients have been found to be a useful parameter in determining the wear modes of steels under dry sliding conditions [165].

### 7.9.4 Shakedown and surface texture

In an earlier section 7.2.3.4, the conditions for deformation of the surface were considered. It was established that under certain conditions plastic flow occurred under loading body and under different conditions the deformation was elastic and obeyed Hertz's laws. The discussion then was for the static loading of one body on another with no lateral relative movement of the contacting pair. The criterion (e.g. the plasticity index) was applied once. However, there is a similar situation which occurs when there is a cycle or multiple loading in sliding and rolling. In this case the geometrical and the physical properties of the two rubbing bodies change as a function of operating time. This is called 'shakedown' or in a limited sense 'running-in.' It corresponds to the 'settling down' condition in packaging.

Shakedown is a process whereby the geometrical and physical properties of the rubbing surfaces achieve a steady state condition. In other words mismatches between the dynamic properties of the two contacting surfaces are equalized. In this case the elastic shakedown limit has been reached.

In the case of a repeated contact this could mean that there is plastic movement of material at the contact face until the surface curvature could elastically support the load. Also the yield strength of the softer material increases due to strain hardening until it matches the intrinsically harder surfaces.

Because the natural process of shakedown occurs after some time, efforts have been made to bypass it.

The two changes are in geometry and in the physical properties.

Changes in the geometry of the surface can sometimes be achieved by the manufacturing process. One example is in honing.

In figure 7.62 the profile is of a rough honed cylinder liner. After the lining has been 'run-in' the surface looks like (*b*) which has plateau to support the load and valleys to trap oil and debris. In the past the roughness process looked as if it was the full curve in (*a*). By machining, by fine grinding, or honing the steady state profile (*b*) can be achieved without the running-in time thereby reducing time and cost: the 'shakedown' has been bypassed.

**Figure 7.62** Example of multi-process profile.

This method has worked. The problem is that although the geometry is such that elastic contact takes place the underlying physical properties are not much changed from the original piece. This has to be addressed separately.

One way to get the required hardness of the soft substrate is to coat it or harden it by heat treatment.

When the surface has been coated it is the coating which has the required 'shakedown' properties and not the substrate surface below it. Hence one very important consideration is to make the coating thicker than the roughness of the substrate. By coating the softer surface physical properties of the substrate (softer surface) can be relaxed considerably thereby cheapening the process. If it is too thin the overall effect can be detrimental.

The classification of friction and wear can serve many clarifying purposes but it has to be continually upgraded as new knowledge becomes available. In dry friction, wear and contact situations the influence of the surface texture has been linked with other properties of the solid, usually the yield strength, elasticity and occasionally dislocations and other crystallographic properties.

One of the weak aspects of the theory inhibiting advancement is the continued use of asperity models. Usually the asperities are considered independent, but not always. The reason for such a model is for simplicity, not choice, to make the problem tractable. It has been said earlier that surfaces and their geometric properties have to be treated differently from ordinary electrical signals because the function is invariably concerned with the statistical averages of top-down or bottom-up characteristics which occur in parallel mode. These properties are subtly different from those of electrical signals; non-linearities are relative to different datum points etc. The difficulty of dealing analytically with such signals is made more difficult still if tangential movement is involved. It is no longer possible to infer simple composite characteristics of the gap when movement occurs, although it is often allowable when static behaviour is considered. A simple example of this is shown in figure 7.63 which shows the way in which the composite gap between two surfaces gives a completely misleading picture of the direction of forces acting



**Figure 7.63** Gap properties condensed to rough surface on flat.

between the two real independent surfaces. Instead of the *tangential* forces having the common vector, the making of a composite gap dictates that the *reaction* force has the common vector, which is a different situation altogether! Movement in it is the probabilistic considerations of asperity contact which determine wear as opposed to statistical considerations which apply to frictional behaviour. In many other cases the influence of the geometry can be considered much more independently of the physical and material properties. The scattering of waves, for instance, or in lubrication. It is the latter which will be considered next.

However, the subject of lubrication is very large and complex and covers many possible configurations of two surfaces. Wear mechanisms involved in lubricated situations such as pitting will be considered in lubrication.

### 7.4.5   Lubrication

#### 7.4.5.1   General

This is basically a device to reduce the friction between, and the wear of, two surfaces usually of non-conforming geometry and usually moving relative to each other [166]. The lubrication is some liquid, film or gas used to keep the surfaces apart, yet at the same time being able to sustain a load [167]. That the surface geometry has an effect cannot be denied. What is less sure is the quantitative effect it has on the lubrication properties. Also, the surfaces become important when the lubrication has broken down for some reason. Then wear mechanisms come into play such as scuffing and pitting in which contact, wear and lubrication all play a part.

The lubrication regimes are shown in figure 7.64.



**Figure 7.64**  Lubrication regimes.

The issue in these regimes is concerned with the proportion of the load being carried by the surfaces or the fluid film or perhaps a deposited thin film or boundary film. In the case of the hydrostatic fluid or gas bearing, all the load should be supported all the time by the oil film so that surface effects are relatively small and will only be considered in passing. Then the case of hydrodynamic lubrication will be considered. Following this the specialized cases in which material properties again became important, so as a EHD

(elastohydro-dynamic lubrication) and PHD (plastohydro-dynamic lubrication), will be considered, followed finally by boundary lubrication in which the geometry is of primary importance.

### 7.4.5.2 *Hydrodynamic lubrication and surface geometry*

A very good survey of the influence of surface roughness in lubrication has been given by Elrod [168] who goes into some detail about the old and recent ideas of surface roughness and its effect.

Basically the surface geometry can have two effects: either it helps to generate and maintain the fluid film between the surfaces or it tends to reduce or destroy it. Usually it is average parameters such as $R_y$ which are important in generating the fluid film and extremes such as $R_a$ which destroy the film. Many aspects of the surface geometry have to be taken into account, not just the height. Waviness and the longer wavelengths have to be considered separately from the short wavelengths because they can have different effects. Also, for real credibility the 3D case must be considered. Here the effects of directionality and lay can be important. The lay can have incidental effects such as in mechanical seals or it can be intentional, as in the case of spiral groove bearings, where it can act as an Archimedean screw to pressurize the bearing.



**Figure 7.65** Effect of directionality on film thickness.

It is simpler if, as a start, the roughness is considered to be of short wavelength and can be taken in two forms: transverse to the direction of movement and longitudinal to it. It can also be random or periodic.

The effect of surface roughness and waviness on linear and journal bearings has been investigated many times, as recorded by Elrod. What is clear is that such investigations have been somewhat unsatisfactory in the sense that they have been incomplete. What is needed is to find:

1. Whether texture improves or worsens the performance.
2. Which parameters of roughness amplitude are important.
3. Which parameters of roughness wavelength are important.
4. Does it matter if the predominant roughness, if any, is on the stationary or moving workpiece (one-sided roughness) or on both surfaces to an equal extent (two-sided roughness)?
5. What is the effect of 'lay'? Is isotropy or, more to the point, anisotropy important? If so, is there a preferred value?
6. Should the roughness be random or periodic (in the sense that it is deterministic)? Should it be a combination of the two? Alternatively should the roughness be that of a multiprocess in which the skew is non-negative as in plateau honing?
7. What is the influence of directionality (figure 7.65)? This is a parameter which is rarely specified but which is well known practically in bearing manufacture and also in the manufacture of brakes, clutches and seals.
8. How does the odd large scratch or raised hump which is non-typical of the process affect performance (figure 7.66)?

That these issues have not been answered completely satisfactorily is due mainly to the fact that it is extremely difficult to manufacture surfaces having the desired characteristic to prove the point in an experiment. Invariably other features get included which tend to mask the investigation under way. The alternative

**Figure 7.66** Effect of scratch and groove on film thickness.

as always is to put forward a theory and check it by computer simulation, which has been done many times. It often then becomes a case of the theory running away from reality. In what follows some of the most important observations that have been made from experiments and simulations will be described. It is by no means exclusive.

The first thing to remember is that the pressure is built up as a consequence of the lubricant being pushed into a shrinking geometry by the movement of one component relative to the other. The pressure generated must counteract this tendency in order that the oil can escape. This is the starting point for considering roughness.

Clearly, therefore, the more the lubricant flow is hampered by the texture the greater the pressure build up (figure 7.67(a)). In infinitely wide bearings all the escape routes for the lubricant are found in the $x$ direction. From this it seems plausible that longitudinal structures will help the flow since they represent the same basic flow area as a smooth surface [169] (figure 7.67(b)) while the furrows contribute more in easing the flow than the ridges in hampering it. (This is easily verified by considering the nature of the $h^3$ term in the flow equation—equation (7.157) This will be illustrated later.



**Figure 7.67** Effect of surface lay.

The transverse structure will invariably hamper the flow in the $x$ direction because all the lubricant must pass through the narrow gaps caused by roughness ridges. Hence, since all the lubricant is hampered by the transverse structure and only part of the lubricant is hampered by the longitudinal, it is intuitively obvious that the transverse roughness structure has a stronger effect on the $x$ flow than the longitudinal. This argument holds for wide bearings which can have little '$y$' leakage.

In the case of finite width bearings it is to be expected that there will be some flow in the $y$ direction, so a transverse structure may allow new escape routes whereas a longitudinal structure would not. This would lead to the conclusion that if the bearing is very narrow the transverse roughness may start to reduce the load-carrying capacity: the roles of transverse and longitudinal roughness would be reversed. According to Tender [170] this is true. For the consideration of which of the components should have the roughness on them intuitive feelings often help.

If the transverse roughness, say, is situated on the moving element rather than the fixed element then the effective forcing function of the lubricant into the gap is enhanced. Other parts of the flow equation should not be affected. This should alter the pressure curve and hence permit, an increased load-carrying capacity and probably a load centre shift. The effect of this is found to be more pronounced as the film thickness reduces [171]. It is also plausible that if the moving element with transverse roughness is part of a narrow bearing the

enhanced load-carrying capability due to the forcing function outweighs the $y$ leakage. Accordingly a transverse-moving structure seems desirable in many cases. It shifts the pressure centre, increases the load and for narrow bearings it hardly affects the friction while at the same time, owing to the transverse nature of the structure, it must enhance side flow and thereby improve cooling. Directionality will increase the effective restriction, so a positive directionality is likely to increase the load-carrying capacity relative to a negative one. Yet for longitudinal roughness, directionality in the $+y$ or $-y$ direction will have minimal effect. As a general rule improving load carrying capacity increases friction (figure 7.68).



**Figure 7.68** Friction and surface lay.

For clarification the basic condition using the Reynolds equation will be given, as shown in figure 7.69. This will be developed with, first, transverse roughness, then longitudinal roughness and then with a general surface.



**Figure 7.69** Basic condition for the Reynolds equation.

If there is a local pressure $p$, a pressure gradient $\delta p / \delta x$ and a film thickness $h$ and $x$, and the fluid is incompressible, then the flow must be constant. In figure 7.69 it is assumed that the width of the bearing $B$ is great compared with its length so that the flow is substantially along $x$.

From the continuity of flow (see e.g. [112])

$$\frac{\partial}{\partial x}\left[\frac{h^3}{12\eta}\left(-\frac{\partial p}{\partial x}\right)+\frac{h}{2}U_1\right]=0. \tag{7.156}$$

Integrating and letting the constant of integration be $\frac{1}{2}h^*U_1$

$$\frac{\mathrm{d}p}{\mathrm{d}x}=6\eta U_1\left(\frac{1}{h^2}-\frac{h^*}{h^3}\right). \tag{7.157}$$

When this is averaged over a length of film which is short compared with the variation of thickness of the film along the length $L$, it becomes

$$\frac{\mathrm{d}\bar{p}}{\mathrm{d}x}=6\eta U_1\left(E\left[\frac{1}{h^2}\right]-h^*E\left[\frac{1}{h^3}\right]\right) \tag{7.158}$$

where $E$ is the expectation. This equation is the basic and fundamental Reynolds equation giving the pressure distribution in terms of the film thickness as a function of $x$ along the bearing. The derivation of the Reynolds equation (7.158) is for smooth surfaces, not rough ones [166].

Roughness can be considered simply by adding it to $h$ (in the first case in the transverse direction) as shown in figure 7.70. Then

$$h = h_{\text{shape}} + h_{\text{roughness}} = h_s + b_r. \tag{7.159}$$



**Figure 7.70** Effect of rough surface.

From equation (7.158) assuming $h_r \ll h_s$ and expanding the factors in $h$ by the binomial expansion for the first two terms, where the expected value of the roughness is zero and the variance is $\sigma^2$ (i.e. $R^2_q$ in surface parlance).

The effect of surface roughness in the transverse direction can then be estimated by making suitable assumptions about the value of $\sigma$ relative to the average film thickness. In this equation no attempt has been made to consider the wavelength effects of the roughness, only the amplitude. Thus

$$\frac{dp}{dx} = 6\eta U_1 \left( \frac{1}{h^2} - h^* \frac{1}{h^3} \right) + \frac{3\sigma^2}{h^2} 6zU_1 \left( \frac{1}{h^2} - 2h^* \frac{1}{h^3} \right) \tag{7.160}$$

$$\simeq 6\eta U_1 \left( \frac{1}{h^2} - h^* \frac{1}{h^3} \right) \left( 1 + \frac{3\sigma^2}{h^2} \right). \tag{7.161}$$

This is achieved by neglecting the 2 in favour of unity in the second bracketed term on the RHS, which is acceptable because the sign of the next term is negative.

The importance of the roughness is judged by its value relative to $h/\sqrt{3}$.

The first-order equation (7.161) can be solved with its boundary condition $\bar{p} = 0$ at both ends of the plate and the minimum film thickness found. This is the parameter which is most often sought in any application because it is the minimum value, $h_{\min}$ between the nominal shape reference line and the horizontal member reference line. This then determines the true minimum gap of $h_{\min} + \max(h_r)$ which determines the susceptibility of the film to breakdown.

There is also another approximation which is invariably included in any analysis of tilted-pad geometry or indeed in any instance when a long and short wavelength of roughness are added. This is that although the smaller wavelengths are by their nature formed perpendicular to the nominal shape (or perpendicular to the instantaneous slope in the case of waviness), they are invariably considered to be normal to the horizontal, as shown in figure 7.71.

Fortunately, this angular effect is relatively small in practical cases, but it might be interesting to see what happens to even a simple calculation. No one seems to have done this. Software always ensures that the superposition of tilt and roughness is as in figure 7.71(b). This can only be assumed for small slopes. The real

**Figure 7.71** Addition and rotation of roughness on flat.

situation is shown in figure 7.71(a). It has the effect of increasing the forward slopes (which probably increases the gap), that is $E[h_r] = m^2$, where $m$ is the pad tilt rather than zero which is always assumed. This effect is similar to a surface which has directional properties in which

$$\int_0^L (\,|\,\dot{z}\,|_+ - |\,\dot{z}\,|_-)\mathrm{d}x \neq 0 \cdot \qquad (7.162)$$

Christensen [173] made calculations and tests for rough surfaces which were considered rough, random and Gaussian. He found that for $\sigma = 0.3h_{min}$ for a given load and speed (see [174]) the minimum film thickness is raised by 30%.

*(a) Longitudinal roughness (i.e in the 2L direction.)*
Longitudinal marks often do not usually arise as a result of the manufacturing process, as in the case of the transverse roughness case. Instead this is often produced incidentally as a result of wear in the sliding action, perhaps at the start-up of the relative motion.

The surface roughness acts across the slider rather than along it as shown in figure 7.72. In the case of longitudinal marks it would be plausible to assume that there would be no transverse flow of oil or lubricant because on average the pressure differentials would be evened out. Hence $\bar{p} = p$ and

$$E\left[\frac{h^3}{12}\right]\frac{\mathrm{d}p}{\mathrm{d}x} = \frac{1}{12}E[h^3]\frac{\mathrm{d}\bar{p}}{\mathrm{d}x} \qquad (7.163)$$

so the Reynolds equation from equation [7.157] becomes

$$\begin{aligned}\frac{\mathrm{d}\bar{p}}{\mathrm{d}x} &= 6\eta U_1(E[h] - h^*)/E[h^3]\\ &= 6\eta U_1(h_s - h^*)/E[h^3]\end{aligned} \qquad (7.164)$$

and $E[h^3] = E[h_s + h_r]^3 = E[h_s]^3 + 3h_s^2 E[h_r] + 3h_s E[h_r^2] + E[h_r]^3$ If the surface is symmetrical $E[h^3_r] = E[h_r]^3 = 0$ and $E[h_s^3] = h_s^3$. Hence $E[h^3] = h_s^3 + 3\sigma^2 h_s$ and

$$\frac{\mathrm{d}\bar{p}}{\mathrm{d}x} = 6\eta U_1(h_s - h^*)/(h_s^3 + 3\sigma^2 h_s). \qquad (7.165)$$

Thus it is possible to derive two simple expressions to explain the effect of random roughness in the transverse and longitudinal directions when taken separately and neglecting wavelength effects. Therefore, the pressure gradient with transverse roughness obtained from equation (7.161) is

**Figure 7.72** Tilted rough pad-longitudinal texture.

$$\frac{d\overline{p}}{dx} \simeq \frac{6\eta U_1}{h_s^3}(h_s - h^*)\left(1 + \frac{3\sigma^2}{h_s^2}\right) \tag{7.166}$$

and the pressure gradient with longitudinal roughness obtained from equation (7.165) is

$$\frac{d\overline{p}}{dx} = \frac{6\eta U_1}{h_s^3}(h_s - h^*)\Bigg/\left(1 + \frac{3\sigma^2}{h_s^2}\right). \tag{7.167}$$

From (7.167) the corresponding pressure gradients are decreased rather than increased as in (7.166) by the presence of roughness, both by comparable amounts.

The flow rather than the pressure differential in the two situations is given approximately by

$$\text{flow}\left(\begin{array}{c}\text{transverse}\\\text{roughness}\end{array}\right) \propto \frac{1}{h^3 E[1/h_r]^3}$$

$$\text{flow}\left(\begin{array}{c}\text{longitudinal}\\\text{roughness}\end{array}\right) \propto \frac{1}{h^3} E[h_r]^3. \tag{7.168}$$

In equation (7.156) notice the apparent importance of the roughness skew!

Also if, from equations (7.166) and (7.167),

$$\frac{3\sigma^2}{h^2}\frac{6\eta U_1}{h^3} \tag{7.169}$$

is called the pressure effect of roughness $P_r$, then providing it is relatively small, the effect of transverse roughness is

$$\frac{dp}{dx}_{\text{trans}} = \frac{dp}{dx}\bigg|_{\text{smooth}} + P_r \tag{7.170}$$

and the longitudinal roughness is

$$\frac{dp}{dx}_{\text{long}} = \frac{dp}{dx}\bigg|_{\text{smooth}} - P_r \tag{7.171}$$

if $3\sigma^2 \ll h^2$.

It is not surprising therefore that, if both transverse and longitudinal roughness are present, the resultant pressure gradient (by adding equations (7.170)) is often closer to the smooth bearing slide solution than with roughness in only one or the other direction.

These two equations converge to the simple Reynolds equation for a tilted pad when the roughness $\sigma$ in both cases reduces to zero, subject to a small liberty with the factor 2 in equation (7.167).

The general pattern is clear. The surface roughness expressed as an RMS height plays a part with both transverse and longitudinal roughness directions. The shape of the amplitude distribution in the form of the skew value may also be important.

Both cases here are slightly misleading because although the presence of 'U' affects the film thickness, it is the peak parameters which will break the film down and cause any solid-solid contact through the film, which could initiate failure. The skew value would to a small amount reflect this condition. This situation is more likely to be true at the outlet of the gap. Peak height is the critical parameter, not any more refined parameter such as peak curvature. It is the presence or absence of high peaks which can determine the condition for fluid flow, not the various properties of the peaks themselves.

The very high peak causes an effect which is different from that produced by the mere roughness. The treatment given above for average roughness would indicate that to a large extent surfaces having both transverse and longitudinal roughness would have little overall effect on the film thickness. This is so. It is only when there is a distinct 'lay' that the texture has an effect. This represents changes in the statistics of the roughness with direction. Nevertheless, it is still the average roughness which is represented in the Reynolds equation. The very high peak is different in that it simply breaks the film down. Because it is hard to find or predict, the rogue peak effect is difficult to quantify, so in the case of hydrodynamic lubrication it is the tilt of the pad which is most important. This usually means the first term in a Chebyshev or Legendre form for the gap shape. The surface texture is very much of secondary importance and the material properties are hardly relevant at all.

Equations (7.166) and (7.167) give some idea of the quantitative effect of the surface roughness height and its importance relative to the smooth slider solution. For a very practical case where $\sigma \sim h/3$ the pressures can be changed by 30%. Even so, care has to be used in interpreting these formulae involving the Reynolds equation and the presence of extra constraints to the fluid flow. This is because it is somewhat unsatisfactory in the sense that the roughness affects $dp/dx$ which in turn changes $h$. Nevertheless in terms of scales of size the equations are valid and informative.

Material properties become more prominent in the high-pressure cases of EHD and PHD to be considered later. Also, the variation of viscosity needs to be taken into account [175].

The cases of severely anisotropic surfaces have been sketched out above with respect to transverse and longitudinal lay. Isotropic surfaces are another possibility and yet another case involves surfaces which are anisotropic with the predominant key at an angle to the direction of lubricant flow.

One point not brought out so far in this simplified solution of the Reynolds equation is that, whereas transverse surface roughness increases the load-carrying capacity, it also increases the frictional force. However, the increase in the frictional force is less than the increase in the load-carrying capacity, so the effective coefficient of friction is reduced. This produces the remarkable result that the bearing performance in terms of friction improves hydrodynamically with the addition of roughness. One result of this finding is that any tendency to make the surface roughness smoother in order to improve the bearing performance may be counterproductive.

Consider the more general case when the surface is three dimensional. In particular the surface is random. The Reynolds equation has to be suitably modified in two directions simultaneously to cater for total roughness. In general the $x$ and $y$ roughness are not independent as treated above; cross-correlation terms will usually abound.

Because the surface roughness is random, the solution of the Reynolds equation for pressure distribution and flow will be variable. To take full account of the roughness therefore, the solutions have to be averaged somewhere in the calculation.

Investigators therefore have tended to use one of two methods for general roughness investigation. One method basically obtains many solutions to the Reynolds equation with different roughness embodiments

and then averages the solution. The other method averages the roughnesses beforehand in the Reynolds equation and then solves the averaged equation. Both methods are valid.

The first method was used by Tzeng and Saibel [115]. It is only possible when a solution of the Reynolds equation can be obtained in an analytic form. As a result transverse roughness solutions have been possible. The other method has a numerical advantage over the first. After averaging, the Reynolds equation contains only slowly varying quantities and so is amenable to numerical analysis; also the solution is itself a mean value and hence does not need averaging. It has therefore been most widely used [173]. The major step in this method is the separation of the stochastic quantities which appear as products in the Reynolds equation. How this separation is achieved without violating the properties of the roughness has been addressed [177].

Another factor when the roughness is anisotropic is the quantitative importance of the lay and how effective this is in pumping up the lubricant film or destroying it.

Apart from whether or not the roughness is random or periodic there is the question of wavelength (figure 7.73). How does the dominant periodic wavelength or RMS roughness wavelength $\lambda_r$ compare with the dimensions of the film thickness?

It is conventional [168] to call the roughness 'Reynolds roughness' if A > h, in which case the Reynolds equation can be used. On the other hand if $\lambda$ is very small then Elrod [168] and Chen and Sun [178] all think that Stokes' full equation for viscous flow should be used [178]. For these short wavelengths it is suggested that especially in longitudinal flow the curvatures and slopes of the roughness, not just the height, become important. Note here that the wavelength is used as a criterion rather than a functional parameter in its own right. The roughness under these conditions is sometimes referred to as 'Stokes roughness' to distinguish it from the simpler 'Reynolds roughness' (figure 7.76).



**Figure 7.73**   Areal roughness with periodic-like lay.

The wavelength of the roughness in a way acts for lubrication in a similar way to the wavelength of light in optical scatter. Instead of the criterion of behaviour being in terms of the ratio of $\lambda_r$ to $h$ in lubrication

$$\lambda_r \gg h \qquad \text{Reynolds roughness effects}$$
$$\lambda_r < h \qquad \text{Stokes roughness effects}$$

in optics it is

$$\lambda_r > \lambda \qquad \text{geometric optics}$$
$$\lambda_r \ll \lambda \qquad \text{diffraction optics}$$

where $\lambda$ is the wavelength of light.

The consideration of Stokes roughness [179] has been highlighted because the application of the Reynolds equation to the hydrodynamic lubrication of rough surfaces has been challenged [180]. Certain of the terms which are omitted in the derivation of the Reynolds equation may not be negligible in certain limits.

However, the Reynolds equation can usually be taken to be an asymptotically acceptable model to get an idea of quantitative effects.

Dealing with roughness by ensemble averaging the Reynolds solution is very tedious and difficult. In fact, given that the Reynolds equation is

$$\Delta h^3 \nabla p = \partial h / \partial x \tag{7.172}$$

and where $\iint xy\, dx\, dy = \Omega$ is the area, $p$ the non-dimensional film pressure and $h$ the non-dimensional film thickness, then the characteristics

$$W = \iint p\, dx\, dy \tag{7.173}$$

of the slide bearing (smooth) are $W$, $F$ and $Q$ given in equations (7.173) and (7.174) and (7.175) below. The bearing frictional drag on the runner and the inflow to the bearing are given by

$$F = \iint_\Omega \left( \frac{1}{h} + 3h\frac{\partial p}{\partial x} \right) d\Omega \tag{7.174}$$

and

$$Q_{\text{in}} = \int \left( -\frac{h}{2} + \frac{h^3}{2}\frac{\partial p}{\partial x} \right) dy. \tag{7.175}$$

For the averaging over a set $h_a$ of film thicknesses with different attempts at the roughness $n$,

$$\langle p \rangle = \int P_\alpha p(\alpha)\, d\alpha \tag{7.176}$$

<div align="center">ensemble of<br>the film thickness</div>

$$\langle W \rangle = \int W_\alpha p(\alpha)\, d\alpha$$
$$\langle F \rangle = \int F_\alpha p(\alpha)\, d\alpha \tag{7.177}$$
$$\langle Q \rangle = \int Q_\alpha p(\alpha)\, d\alpha$$

if $p(\alpha)$ is the probability density of $\alpha$ corresponding to one set of roughnesses (and one value therefore of film thickness). It has been pointed out, however [179], that it is not straightforward to go from equation (7.173) to (7.175) using the ensemble pressure $\langle P \rangle$ except in the case of the load $W$, where

$$\langle W \rangle = \iint \langle p \rangle d\Omega \quad \text{and} \quad P_\alpha = \sum_{i=0}^{\infty} \varepsilon^i p_i \tag{7.178}$$

because similar results do not hold for $\langle F \rangle$ and $\langle Q \rangle$ for the reason that the expected value of products is not necessarily the same as the product of expected values! In other words, the whole exercise for $F$ and $Q$ has to be carried out time and time again using different roughness values in order to get $\langle F \rangle$ and $\langle Q \rangle$ ~ which is a time-consuming technique.

In fact

$$\langle F \rangle = \iint \left\langle \frac{1}{h} \right\rangle d\Omega + 3\iint \left\langle h\frac{\partial p}{\partial x} \right\rangle d\Omega \tag{7.179}$$

$$\langle Q_{\text{in}} \rangle = \int \left\langle -\frac{h}{2} \right\rangle dy + \int \left\langle \frac{h^3}{2}\frac{\partial p}{\partial x} \right\rangle dy. \tag{7.180}$$

The solution in terms of perturbations is found by expressing equations (7.179) and (7.180) and $\langle W \rangle$ in terms of small changes in the film thickness $h_a$ $h_0 + \varepsilon n_a$ where $h_0$ is the mean film thickness, $\varepsilon$ is a small increment and $n_a$ is the specific roughness, for example.

For the normal perturbation methods where $i = 4$ in equation (7.178), $\langle F \rangle$ and $\langle Q \rangle$ this gives

$$\langle F \rangle = F_0 + \varepsilon^2 \left[ \iint \frac{\langle n_\alpha \rangle^2}{h_0^3} \, d\Omega + 3 \iint \left( \left\langle n_\alpha \frac{\partial p_1}{\partial x} \right\rangle + h_0 \frac{\partial \langle p_2 \rangle}{\partial x} \right) d\Omega \right]$$

(7.181)

$$\langle W \rangle = W_0 + \varepsilon^2 W_2 + O(\varepsilon^3) \qquad \langle p \rangle = p_0 + \varepsilon^2 \langle p_2 \rangle$$

where

$$W_i = \iint \langle P_i \rangle d\Omega$$

$$\langle Q_{in} \rangle = Q_0 + \varepsilon^2 \left( \frac{1}{2} \int h_0^3 \frac{\partial}{\partial x} \langle p_2 \rangle dy + \frac{3}{2} \int h_0^2 \left\langle n_\alpha \frac{\partial p_1}{\partial x} \right\rangle dy + \frac{3}{2} \int \eta_0 \frac{\partial p_0}{\partial x} \langle n_\alpha^2 \rangle dy \right) + O(\varepsilon^3).$$

(7.182)

(Terms in $\varepsilon_l = 0$ as mean values of the perturbations are zero.)

Explicit representations for the expected values and the probability distributions of the bearing characteristics are difficult if not impossible to find even using the simplest of statistics for the roughness. The usual method of working out the effect of roughness is to use perturbation methods or to simulate the surfaces as with Christensen and see what the result is.

Regular perturbation methods and smooth perturbation methods as well as Monte Carlo techniques have been used [120]. Some results of the comparisons are shown in figure 7.74. From the figures with $\sigma^2 = 1/3$ it is clear that there is not much difference between the methods.



**Figure 7.74** Properties of bearing in presence of roughness.

In the last decade most development in the roughness effect on lubrication has been in the breakdown of wavelength effects, usually into more than one scale of size. When it is realized that the solution of a 3D example of two random surfaces separated by a film and with one moving is very difficult to obtain, it is not surprising that many conflicting and contradictory results are being obtained. One thing, however, is certain. This is that the classical Reynolds equations, averaged or not to account for the roughness, are not sufficient to explain rough surface behaviour. As the roughness wavelength becomes small compared with the film thickness, however calculated, the basic assumptions used to derive the Reynolds equation no longer hold.

One usual criterion is that variations across the film are zero, that is $\partial p / \partial y = 0$. This implies, of course, streamlines which are parallel with the slider. This situation cannot happen when the surfaces are rough and sharp as shown in figure 7.74(b). The flow is more turbulent.

One fairly recent attempt to clarify the situation is that by Bayada and Chambat [181] who use a recent advance in the mathematics of small-parameter equations. Although they are not able to apply the theory to random surfaces they do apply the treatment to multiscale roughness of a quasiperiodic nature. This seems to be close enough to produce at least some idea of the general trends. They are in fact able to show how, by changing the wavelength of the roughness relative to the film, the flow conditions fall neatly into three regimes which are given below.

Summarizing their findings, if $\zeta$ is the film thickness, $\varepsilon$ is the roughness wavelength and $\lambda = \zeta / \varepsilon$ then, using their notation, the following apply:

1. *As $\lambda$ = constant, $\zeta \to 0$ and $\varepsilon \to 0$:*
    (a) The limiting pressure is two dimensional
    (b) It satisfies a new generalized Reynolds equation whose coefficients are calculable and depend on the microstructure of the surface and $\lambda$.
    (c) The height of the roughness has no influence on the result. (Obviously the lack of streamlines is caused physically by local changes in slope and curvature.)
2. When $\varepsilon << \zeta$ (the Reynolds equation regime), $\lambda \to 0$ and a very simple limit equation results which they ascribe to restricting flow in the surface valleys
3. When $\zeta << \varepsilon$, $\lambda \to \infty$ in the generalized Reynolds equation. This justifies what was previously obtained by 'averaging' the classical Reynolds equation.

An increasingly large number of researchers are attempting to analyse rough bearings based on the Stokes equations

$$\Delta u = \nabla p \qquad \text{div } u = 0 \tag{7.183}$$

where $u$ is the real velocity field and $p$ is the real pressure. Typical values for the roughness [180] are a roughness wavelength of 12.5 $\mu$m (and the film thickness) and a surface amplitude of 1.25 $\mu$m. What they do in practice is to obtain an approximate solution of the Stokes system using the classical ratio between the nominal gap and the length of the bearing as a small parameter and then complete an expansion of this solution using the roughness height as a new small parameter. This gives a Stokes-type solution. Then Sun and Chen [180] do a calculation on the Reynolds equation.

Finally, the two answers for the pressure distribution are compared. Phan-Thien has also investigated the differences in the results obtained using the Stokes and Reynolds roughnesses [182]. Two points are generally made:

1. If the Stokes system is considered, nearly all results are devoted to small height roughness and very little in the way of general trends emerges.
2. If the Reynolds equation is considered, computational procedures are proposed to obtain an averaged equation regardless of the roughness height. This gives many results but no one is sure of their validity.

As an example Phan-Thien gives the interesting comparison that if $\lambda h < 0.5$, where $\lambda$ can be considered to be the average wavelength of the surface and $h$ the mean film (thickness), the errors are less than 10% for parallel surface roughness if the surface corrugations on the two bearing plates have uncorrelated parameters. If $\lambda h > 1.91$ then the Stokes solution suggests that there is or should be an increase in the load-carrying capacity, whereas for the same data the Reynolds equation predicts a negative load capacity that depends on $\lambda$ (to the order of $h / L$, where $2L$ is the bearing length).

This is clearly not a satisfactory state of affairs, when two conflicting results are being presented. Fortunately, the quantitative nature of the correction for the surfaces is relatively small anyway, that

is the surface effects in general, although important, are second order, at least as far as the geometry is concerned.

Criteria similar to those of Phan-Thien have been found by Bayada and Chambat [181]. They also use the criteria $h \gg \lambda$, $h \geqslant \lambda$, $h < \lambda$, thereby at least getting some sort of agreement between researchers on the usefulness of the spacing parameter. The important point is that the $h/\lambda$ value seems to determine the mode of solution. Also solutions of transverse random roughness for the Stokes condition are mostly intractable. On the borderline between $0.5 < h$ and $\lambda < 1.91$ differences of only about 10% between the Reynolds and Stokes solutions result, which is somewhat reassuring.

The problem with investigators in this subject at the present time is that of falling between two difficult ways of approach. One is via the practical experimentation on real surfaces in typical conditions, the other is a comprehensive mathematical analysis of the situation. Both seem to be very difficult options. New ideas on the mathematical approach continue to appear, for example Phan–Thien [183] produces a smoothed version of the Reynolds equations in three dimensions. He ultimately finds that the new method taken to the limit reduces to Christensen's results of a decade earlier!

It may be that these two dead ends in theory and practice have accounted for some of the stagnation in the past decade in the subject. This has also led to numerous quasisolutions involving partial modelling and computer simulation which tend to make comparison between researchers' results difficult.

The results obtained so far by whatever means show that surface roughness does make a difference to the performance of tilted slider bearings and consequently journal-type bearings. Differences of a few per cent are typical rather than orders of magnitude. It could be argued that if the outcome of the effect of roughness is questionable then a better manufacturing process should be adopted to ease the problem. Probably the most obvious of the effects of roughness which is not difficult to quantify is that of intermediate values of 'lay'. In this transverse and longitudinal roughness are the extremes and are usually taken to be orthogonal or zero correlated if both are present together, Figure 7.75 gives a pictorial view of the hydrodynamic situation.

The foregoing has looked at the effect of roughness and methods of quantifying it. Apart from just getting an idea of how important it is, another possibility arises. This is to determine whether or not the roughness texture can be used as a positive way to improve the performance of the bearing. It is well known, for example, that deep Archimedean spirals in a surface can produce a cross-pressure which is useful in bearings and mechanical seals. In these cases the spiral is deliberately machined in as a part of the design. Is the same effect made possible by means of having an angled lay on the moving surface, the stationary surface or both? T$\phi$nder has considered this [184] and comes to some interesting conclusions.

The main effect is a crosswise shear flow component due to the channelling action of the angled roughness lay. It affects the pressure formation modestly but the flow can be large and used for cooling. It can be stopped off using a herring-bone appearance which can be used to generate extra pressure (see figure 7.75).

T$\phi$nder suggests that the cross-hatched type of bearing could make a small increase in the pressure distribution for less cost simply by changing the direction of the machining at a relatively early stage in the rough finishing. This approach is good. It is attempting to match the process to the function.

A groove angle of $60°$ appears to be the optimum. This is an interesting observation because $60°$ is about the cross-hatch angle adopted for the cylinder liners of cars. The best cross-hatch angle of honing in the automotive industry has been obtained empirically.

It seems that an angled roughness pattern on the stationary element of the bearing is likely to be better than if it is on the moving part. However, if the pattern is put on both, the improvement in pressure distribution is even larger. Also, the pressure centre distribution remains more symmetrical.

Because the surface roughnesses have an effect on bearing performance it seems logical to propose that the surfaces are generated specifically to improve the bearing characteristics. By this it should be possible to design a surface texture configuration which changes the pressure distribution to that required. The use of CNC machine tools allows for more exotic machining methods. A microstructured surface could easily

change the pressure distribution. Perhaps this could be another example of linking manufacture and function? A simple pictorial synopsis is shown in figure 7.75(b).



Figure 7.75 (*a*)Effect of lay in generating pressure; (*b*) hydrodynamic lubrication.

### 7.4.6 Interaction between two surfaces via a film

#### 7.4.6.1 General

In the section on hydrodynamic lubrication, the lubricant separates the two surfaces more or less completely. It is assumed that there is to be no contact and that the pressure distribution and friction (traction) can be worked out from the flow equations based on either the Reynolds or the Stokes equations. The surface roughness has a second-order effect on the flow. This effect is small but nevertheless significant. Also the surfaces are moving relative to each other. As the load is increased the probability of contact becomes possible and probable. Also in these circumstances an element of rolling is introduced in the relative motion between the surfaces. It is important to realize, however, that the relative movement when contact occurs can be regarded as sliding.

Under these circumstances the pressure distribution has to be determined simultaneously from solutions of the Reynolds equations on the one hand and the elasticity equations on the other. This regime is called elastohydrodynamic lubrication (EHD) (or sometimes EHL) and often occurs in gears, cams, etc. Another lubrication regime occurs when the two surfaces are separated by a very thin film whose thickness is usually much smaller than the surface texture. This type of lubrication is often called 'boundary lubrication.' It will be considered after the section on EHD.

#### 7.4.6.2 EHD lubrication and the influence of roughness

This is important because various forms of wear develop as a result of the regime, such as scuffing and pitting.

The basic equations which have to be solved are listed below:

(1) Hydrodynamic equations:
  continuity equation
  equations of motion $\left.\right\}$ Reynolds equation
  shear stress rate of strain relationship
(2) Elasticity equations:
  compatibility relationship
  equilibrium equations $\left.\right\}$ surface displacement equation
  stress-strain relationship
(3) Thermodynamic equations:
  energy equation for the lubricant
  heat conduction equation for the solid
(4) Lubricant properties:
  equation of state
  relationship between viscosity, pressure, temperature and shear rate.

The full mathematical form for the above can be found in reference [120].

*(a) Elastohydrodynamic lubrication*
The general solution involves the following considerations [153]:

(1) the Reynolds equation;
(2) the equation defining lubricant behaviour;
(3) the lubricant film thickness equation incorporating roughness;
(4) the load equilibrium equation.

(1) is usually taken as equation (7.184) where $U$, $V$ are horizontal components of surface velocity and $W$ is the vertical, $p$ is fluid density and $\eta$ viscosity

$$\frac{\partial}{\partial x}\left(\frac{\rho}{\eta}(H_2 - H_1)^3 \frac{\partial p}{\partial x}\right) + \frac{\partial}{\partial y}\left(\frac{\rho}{\eta}(H_2 - H_1)^3 \frac{\partial p}{\partial y}\right)$$

$$-6\frac{\partial}{\partial x}[\rho(U_2 + U_1)(H_2 - H_1)] - 12\rho U_2 \frac{\partial H_2}{\partial x} + 12\rho U_1 \frac{\partial H_1}{\partial x} + 6\frac{\partial}{\partial y}[\rho(V_1 + V_2)(H_2 - H_1)]$$

$$-12\rho V_2 \frac{\partial H_2}{\partial y} + 12\rho V_1 \frac{\partial H_1}{\partial y} + 12\rho(W_2 - W_1) + 12(H_2 - H_1)\frac{\partial \rho}{\partial t}.$$

(7.184)

If $H = H_2 - H_1 = h + \varepsilon_1 + \varepsilon_2$, where $\varepsilon_1 = h_1 - H_1$, $\varepsilon_2 = H_2 - h_2$ and $h = h_2 - h_1$, $V_1 = V_2 = 0$ $W = -U(\partial h/\partial x) + (\partial H/\partial t)$, $(\partial \rho/\partial x) = (\partial \rho/\partial x) = (\partial \rho/\partial y) = 0$, then equation (7.184) becomes

$$\frac{\partial}{\partial x}\left(\frac{H^3}{\eta}\frac{\partial p}{\partial x}\right) + \frac{\partial}{\partial x}\left(\frac{H^3}{\eta}\frac{\partial p}{\partial x}\right) = 6(U_1 + U_2)\frac{\partial H}{\partial x} - 12\left(U_2 \frac{\partial \varepsilon_2}{\partial x} + U_1 \frac{\partial \varepsilon_1}{\partial x}\right).$$

(7.185)

*(b) Pressure dependence density*
This is almost always negligible [186]

*(c) Pressure dependence of viscosity*
Most often the exponential law of Barus for isothermal conditions is used. Thus

$$\eta(p) = \eta_0 \exp(\alpha' p).$$

The reduced pressure (7.186)

$$q = \frac{1}{\alpha'}[1 - \exp(-\alpha' p)]$$

replaces $\rho$ in (7.171) and (7.172). Hence the generalized Reynolds equation becomes

$$\frac{\partial}{\partial x}\left(H^3 \frac{\partial q}{\partial x}\right) + \frac{\partial}{\partial y}\left(H^3 \frac{\partial q}{\partial y}\right) = 6\eta_0(U_1 + U_2)\frac{\partial H}{\partial x} - 12\eta_0\left(U_2 \frac{\partial \varepsilon_2}{\partial x} + U_1 \frac{\partial \varepsilon_1}{\partial x}\right).$$

(7.187)

*(d) Surface roughness introduction*
For the smooth case $\varepsilon_1 = 0$ and $H = h$. Hence

$$\frac{\partial}{\partial x}\left(h^3 \frac{\partial q}{\partial x}\right) + \frac{\partial}{\partial y}\left(h^3 \frac{\partial q}{\partial y}\right) = 6\eta_0(U_1 + U_2)\frac{\partial h}{\partial x}.$$

(7.188)

(*a*) Reynolds scale          (*a*) Stokes scale



Smooth surfaces          Rough surfaces

**Figure 7.76** Different scales of size of roughness in hydrodynamics.

For longitudinal roughness

$$\frac{\partial \varepsilon_1}{\partial x} = 0 \quad \frac{\partial \varepsilon_2}{\partial x} = 0. \tag{7.189}$$

Then the Reynolds equation is the same as that used for the smooth case with $h$ changed to $H$:

$$\frac{\partial}{\partial x}\left( H^3 \frac{\partial q}{\partial x}\right) + \frac{\partial}{\partial y}\left( H^3 \frac{\partial q}{\partial y}\right) = 6\eta_0 (U_1 + U_2)\frac{\partial H}{\partial x} \tag{7.190}$$

with transverse roughness

$$\frac{\partial}{\partial x}\left( H^3 \frac{\partial q}{\partial x}\right) + \frac{\partial}{\partial y}\left( H^3 \frac{\partial q}{\partial y}\right) = 6\eta_0 (U_1 + U_2)\left(\frac{\partial h}{\partial x} + \frac{U_1 - U_2}{U_1 + U_2}\frac{\partial}{\partial x}(\varepsilon_1 - \varepsilon_2)\right). \tag{7.191}$$

*(e) Lubricant film thickness equation*

$$H(x, y) = h_{\mathrm{c}} + H^0(x, y) + \{U_{\mathrm{n}}(x, y) - \mathrm{Pen}\}. \tag{7.192}$$

Here $h_c$ is the separation between the two surfaces at the centre of contact, $H^0$ is the normal non-deformed distance between surfaces, and $U_n$ is the difference between normal displacements between two opposing points of each surface. Pen is the supposed penetration between the two solids, that is

$$H(x, y) = h_{\mathrm{c}} + H^0(x, y) + \frac{1 - v^2}{\pi E}\int_{-y_{\mathrm{c}}}^{y_{\mathrm{c}}}\int_{-x_{\mathrm{s}}}^{x_{\mathrm{s}}}\frac{p(x', y')\mathrm{d}x'\,\mathrm{d}y'}{[(x - x')^2 + (y - y')^2]^{1/2}} - \mathrm{Pen}. \tag{7.193}$$

This equation is a typical example (in this case for point spherical contact). It has to be evaluated, depending on $p(x, y)$, to give the elastic or otherwise mode of solution.

*(f) Equilibrium equation*

The load $F_n$ applied to the contact is transmitted to the lubricant and must be in equilibrium with the lubricant and also in equilibrium with the hydrodynamic pressure distribution. Thus

$$F_{\mathrm{n}} = \int_{-y_{\mathrm{c}}}^{y_{\mathrm{c}}}\int_{-x_{\mathrm{s}}}^{x_{\mathrm{s}}} p(x, y)\mathrm{d}x\,\mathrm{d}y. \tag{7.194}$$

These expressions represent the considerations which have to be taken into account in order to get some solution of the Reynolds equations for rough surfaces.

Elastohydrodynamic lubrication brings together all the concepts discussed earlier. It can be considered as a contact in which a film is trapped between the contacting part of the surfaces and no relative movement is present. This mode brings in rolling bearings. On the other hand if sliding is allowed, modes of wear and friction similar to dry bearings occur. Hence EHD can be considered to be the meeting point in tribological terms of many engineering situations. (The whole development of hydrodynamic lubrication is lucidly given by Pinkus [187] and Dowson [188]). Many of engineering's most important mechanical parts relate to EHD, although it was not always obvious because of the very small films generated [188]. It was the physical evidence of the retention of surface texture marks on gears, even after prolonged usage, that pointed to the possibility of a hydrodynamic film being present, although it was only when the concept of elastic deformation in the contact region and the change of viscosity with pressure were developed by Grubin [167] that realistic film thicknesses were predicted.

It was the advent of the modern computer that led to solutions of the Reynolds equation, in particular by Dowson and Higginson [186]. They broke new ground which still holds true. It is interesting to note that they achieved the first horrendous numerical calculations on EHD by hand.

In fact, most applications in engineering, such as gears, cams, ball bearings, etc, work under high pressure at the point of contact. EHD really resulted from the realization that fluid properties through the film as well as along it were important. Dowson [189] developed a generalized Reynolds equation for the fluid film lubrication of smooth bearings to investigate just these effects, in particular the effects of viscosity. Roughness was later included in the considerations [190, 191]. Unfortunately neither the measurement technology nor the theory of surfaces was as well understood then as it is now, in particular that relating to asperity interaction and asperity persistence under extreme load. The fact that surface texture could still be visible after much usage was used as a vindication that no solid-solid contact occurred and hence a full fluid film still separated the surfaces. It is now known that under conditions of no slip, it is likely that the asperities as well as the fluid film bear some of the load.

In its simplest form the issue is as shown in figure 7.77. In *(a)* no distortion of the surfaces occurs. In *(b)* the general form is elastically squashed according to the Hertz equations. The fluid pressure is sufficient to deform a general shape of radius $R$ but not the small radius of curvature of the asperity peaks. A distinct separation in curvature distribution is tacitly assumed. This is very probably in order considering the different requirements of the design and hence the manufacture. In *(a)* and *(b)* there may be no relative velocity of the surface producing the film (i.e. no slip), only a squeeze effect produced by the moving contact spot across the surface. The film is present nevertheless. In *(c)* the asperities contact in some places through the film and the pressure of contact causes elastic deformation. Increasing the load even further could cause plastic deformation, in which case plastohydrodynamic (PHD) conditions could be said to exist.



**Figure 7.77** Distortion of surface geometry as load increases *(a)* no distortion, *(b)* general form elastically squashed, *(c)* asperity contact and elastic deformation of asperities occur.

When the failure of bearings due to the wear regime called scuffing occurs, it is often postulated that there is a temperature run-away caused by asperity contact through the film. Under these circumstances the elastohydrodynamic lubrication calculation is often accompanied by thermal considerations [192] and sometimes also by the effect of non-Newtonian lubricant behaviour [193].

In EHD the treatment of the surfaces is very complicated. Asperity contact or interaction is not a simple occurrence; the whole of the gap geometry is relevant. It is across the gap that the pressure is borne by the fluid, not simply in the vicinity of incipient contacts. This makes the analysis of film thickness, especially in

the presence of pressure-dependent viscosity and non-Newtonian rheology of the fluid, difficult to carry out [193]. Most probably, what happens is that the instantaneous 3D curvature of the gap is the important parameter. The two opposing surfaces, especially in rolling as in gears, will have a general smoothing of all the geometry of both surfaces elastically through the film as shown in figure 7.78.



Rough surfaces                Smooth surfaces

**Figure 7.78** Asperity persistence under pressure through film.

This model still preserves the asperity persistence and also maintains the critical surface texture values at the inlet, which determines to a large extent the minimum film thickness.

This effect is likely in both transverse and longitudinal roughness, although it is more likely to be longitudinal in gears because this is the usual direction of grinding. A point to note and make use of in analysis is that even if the grinding marks are not Gaussian, the gap between two such surfaces is often more so than either surface individually. Under these conditions all of the surface geometry is important rather than just the asperities. This point is similar to the distinction between electrical conductance and heat conductance. In the latter case radiation across the gap as well as convection can be mechanisms of the transfer.

When the gap is important the very idea of asperity models for use in surface behaviour seems to be questionable. Despite this, asperity models are still being used, although admittedly in line contacts to simulate the gear or follower [194]. Tractable solutions to the smoothed Reynolds equation and the energy equation seem to be numerically possible providing such problems as transients are ignored. From the results of such analysis [194] it can be concluded that the surface roughness does affect the film temperature. It also causes significant changes in film thickness etc. In many papers there has been a noticeable reluctance to quantify the roughness effects. Much of the earlier work by Dyson [195, 196] and colleagues has yet to be surpassed in credibility. Part of his simulations resorted to storing many real tracks taken from discs, storing them in a computer and literally carrying out the running-in on the computer—similar to the numerical plasticity index techniques later used by Webster and Sayles in contact theory [66].

The argument generally adopted is that there will be substantial contact when the EHD system is incapable of generating sufficient pressure to keep the surfaces apart. In order words, that there is a failure is a consequence of the approach of the two surfaces being restricted by the pressure exerted by the elastic deformation of the asperities. In principle therefore, if the force compliance characteristics of the two surfaces and the gap height distribution are known the conditions required for a failure could be predicted.

Point contact EHD has been considered, notably by Zhu and Cheng [197] and de Silva *et* al [198].

Patir and Cheng [199] have made significant strides in attacking the effect of surface roughness by ascribing an additional 'flow factor' to the fluid film to account for the surface roughness. Also, other workers have tried to simplify the calculations by considering only single asperities [200] and then partially extrapolating to wavy surfaces. In the latter case pressure distributions were found which produced pressure oscillations along the film of the same wavelength as the waviness—a reassuring result if no more.

Recent work [201] to verify early analytic results by Christensen [202] shows good agreement for two types of bearing. A sliding bearing and a journal bearing with and without surface roughness. Two cases of roughness were considered, one transversal roughness and the other longitudinal (Fig. 7.79).

**Figure 7.79** Slider bearing (schematic).



**Figure 7.80** Journal bearing.

The results obtained using finite element methods agreed within a few per cent of the journal bearing. The film thickness was made to be the same order as the roughness.

The real problem in the case of EHD and even hydrodynamic lubrication itself is that the analysis is not straightforward. Even the Reynolds equation is very complicated mathematically, being a non-homogeneous partial differential equation. When this has to be solved with elastic equations, energy equations, and even random roughnesses, it is not surprising that a large divergence of results appears. Single-point, line and area contacts have to be treated as well as rolling and sliding in various relationships to each other.

It seems that there is some way to go before a full solution is forthcoming. In most papers the effect of a stationary roughness is considered. This can be visualized by considering a pure slip situation where the rough surface has zero speed while the smooth surface moves (figure 7.81).

The full problem where both surfaces are rough leads to time-dependent calculations, which even today require excessive computing times [203].

That such effort has been expended on this particular area of tribology is not surprising. Many of the economically important engineering examples where loads are high and the transfer of energy is critical involve an EHD problem. Already mentioned are gears, cams and seals.



**Figure 7.81** Usual configuration for considering EHD.

From what is known, the effect of roughness, whilst not yet being adequately understood quantitatively, is still known to be important, especially in mechanisms which give rise to failure. Probably the most important of these are scuffing and pitting wear. One is related to surface contact and resultant thermal effects, the other is related more to the fatigue effects produced by the collisions and interactions between the surfaces.

A main conclusion is that the surface texture is important in faction but tends to be second order [204]. It also tends to determine the average pressure from which the rheology of the oil is determined from the viscosity-pressure curve.

An interesting point is that perturbation theory produces terms in $\varepsilon$, $\varepsilon^2$, etc, where $\varepsilon$ is $\sigma/h$. The first term is non-linear and difficult to deal with, but luckily this turns out to be precisely the case for smooth surfaces, which has already been analysed. The second-order perturbation terms are linear and this is fortuitous because it is in this second-order consideration that the surface roughness can be introduced [205]. (A pictorial breakdown of EHD situations is given in figure 7.82 where the roughness comes into it [205].)

*(g) Hydrodynamic lubrication and gas lubricants*
Many papers have been written on the effects of roughness on incompressible or virtually incompressible lubrication. They have resulted in the modified Reynolds equations or sometimes Stokes equations relating the averages of pressure gradients with the averages of film thickness. However, little has been published on the compressible case where gases or air are involved. In practice it has been assumed that the averages obtained for incompressible lubrication apply to gas bearings at least for small bearing numbers $\Omega$ (ratio of viscous to inertial forces). The problem is important since many gas bearings run on very thin films. It is also very important in connection with computers and magnetic recording equipment. Researchers Tønder [206] and White and Raad [207], in particular have addressed this problem. They have tended to consider sinusoidal roughnesses in one or more direction. It has been pointed out [208] that real simulations required random surfaces on both of the elements. Computer simulation of this case requires that very small grids are

| Condition | Description | Treatment |
|---|---|---|
| $R$ | Line contact<br>Smooth surfaces<br>Rolling | Dowson<br>Higginson<br>Numerical |
| | Line contact<br>Smooth surfaces<br>Rolling + sliding | Reynolds<br>+Hertz |
| | Line contact<br>Rolling<br>Longitudinal roughness | Dyson<br>Reynolds<br>Hertz<br>+energy |
| | Line contact<br>Rolling<br>Transverse | Reynolds<br>Hertz<br>+energy<br>+rheology |
| | Line contact<br>Rolling + sliding<br>Isotropic roughness | |

**Figure 7.82** Elastohydrodynamic lubrication.

used in order to get the very fine detail necessary. A new computational method which allows parallel processing element by element using (EBE) is now being used [209]. This in theory at least should enable simulations to be carried out using very complicated surfaces.

The general consensus is that [210] the fluid film height averages that apply to the incompressible case do not apply to the compressible case even for an infinitely high bearing number $\Omega$ $(6\mu UL/P_a h^2)$ Basically the following conclusions seem to be relevant:

1. With transverse roughness (stationary) the load-carrying capacity peaks at a finite bearing number. This is unlike the smooth bearing case where the load increases monotonically and approaches a limit as the bearing number becomes very large.
2. At very high gas bearing numbers, the load is governed by the limiting form of the Reynolds equation

$$\frac{\partial(pH)}{\partial x} = 0. \tag{7.195}$$

The mean pressure produced from equation (7.195) is of the form

$$\tilde{p}(x) \sim \overline{H^{-1}(x)} \tag{7.196}$$

which is different from that proposed by Tønder, which appears as

$$\tilde{p}(x) \sim \overline{H^{-3}(x))}/\overline{H^{-2}(x))}. \tag{7.197}$$

Although this is dimensionally the same, it is not in practical terms.
3. Surface roughness influences the bearing load over a much wider range of gas bearing number than for the smooth bearing case. This is due to the influence of the short-wavelength scales of roughness on the higher-pressure diffusion terms in the Reynolds equation.
4. Note that with the gas bearing the load-carrying potential developed with the moving rough surface is less than the load developed with the stationary surface bearing.

This last point is now important because it could determine which is most important—the roughness of the head or the roughness of the tape/disc—when considering the maximum pressure generation to aid lift-off when starting up a magnetic or CD device.

It seems on this argument that having roughness on the disc is definitely to be preferred, not only because it establishes the air film sooner but also because it might, in addition, reduce friction on start-up. Obviously wear would be less in this case. It is also relevant that because the film thicknesses in gas bearings can be thinner, the roughness is even more important than for the corresponding fluid film bearing despite the fact that the viscosity variation with pressure is less.

The Reynolds equation can be used, after some modification, to consider the conditions in an ultra-thin film gas bearing with various roughness orientations. The application is in magnetic recording systems where the film thickness has to be small to get maximum recording density. Here thickness is taken to be 0.1 $\mu m$.

The effects of the topography on the behaviour of thin film flow is complicated. Two possible methods already mentioned are averaging the film thickness or averaging the flow. The average flow method is easy to apply. It is the rate of flow passing through the averaged gap.

Although results are available concerning the effect of roughness on flow and load carrying capacity quantitative valves are rare. However, as previously stated for hydrodynamic lubrication longitudinal roughness helps flow and reduces load carrying capacity in EHD when the film thickness is a fraction of the roughness value equation 7.200 and 7.201. It is interesting to note that for EHD situations where many asperity contacts are made that the two surfaces conform to a degree where the peaks on one surface resemble the valleys in the contacting surface. The idea is that the one surface is in effect the negative of the other.

The actual degree of conformity was taken to be the value of the cross correlation of upper surface $S_1$ with the inverted $S_2$.

**Figure 7.83** Conformity of two surfaces using cross correlation.

The cross correlation maximum was used to align the two graphs together. Then the effective composite surface was made from $S_1 - S_2$. The difference replaces either of the original profiles.

Parameters such as the $R_q$ value and curvature of the new composite surface now act as the significant parameters. Values in longitudinal and transverse directions are offered as control parameters for wear.

An interesting observation concerned the isotropy of the surfaces $\Lambda = \lambda_x / \lambda_y$ where $\lambda_x$ and $\lambda_y$ are the independence lengths in orthogonal directions.

It turns out that $\Lambda$ is 3 rather than 0.05 found when using the original surfaces.

Tyagi and Serhuramiah approach the investigation in a sensible way in that they actually see what happens in a real experiment, and draw some conclusions about the best parameters.

The big problem is that the whole idea behind finding functionally useful parameters is the prediction of them from the original surfaces e.g. what is the number of and nature of the contacts for a given load and speed given the unworn surface data? What needs to be done is to develop a functional operator to apply to the original profiles or area data to get the desired output! It is this last step which is the difficulty.

EHD has suffered inexplicable failures which accepted theory does not account for. Cutting off the lubricant or allowing inlet shear can reduce the film or make the film fail but only after a time lapse. Smooth surface theory does not adequately explain some fast failures.

This is an obvious case for rough surface models. The issues as far as roughness is concerned is first, how the minimum film thickness is affected by rough surfaces entering the inlet of the gap and, second, what happens to the film and the pressure distribution when roughness enters the high pressure zone. There is another issue which is whether or not the roughness itself undergoes distortion when it enters the high pressure region. In a series of papers on the effect of roughness and waviness in EHD.

Morales–Espejel and Greenwood [211, 212, 39] and Venner and Lubrecht, [213, 214, 215] have tried to quantify what happens when rough surfaces are involved.

The normal Reynolds Equation for line contact is

$$\frac{\partial}{\partial x}\left(\frac{ph^3 \partial p}{12\eta \partial x}\right) = \bar{u}\frac{\partial(ph)}{\partial x} + \frac{\partial(ph)}{\partial t} \tag{7.198}$$

For the case of where roughness is on the upper and lower surfaces the film thickness is

$$h(x,t) = \xi_{00}(t) + S_2(x,t) - S_1(x,t) \tag{7.199}$$

where $S_1$ is given by

$$S_1(x,t) = \frac{x^2}{4R} + \frac{1}{2}v(x,t) + Z_1(x,u_t)$$

$$S_2(x,t) = \frac{x^2}{4R} + \frac{1}{2}v(x,t) + Z_z(x,u_2 t) \tag{7.200}$$

The elastic displacements

$$v(x,t) = \frac{4}{\pi E'} \int_{-\infty}^{\infty} p(x',t) P_n |x - x'| dx'$$

(7.201)

and the viscosity

$$\eta(p) = \eta_0 \exp^{ap}$$

(7.202)

The equilibrium equation supporting load $F$ per unit length

$$F = \int_{-\infty}^{\infty} p(x,t) dx$$

(7.203)

$$Z(x,t) = d \, Sin\left(2\pi \frac{(x - xd)}{w_1}\right) \text{represents the waviness}$$

(7.204)

i.e. the spatial variation on the surface.

For an individual surface feature

$$Z(x,t) = 10^{-10[(x-x_a)/w_1]^2} . d\cos\left(2\pi \frac{(x - x_d)}{w}\right)$$

(7.205)

$xd$ is the initial point of the waviness at time $t$ and the feature centre at time $t$

In practice the viscosity is so high that

$$\frac{ph^3 \partial p}{12 \zeta}$$

(7.206)

So a simple linear equation results

$$\bar{u} \frac{\partial (ph)}{\partial x} + \frac{(ph)}{\partial t} = 0$$

(7.207)

this is solved for the complementary function, the transient response and the particular integral.

So for relatively simple surface structure on both surfaces the qualitative solution for the transients are mode up of two components, a moving, steady state solution travelling with the velocity of the surface and a complementary function which comes from the entering roughness travelling with the velocity of the lubricant $\bar{u}$. In two sided waviness there are two complementary functions and two moving steady state solutions.

When the surfaces travel with different velocities a new unknown occurs. This is the phase shift in the induced waves produced because the roughness from each surfaces hits the high pressure zone at different times.

The problem with this approach is that the inlet of the contact is not in the analysis so that the amplitudes of the transients is unknown. The way usually adopted is to solve numerically and try a few values until the film conditions become clear. The analytic solution could possibly have been worked out by Laplace transforms where initial conditions are inherent and the inclusion of freak surface conditions is easier.

The main conclusion from recent analytical attempts is that the roughness plays a prominent role. What would be of interest is the nature of the roughness deformation in the pressure zone. Numerical solutions of the modified Reynolds equation seem to be the way forward. It would also be instructive to have random roughness rather than periodic.

### 7.4.7 Boundary lubrication

#### 7.4.7.1 General

There are a number of critical engineering situations where EHD conditions cannot be sustained. Constant-velocity joints, for example, are heavily strained machined parts where this is true. High Hertzian pressure, low relative velocity and high temperature are the conditions where EHD fails. In this regime boundary lubrication conditions hold where the use of materials and surface conditions with a low tendency to adhere (and subsequently weld with the associated high wear) and high resistance to fatigue is required. There is also a requirement for thin-film lubricants which are able to generate reaction layers in any local contact areas to reduce adhesion further [216, 217].

Basically boundary lubrication is the formation and use of very thin durable films on the contacting surfaces. The emphasis is very much on the chemistry of the film. The contacting properties and modes of contact are very similar to those for dry conditions in which there is little relative movement or when there is very high pressure.

The way in which the regime changes from hydrodynamic lubrication to boundary lubrication is well exemplified in a Stribeck diagram shown in figure 7.84, where A is boundary lubrication, B mixed lubrication, C hydrodynamic lubrication and D elastohydrodynamic lubrication.



**Figure 7.84** Lubrication regimes. Stribeck diagram.

These regimes are shown somewhat simply in figure 7.84. Moving as the speeds reduce, or the load increases, or the viscosity is decreased, the surface condition goes from A to D.

The well-known minimum in the coefficient $\mu$ is shown in figure 7.84. In early work boundary lubrication was considered to have been the dominant regime when p was independent of $zV/p$ (i.e. regime A in the figure) but these have been modified since the work of Dowson and Higginson showed the characteristic nip at the rear of the potential contact region (figure 7.84(b)).

In boundary lubrication, because of the extreme pressures the film can be compared with that of an elastic solid rather than a viscous liquid; hence the reference back to dry contact behaviour as far as asperities are concerned.

There has been considerable speculation as to what constitutes boundary lubrication and what constitutes the thin-film region of EHD.

Often the higher shear rates generated in the very thin films raise the temperature and lower the viscosity, thereby counteracting the increase in viscosity due to pressure. Often a temperature-viscosity relationship

**Figure 7.85** Lubrication regime transition: (*a*) lower limit hydrodynamic; *(b)* elastohydrodynamic; *(c)* mixed; (*d*) boundary lubrication.

is assumed to be exponential in the same way as the pressure. The difference is that the variable is in the denominator of the exponent. Thus

$$\eta = \eta_0 \exp(Q/RT) \qquad (7.208)$$

where $Q$ is an activation energy (~ 16 kJ mol$^{-1}$ for hydrocarbons), $R$ is the gas constant and $T$ the absolute temperature. These EHD films are capable of supporting the normal load even when the thickness is less than the surface roughness. Also, relatively high coefficients of friction have been observed for these films even when their thickness is greater than the mean roughness.

So there is still a somewhat nebulous region in which the regimes of boundary lubrication, mixed lubrication and EHD are blurred. What used to be considered boundary lubrication can still be in the province of EHD. The large increase in viscosity due to pressure in EHD with its effective solidification of the film seems to remove the need for a solid-like layer of adsorbed molecules at the surface interfaces. This does not mean that oil technologists do not still add small quantities of boundary lubricants to their products to improve the wear and damage capability. These additions have been shown to help in the performance of cam and tappet systems.

The three issues [218], simply put by Briscoe are as follows:

1. How are the boundary lubricating films formed? Can they be artificially formed?
2. What are their structures and properties (chemical and mechanical) under the sliding conditions which prevail at the point of asperity contact?
3. Under what conditions do they fail?

In principle the film has to be governed by the nature of the solid interface and the lubricant. A point to note is that boundary lubricant layers are formed at 'normal' working conditions (i.e. at temperatures less than 150°C) which distinguishes them from those layers produced by extreme pressure additives. The usual distinction has been that the former are primarily organic, whereas extreme pressure additives are mainly inorganic.

Typical boundary lubricants are stearic acid, which has a classical long linear hydrocarbon chain which acts as the durable 'separator' of the solids, and a carboxylic acid group which acts as the anchor to fix it to the 'metal' surface. Other active functional groups are amines, esters and alcohols. The molecular structures of some simple boundary lubricants are shown in figure 7.86.



**Figure 7.86** *(a)* General lubricant; *(b)* active group.

The classical boundary lubricant, stearic acid, consists of a chain of 17 methylene groups. Usually the chain is linear and may zigzag, but it has no side groups because these reduce the efficiency. However, sometimes one of the $CH_2$-$CH_2$ bonds is replaced by unsaturated bonds, as for example in the oleates.

The active groups have varying degrees of reactivity on the metals but the general picture is as shown in figure 7.87.

In simple boundary lubrication it could be argued that the physical and mechanical properties of these oil film materials are not very important because, at the actual junction, the material formed may have very different properties.

The actual nature of the material surface is most important, the outermost layers, which come into contact with the lubricant controlling this chemistry. For chemically inert metals this may be just adsorbed gas or



**Figure 7.87** Boundary film configuration.

water vapour. However, for reactive metals like iron which may be covered in oxide the situation is not so clear. The lubricant has a complex surface onto which to build its film.

The films themselves can be physically adsorbed, chemisorbed or chemically formed. The last involves van der Waals forces and, because it is not an activation process, is independent of temperature and reversible. Chemisorption is an activated process and is not reversible. Chemically formed films are generally much thicker and are formed by a wide range of lubricant-material pairs depending on conditions. Usually the adhesion is polar bonding. Often the strength of attachment of the molecules is very large and that attachment of the boundary lubricant occurs through the chains as well as through the active groups.

The classical chemically-formed boundary is that formed between fatty acids and electropositive metals. These films may grow appreciably, perhaps up to 1000 nm.

Comparatively little is known about these films which have often been compared with grease. The films can be in a mesomorphic state and are usually not capable of withstanding even small shear forces. This obviously makes them a basis for very low-friction coatings like the Langmuir–Blodgett films.

### 7.4.7.2 Mechanical properties of thin boundary layer films

When sliding begins the lubricant film is subjected to a shear stress and the idea is that the film yields along a certain shear plane. Much thought has been given to the shear strength and, in particular, how it varies with pressure. Again the real problem is that the properties of these thin films may well be very different from those of the bulk because of the ordering or structuring of the layers under the extreme pressures which exist. Shearing of the film is thought to involve the sliding of extended molecular chains over each other, which is a relatively low-energy process when compared with the more complex shear motions which occur in bulk shear. The bulk properties of some of the organic polymers have been quite a lot higher than those of the film. An order of magnitude is not unusual.

### 7.4.7.3 Breakdown of boundary lubrication

Boundary layers are very thin and hence they cannot greatly influence the normal stresses at the points of closest approach of the asperities. Thus fatigue-wear mechanisms are little affected. What is more, the interaction of surface active materials with cracks and flaws could well decrease the strength of the surface interface and increase brittle fatigue. The major role of such boundary layers must almost by default therefore reduce the severity of the adhesive types of wear.

One question which should be asked concerns the influence of the surface topography on boundary film action. Do the features of surface geometry, such as large surface slopes or high curvature, in any way inhibit the ability of the molecular chains to separate the surfaces? The answer to this and similar questions is not straightforward. The reason for this is the difficult one of trying to justify the existence of boundary layers at all. There is a growing overlap between the boundary layer properties and those properties which seem to be possible because of micro- or even nanocontacts and the resulting phenomena of micro EHD. At these small scales of size it becomes difficult to differentiate between the behaviour of solid films and the rheology of interface fluids under EHD conditions. It is therefore difficult to judge what the effects of the topography are. There is little evidence that topography improves reactive chemistry, yet on the other hand there seems to be no doubt that the 'micro EHD' and hence quasi-boundary properties are topographically controlled. Perhaps this area should be a topic for future investigation. Does boundary lubrication exist? The implication is that surface topography on the nanoscale would not be very important if it does, yet on the other hand if boundary lubricant properties are micro-contact controlled then surface topography definitely is an important factor. The problem at present is that of doing practical experiments and monitoring them at molecular scales of size to find out the true situation.

## 7.5 Surface roughness, mechanical system life and designer surfaces (see also chapter 6 for structured surfaces)

### 7.5.1 Weibull and dual-frequency-space functions

Surface texture can enhance or detract from the performance of systems involving relative motion between two surfaces and a lubricant. However, so far what has been considered are the steady-state performance characteristics. These include the load-carrying capacity and the traction, the coefficient of friction and so on. All these are properties which are maintained during the normal life of the bearing or gear. There are other factors which need to be considered. One is concerned with how the system builds up from its initial state to its steady state—the run-in period. Another is concerned with the influence that the surface roughness has on failure.

Reverting back to the Weibull curve for failure statistics shows that there are three regions in which it is useful to examine the role of roughness. Because the Weibull curve is conventionally a function of time and hazard, it is logical to do the same with the workpiece in terms of its normal performance and the factors which influence its good working or failure.

Often the tribological characteristic of interest such as friction or wear is time dependent. That can be seen in the Weibull curve shown in figure 7.88, which illustrates the failure modes of components in terms of their geometrical and subsurface characteristics. As a result of these failure mode patterns conventional signal processing methods are unsuitable for characterizing them. They need revising. This was realized by Aouichi *et al* [171] in considering the stages in the wear regime in the dry sticking of steel on steel. They suggest moments of the spectrum based on the works of Marks [161] and the ambiguity function [160]. They actually use the moments of the spectrum in the same way as Whitehouse and Zheng [162] in manufacture, but they do not realize that the Wigner distribution describes these curves much better than the ambiguity function or the use of wavelets.



**Figure 7.88** Effect of surface finish and different wear regimes within time.

Aouichi *et al* use frequency moments

$$m_x^n(f) = \int_{-\infty}^{\infty} t^n F(f,t)\,\mathrm{d}t \tag{7.209a}$$

where $F(f, t)$ is the Fourier spectrum. They maintain the chronology of events by the local frequency moments

$$m_x^n(t) = \int_{-\infty}^{\infty} t^n F(f,t)\,\mathrm{d}t \tag{7.209b}$$

in which $n = 0$ corresponds to the case of total energy $M_x^0(t) = E(t)$, $n = 1$ is central frequency $\zeta(t)$, and $n = 2$ the nominal bandwidth $\Delta f$. These are the same as the moments used in metrology where $F(f, t)$ is replaced by the probability density and $t$ becomes $x$.

They also used the cepstrum (the inverse Fourier transform of the log of the power spectrum) to get an estimate of the mechanical system parameters (i.e. excitation time). The cepstrum is very similar to the auto-correlation function in which large amplitudes are weighted down logarithmically. The use of the Wigner distribution to monitor the way in which the run-in geometry changes is now a practical possibility. The possibility of using dual-space-frequency functions to characterize many aspects of the running-in process, such as the wear scar volume, friction forces or the surface geometry, by using either the local spatial moments or the local frequency moments is now emerging.

### 7.5.2 Running-in process—shakedown

These are imperfectly defined engineering terms and the names indicate that something happens or needs to be done to a system before it can be considered to be in its normal working mode. The run-in condition is usually identified with very little wear and/or low friction, whereas the initial condition of working may well be characterized by the presence of severe wear or wear debris and high friction. Another characteristic of the run-in machine or system is that of having a high load-carrying capability. Often a machine is put through a sequence of run-in periods with increasing severity in order to increase stepwise the capability and at the same time reduce the risk of catastrophic failure due to scuffing where the localized overloading of workpieces takes place.

In essence, the running-in period is one where the various attributes making up the normal working characteristics are allowed time to adjust to an equilibrium balance. The changes which result from these adjustments are usually a combination of many things, such as the alignment of axes, shape changes, changes in the surface roughness and the equalizing of various mechanical and chemical properties between moving members, such as the microhardness produced by selective work hardening or the formation of oxide and other boundary layers. All these changes are adjustments to minimize energy flow, whether mechanical or chemical, between the moving surfaces. In the words of experimenters the parts 'accommodate' to each other [154, 163].

Although the subject of running-in is somewhat nebulous there can be no question that the most obvious changes occur on the surfaces or just under them. For example, large plastic strains develop in response to the applied loads. These are necessarily accompanied by significant changes in the microstructure and the crystallographic texture, that is the preferred orientation of crystals. Local changes in composition also occur, and surface layers will be formed, modified or removed [155]. These layers can include contaminant layers, well-bonded coatings or layers of reaction products such as oxides, sulphides, solid lubricants and sometimes debris or transferred material.

### 7.5.3 Surface roughness 'one-number' specification and running-in

Running-in is regarded as the transition from the initial start-up of a system characterized by a high coefficient of friction and general reluctance to motion, heat generation and mild wear of the mating workpieces. The condition of run-in is characterized by a low or lower coefficient of friction and minimal wear. It has also been referred to as 'shakedown'.

It is now generally accepted that one process of running-in could only be explained by the transition from plastic conditions of contact to elastic conditions. In the former material is removed or moved at contact until a steady-state regime can be maintained. This means then that the contact mode has become elastic—the geometry recovers its initial shape after the load is removed and this then happens time after time. Obviously in this regime the geometry of the two surfaces and to some extent their material properties must conform to one another, otherwise the steady-state condition would not hold. It could be argued that there is no reason to

measure the surface texture at all during the running-in process: once run-in, let it run. The running-in process, although ultimately achieving a worthwhile result, is time consuming and therefore expensive and the aim is always to get to a working situation as quickly and cheaply as possible. A good possibility is to try to mimic the geometry and material properties of the rubbing, sliding or rolling parts directly by means of the manufacturing process. In order to be able to do this, it is necessary to have a very good idea of what the best geometry of a run-in workpiece actually is. Then, by examining the texture, form and material properties, it should be feasible to design or develop or even modify the manufacturing process to resemble it as near as possible. At one time it was thought that this was possible using the simple $R_a$ value (or CLA, AA) as a measure of the surface, but in view of what has been said earlier it is not surprising that this simplistic method of describing the surface proved to be unsatisfactory. This does not mean that the blame was necessarily on the design or the manufacturing engineer. It often meant that no suitable means of measuring, let alone characterizing, the surface existed. Today, with the techniques at the disposal of the engineer, a much better attempt can be made of finding the ideal workpiece characteristics which correspond to the run-in or steady-state performance. The problem to some extent resembles the inspection procedure outlined earlier. In the first case the static conditions are satisfied (i.e. alignment of dimensions) and then the dynamic conditions have to be satisfied, which presumes shape fidelity or at least some conformity, although alignment can be changed by dynamic forces. The long-term steady-state performance requires more of a balance in the energy equation; running-in cannot help the static situation but it can help to a small extent in the dynamic situation and greatly in the long term, as shown in figure 7.89.

The future trend is shown in figure 7.89. By far the easiest intervention starts from the left in the figure and gets progressively more difficult to reduce by purely manufacturing means towards the right of the figure.

More often than not, although the benefits of an as run-in part are appreciated, many engines and systems are still allowed to run-in with chance changes in loads and speeds during their early life. It is possible to help the running-in process by means other than step-by-step increases in loads and speeds rather than the random but cheap method mentioned above. One is to include certain oils of, say, low viscosity which permit the run-in to occur under milder conditions. The oil tends to prevent metal transfer, for example. Chemically active additives in the special run-in oil can provide further protection and at the same time enhance the production of smoother surfaces.

Before considering what happens when run-in fails and scuffing or welding occurs between surfaces because of thermal run-away, it is probably useful to consider the way in which the geometry of the surface



**Figure 7.89** Requirements for functional use.

actually does change during the running-in. The cases most often cited as typical examples are cylinder liners, piston rings, gears, cams and rolling bearings. Most often the running-in is associated with a high degree of slide, but not always. Figure 7.90 indicates the way in which running-in relates to other tribological regimes.



**Figure 7.90**  Running-in and boundary lubrication.

### 7.5.4  Designer surfaces—running-in

There are many suggestions for the mechanism of running-in. Most consider that the asperities are smoothed. As suggested before, the steady surface could be that one in which the curvature at all points is low enough to preclude plastic deformation. Other possibilities have been suggested in which the surface gets effectively smoothed by the action of wear, debris filling up the valleys [221], although it is still not clear how the debris could get so firmly attached to the valleys.

Sometimes the smoothing of the asperities is attributed to a glaze [223] of $Fe_3O_4$ on the surface in the case of steel or iron surfaces. Another non-geometric effect is the phase condition of the surface. It is noticeable, for example, that two-phase materials offer better scuffing resistance than do single phase. Whatever the reason, most people accept the premise that the surface texture is important. How important or what exactly needs to be measured is virtually unanswered.

Thomas [224] tried to indicate what researchers have actually measured. The breakdown is shown in table 7.5.

**Table 7.5**

| Parameter | % used |
|---|---|
| $R_a$ | 15 |
| $R_q$ | 8 |
| Skew | 8 |
| Kurtosis | 6 |
| Zero crossing density | 8 |
| Mean peak height | 4 |
| Std dev. of surface peak heights | 6 |
| Mean valley depth | 6 |
| Absolute average slope | 11 |
| Mean peak curvature | 6 |
| Std dev. of curvature | 4 |
| Autocorrelation | 8 |
| Bearing ratio curve | 8 |

This table shows that just about everything that can be measured is measured but mostly the $R_a$. This does not reflect its usefulness; it just indicates that the $R_a$ is best known and is available on most surface-measuring instruments.

Figure 7.91 shows the shape of the amplitude distribution curve when the running-in has taken place. In this typical run-in situation, it is clear that the appearance of the surface changes dramatically, although in fact the numbers do not. In an accelerated-wear run-in experiment [225] the way in which the various parameters change is as shown in the normalized graph of figure 7.92.

The parameter which changes the most as the running-in progresses is the average peak curvature, which changes by 60%. The average valley depth on the other hand changes by a mere 10%.



**Figure 7.91** Effect of running-in on profile.



**Figure 7.92** Effect of running-in on profile parameters.

In fact the height parameters do not change greatly. The very fact that the peak curvatures change the most, as seen in figure 7.92, is an indication of the probability of a plastic flow to elastic transition as being one of the most, if not the most, important mechanism of run-in. Some idea of the mechanism can be found from theory; for example, consider a Gaussian surface being run-in (figure 7.93).

This figure shows the relationship between the autocorrelation of the original surface against that of the run-in surface. Compare the correlation function independence length of the original waveform $\rho$ with that of the truncated one $\rho_T$. It is straightforward to find one in terms of the other. Thus [226]

**Figure 7.93** Relationship between correlation function of original and truncated profiles of Gaussian surfaces.

$$\rho_T = (\rho \sin^{-1}\rho + \pi/2 + \sqrt{1-\rho^2} - 1)/(\pi - 1). \tag{7.210}$$

This equation shows that the truncated waveform has a lower effective correlation than the original, that might look a bit unlikely until it is realized that the run-in profile as seen in figure 7.91 introduces very high curvatures and hence high frequencies just at the truncation level, which in effect increases the bandwidth of the signal waveform.

It is only when the truncation falls below the mean line that there is a very large drop of the correlation function. This drop is pronounced near to the origin which implies a high-frequency cut. This in turn indicates that high frequencies are important and that the removal of them constitutes a run-in surface! One way to see the changes in topography with run-in is to examine the changes in the individual parameters; another equally valid way is to use the bearing (or material ratio) curve. This was designed in the first place to help in determining the performance of bearings. What is especially good about this method of visualization is that it really does help to identify the key change that has occurred. This change is that the whole surface has been transformed from being characterised by a single statistical model to that requiring dual statistics. What this means is that the original surface is more or less preserved—in the vicinity of and lower than the mean line. Above this there is a modified surface from the top summits downwards to the truncation level. This has been described by Williamson [3] as a transitional topography. This is seen more clearly in figure 7.94. The material ratio curve is exactly what is to be expected from a plateau-honed surface which is specifically engineered to mimic this type of run-in surface.



**Figure 7.94** Transitional processes.

Running-in and the scuffing which occurs when the running-in process has failed can happen in all sorts of operating conditions which vary from high slide rates to low and can cover wide slide/roll ratios. It is not confined to any one regime of operation.

Using these two approaches, a knowledge of the actual geometry of a satisfactorily run-in surface and the mechanism of the running-in process, it becomes possible to design a surface geometry which would be immediately satisfactory. Thus it has to have peaks that will support load satisfactorily and valleys, preferably deep, to contain any wear debris and carry oil. The cross-section of a unit would look like that in figure 7.95(a).

The unit event of the 'perfect' surface characteristic shown in figure 7.95(a) need not be deterministic in the sense that it is machined as part of a cutting or turning process. It is most likely that it could be an 'average' surface topography feature that satisfies the run-in requirements. In this case it would have well-understood properties of statistical symmetry as well as having average geometric properties.



**Figure 7.95** Designer surfaces.

For example, in deference to its being virtually a two-process finish the skew of the distribution would be highly negative. This is obvious from its shape and would more or less assure the very deep valleys and rounded peaks. However, skew by itself would not guarantee the wide, deep valleys. Here kurtosis would be a help. Low kurtosis would dictate an open type of surface that would be more conducive to trapping debris (figure 7.96(b) rather than 7.96(a)).



**Figure 7.96** Low-kurtosis surface.

There is no necessity for this surface to be random in geometry. In fact to generate the required attribute of the areal picture as in figure 7.95(b) it might be easier to have a deterministic surface, in which case more control could be exercised over the specific requirement such as in figure 7.95(c).

The most successful 'designer' surface is probably plateau honing, but there is no end to the possibilities. It seems possible that the micropatterning of surfaces as advocated in Russia [227] might be the method of surface generation most suited to this flexible approach. Another possibility is that of energy beam machining.

There are examples of design for function. One example is in ceramics [228]. Although ceramic materials are exceptionally well suited for highly loaded functional surfaces, their brittleness is a problem. It has been found that if instead of using just the honing process, polishing after honing is used instead, considerable improvement in wear characteristics results. In the case of silicon nitride, $Si_3N_4$, in addition to the new process a new surface parameter could be used—the valley sharpness. The two together enabled rolling properties with paraffin lubricant to be greatly improved. Quite often the requirement is for a pattern on the surface with highly controlled parameters [353]. Certain increases in wear resistance have been reported. Contact stiffness can also be improved by vibro-rolling as can fatigue wear [354].

Another designed surface is that of laser machining of ceramics [355]. It has been found that by using an excimer laser to modify the surface of ceramics a regular cone-like structure results (microspherical caps) which can be used to improve adhesive bonding [356].

In all methods there has to be a close link between understanding the geometry of the function and the mechanism to produce the function. By using these it is possible to develop an immediate generation process which will reduce the run-in or time taken to reach the functional steady state.

Scuffing is likely to occur as a failure under conditions where there is excessive sliding under large pressures. Under these conditions the film thickness is very small, almost approaching that of boundary lubrication. The very high temperatures generated as a result of asperity interaction, usually at the inlet cause a thermal run-away to occur and catastrophic failure.

One of the problems encountered in doing experiments to see the influence of roughness is that it changes considerably during the actual test. It can be measured at the start and also at the end of a test if it has failed due to scuffing, but not in the middle of an experiment. It is when there is a large slide/roll ratio that scuffing is likely to occur, and there is likely to be severe asperity contact.

Another reason why scuffing (sometimes called scoring or spalling) is not well understood is because it is not as a whole progressive. This is in contrast to fatigue, corrosion and abrasion which tend to be slowly progressive phenomena. Hence designers can predict life expectancy with some confidence. Scuffing is usually treated like the yield strength of a material, i.e. it is often thought that there is a safe regime within which safe operation is guaranteed. The problem with scuffing is that it is typified by very sudden failure; in fact it once used to be defined as 'sudden failure'.

### 7.5.5 Scuffing

It is easy to see how the concept of scuffing or spalling is associated with heat generation and catastrophic failure. Heat generated at an asperity collision is dissipated into the lubricant locally which lowers its viscosity, resulting in a thinning of the film thickness, which in turn causes even more contact and so on until welding occurs.

The way in which the presence of the thin film in a quasisolid state can inhibit weld growth is pictured simply in figure 7.97. This picture is an estimate of what happens in practice.

Another mechanism sometimes advocated to reduce scuffing without the need for running-in is to control the formation of oxides [233].

The actual criteria for scuffing are not well accepted despite the fact that everyone knows when a failure has occurred. Scuffing is a result of failure of the lubricant system. The lubricant is specifically meant to prevent the failure of sliding surfaces in high-speed and high-load service. Liquid or quasiliquid lubricants are mainly required to distribute some or all of the load over the entire nominal area, rather than upon widely separated asperities. Another role is to reduce asperity heating when contact occurs so as to reduce the avalanche mechanism outlined above. This requires a careful consideration of lubricant access and its conductivity. (Here the lay of the texture can be important.)

Another function of the lubricant is to remove as much of the heat as possible from hot regions on the surface that may have been deformed by heat, thereby reducing the possibility of misalignment.

**Figure 7.97** Boundary lubrication to reduce welds.

The incidence of scuffing is undoubtedly related to the asperity contacts through the film. The probability of contact is obviously related to the probability of the gap statistics. The mean gap will be determined by the load and mode of deformation. It follows therefore that scuffing depends on the load and the rate of change of asperity heights (figure 7.98).

A knowledge of the summit distribution can therefore be a useful factor in determining how sensitive the scuffing dependence is on load. The simplest assumption is that this is an exponential relationship. Carrying this further, a knowledge of the joint probability of summits at a given curvature and height can also be relevant:

$$\int_{Z_L}^0 \int_{C_L}^0 p(z,C)\mathrm{d}z \ \mathrm{d}C = p(\text{debris}) \cdot \tag{7.211}$$

Thus the probability of debris being formed in an impact will be related to the height of the asperity and its curvature. This is not a one-to-one relationship, of course, but it is a useful comparative indicator.

As the simple height relationship of the summits is exponential the thickness is constant as a function of load.

It can be seen therefore that the task of matching the manufacturing process to the function in this case can take one or two routes. One is deterministic and involves designing a manufacturing process to give the preferred



**Figure 7.98** Rate of change of asperity height.

run-in geometry. The other is to accept a certain amount of run-in period, and to determine the manufacturing process-probably be of a random abrasive nature—which will reduce or minimize the possibility of scuffing during run-in by having summit properties which change with height in a controlled way determined from the probability density of height and curvature. Many of the critical applications in gears, cams, ball bearings, cylinder liners, etc, could be tackled in this way. It has been shown earlier that designing engineering surfaces to fit a function is not confined to those involved in carrying load; it can equally be extended to optical performance. That these new approaches are being taken seriously is seen because of the growth in the use of new multiprocesses such as plateau honing and new procedures like programmed stepwise running-in [234].

The idea of designer surfaces to produce the optimum geometry for the surface is not the complete answer because the subsurface layers at the interface are also important. A balance has to be reached, based on experience, as to which to address most strongly—the surface geometry or the subsurface physical properties.

Even though scuffing is principally a temperature effect the thermal effects are not completely understood. Early work on non-lubricated contacts by Archard [235] and Welsh [236] is coming under some scrutiny and is revealing evidence of a 'scale of size' effect [237]. The use of steady-state theory on the heat generated at a single point of surface contact seems to give flash temperatures that are considerably lower than practical evidence suggests. Instead of the 7–800°C normally expected, which results in the well-known 'white layer' deposit typical of scuffing, temperatures nearer 1 100°C are found.

The transition temperature in boundary lubrication, above which friction and damage rise considerably, is very much influenced by the surface texture [238]. So, in boundary lubrication and hydrodynamic lubrication the roughness is important indirectly, in that it influences another parameter which in turn affects the performance.

### 7.5.6 Rolling fatigue failure (pitting, spalling)

Figure 7.99 shows how the failure of ball bearings and single-point contact configurations fit into the general scheme of wear tests of highly loaded machines.



**Figure 7.99** Rolling fatigue failure.

### 7.5.6.1 Rolling failure

Failure under rolling contact is a fatigue phenomenon. Whereas in scuffing the failure is due to the asperity contact in the tangential direction producing plastic flow, high friction and a rise in the temperature, the mode of failure in pure rolling is normal and repetitive—called pitting or pitting corrosion.

In well-lubricated rolling-element bearings there should be no progressive visible wear due to adhesion or abrasion mechanisms, and there should be no abrasive particles produced. Such a system is limited by fatigue mechanisms. Large wear particles are produced after a critical number of revolutions. Prior to this there is almost no detectable wear. However, as soon as pitting starts the bearing life can be said to be finished. Thus it is usual to speak of the useful life of a bearing rather than the wear rate, as would be appropriate in abrasive wear.

Although, as in hydrostatic bearings, there should be no solid-solid contact, there is a considerable pressure applied through the film. This compresses the surface but at the same time it produces the well-known shear stresses below the surface where the onset of plastic flow and fragmentation begins. As the rolling proceeds all subsurface elements are subjected to a stress cycle corresponding to the passage of each ball or roller. It is this stress cycle which is the source of dissipated energy. If the stress amplitude is high the bearing will ultimately fail from fatigue. If the materials are perfect and only rolling occurs [239] the position of shear stress would be at the Hertzian position about $0.5a$ below the surface. If sliding is included, as is true for many gear teeth, then the sliding traction would move the position of the maximum shear stress nearer to the surface. Practically, because the material is never perfect, the actual position of failure depends on the presence of microcracks near to the surface produced perhaps by heat treatment or machining.

An empirical relationship between load $W$ and life $L$ is that the bearing life will be reached by 90% of components if

$$W^3 L = \text{constant.} \qquad (7.212)$$

Often the wear, called delamination by Suh [240], has what is considered by some to be a fatigue-like mechanism. It does not attach much importance to surface texture, so it will not be considered in detail here. It applies to conditions of very high shear in sliding conditions. The mechanism of delamination seen as a fatigue mechanism is as follows:

1. When in a two-sliding-surface situation (not rolling) the asperities on the softer member flatten by repeated plastic flow then at each point on the softer surface there is a cyclic loading as the harder, rough surface passes over it.
2. As this continues voids and cracks are nucleated below the surface. Crack nucleation at the surface is prevented by the triaxial compressive stresses existing immediately below the surface.
3. The voids and cracks link up to form long cracks below the surface.
4. When the cracks reach some critical length they break through the free surface, resulting in plate-like particles. This is not the same as pitting although it is fatigue, which is much more localized underneath the surface.

It is interesting to note that it was in pitting that the importance of surface roughness was first noticed and that early work on this was due to Dawson in 1962 [241]. A redrawing of the first diagram showing this is given in figure 7.100.

Notice in figure 7.100 that the correlation between the number of revolutions to pit and the $D$ ratio holds over a very high range—almost three decades. However, the roughness measured is the *initial* surface roughness and not necessarily the roughness that exists when pitting actually occurs. $D$ is ratio of surface roughness/film thickness.

These experiments were repeated [242] for real gears, not discs, and showed similar trends but at much lower revolutions—presumably because of the different way the surfaces were produced.

Failure by mechanisms other than these due to fatigue and the $D$ ratio can occur if the material is impure.

In roller and ball bearings the relative slip between the parts is low (i.e. the relative velocity is low). It is generally accepted that the surface finish under these conditions is critical. Dawson [241] indicated that in

**Figure 7.100** Pitting failure (after Dawson).

order to maintain an EHD film the $D$ ratio should be about 2.5 or more. Tallian *et al* [243] go further and suggest a ratio of 4.

In an attempt to quantify the effect of prolonged rolling Leaver and co-workers [244, 245] carried out some experiments with taper roller bearings. They found that there was a significant change in the surface texture as the experiment continued. Of particular interest was the way in which the filtering of the signal affected the result, making it very clear that the digital analysis of such results had to be carried out with extreme care. Also it seems that carrying out experiments on realistic engineering surfaces—in this case a taper roller bearing—can introduce so many complexities (e.g. friction between the rib and the roller) that the results are confusing, whereas attempting to simplify the rig to, say, a disc machine produces results that do not correlate with practice.

The generally accepted theoretical model of rolling contact fatigue not involving roughness is usually based on the work of Lundberg and Palmgren [246] who examined the possibility of defects being in the material and being the source of failure. Taking this further [247], discs are often used as simulated gears and cams and four-ball machines are often used to simulate ball bearings, as seen in figures 7.101(a) and (b). In the latter the spindle ball is chucked and rotated in a spindle whereas three identical balls are spaced by a separator.



**Figure 7.101** Two-disc (a) and four-ball *(b)* simulation machines.

The four-ball tester has been used extensively in the oil industry to investigate the properties of oil additives and similar subjects such as the rheology of oil etc. It has also been used in investigating the effect of surface roughness on ball bearings, notably by Tallian *et al* [248] very early in 1963 and, although not precisely predicting the theoretical results, it did get order of magnitude agreement.

Their experiment was a classic in the sense that it was unique in utilizing random process theory and the practical determination of contact time using electrical contact methods.

### 7.5.6.2 Roughness effects on 3D body motion

Figure 7.102 is a simple diagram which illustrates the various possibilities in a substantially rolling situation.



**Figure 7.102** Contact phenomena: static and dynamic behaviour.

In the figure it is assumed that a fluid film is between the objects although, for example, in spinning situations this is hardly possible because no input wedge is set up.

The roller situation refers to gears, cams, and taper and roller bearings. Balls here refer principally to ball bearings. Even in these cases the situation is invariably made more complicated by the presence of restrainers like cages and lips, so at best only an approximation to what happens in reality can be achieved.

It is beyond the scope of this book to go into the dynamics of rolling elastic or plastic bodies in any great detail but it is necessary to see at least some of the theory to get an idea of where roughness might be important.

Obviously a ball on a flat has five degrees of freedom (figure 7.103) with two translation and three angular motions. Because of the symmetry of the real system, two of the translations $x$ and $y$ and two of the rotations $\omega_1$ and $\omega_2$ can be regarded as the same for the purposes of applications, leaving slide $x$, roll $\omega_1$ and spin $\omega_3$. For a cylinder translation $y$ is usually inappropriate, $\omega_2$ impossible and $\omega_3$ unlikely.



**Figure 7.103** Contact phenomena: rotations of a ball about a point.

In the case of a ball the situation is quite complicated, even for small elastic deformation due to 'creep' $\zeta$ (figure 7.103):

$$\zeta = \left| \frac{2rmR - L}{L} \right| = \frac{\text{distance slipped}}{\text{distance rolled}}. \tag{7.213}$$

Parts of the circumference travel different amounts in one revolution, as shown in figure 7.104: some 'stick' and some 'slip' figure 7.105. The effect of roughness therefore in these different regions is quite pronounced, although not yet comprehensively analysed.



**Figure 7.104** Creep diagram.



**Figure 7.105** Plan of creep.

In sliding the running-in process was addressed first followed by failure mechanisms. In the case of rolling the running-in is not as important an issue, although it does occur and can be shown to produce much the same resultant surface geometry as in sliding. However, the differences in traction and wear are not so marked and the failure tends to occur a long time after the running-in in the form of pitting rather than the catastrophic scuffing which tends to occur during running-in in sliding.

In practical systems it is not possible to have a body having just one motion (i.e. rolling). In a roller bearing friction occurs in the contact of the rollers with the main races, side flanges and the cage. Rolling together with sliding occur, whereas in the contact of the roller with the side flange only sliding occurs. Ideally, fluid lubrication should hold at all places, but this is not always possible. For example, in spin conditions on balls no wedge geometry occurs to help the formation of the film so that at the point of contact it is solid-solid contact [249]. At best it is complicated, so much so in fact that conventional methods of description are sometimes rather inadequate. In fact the principles of fuzzy logic are now being used to describe such things as the wear of ball bearings in operation [250]. Also it is often difficult to separate out geometric effects from others. It is sometimes the case that there is a close interaction between physical effects. One example of this is work hardening and texture [251]. The accommodation takes place in work hardening—the

softer element gets harder as the initial surface texture of the harder element gets smaller. The overall pitting resistance thereby increases greatly. It is a fact here that although smooth surfaces help in the formation of a fluid film as pointed out by Dawson, and the presence of the film reduces contact and hence pitting, a rougher surface on the harder member increases the work hardening and this also increases the resistance to pitting. Ideally, therefore, these mechanisms should work in series. The initial rougher surface helps the run-in of physical properties. The run-in of the geometry helps the working life of the system after this by encouraging the development of the EHD film by virtue of the fact that the member is now smoother. Of course, the film could be established in the first case by having smooth surfaces, but this situation is not altogether satisfactory because, should the film be broken by debris, for example, pitting will soon develop due to the poor work hardening. Some insurance against film breakdown is advisable, especially in high-speed ball bearing applications where lubricant starvation is an important effect [252].

A very simple, yet elegant way of showing the influence of roughness on rolling elements has been made by Dowson and Whomes [190]. By rolling a rough cylinder down an inclined plane immersed in a viscous fluid they showed that the coefficient of friction was actually less for the rough surface than for the smooth— a fact that is intuitively difficult to anticipate. The roughness tested was in the direction of roll. Unfortunately only very simple roughness shapes on the roller were used, but the principle was demonstrated. The analysis involved incorporating the roughness into Martin's methods for rollers [253] followed by the Reynolds equation. Both predicted the same result.

What happens is that the load capacity, surface traction and horizontal force all increase, but the load capacity increases fastest so that the coefficient of friction reduces. This reduction of the coefficient of friction only works for low $D$ values as seen in figure 7.106.



**Figure 7.106**

In practice, a minimum is reached, after which, if the surfaces become rougher, the greater amount of direct surface interaction increases the coefficient again. To give some idea of the actual magnitude of the change in the film thickness brought about by the presence of the simple waveforms, the results of Dawson [241] were:

$$\text{square} = \frac{h_{sq}}{h_0} = \frac{1 + \frac{1}{2}d/h_{sq}}{1 + d/h_{sq}}$$

$$\text{sawtooth} = \frac{h_{saw}}{h_0} = \left(\frac{h_{saw}}{d}\right)\ln\left(1 + \frac{d}{h_{saw}}\right)$$

$$\text{sinusoidal} = \frac{h_{sin}}{h_0} = \left(1 + \frac{d}{h_{sin}}\right)$$

(7.214)

where $d$ is the peak-to-valley ratio of the roughness. This analysis, based on Martin's, did not allow axial flow whereas the Reynolds analysis did. However, the results give an idea of what the relationships are for simple waveforms.

From equations (7.214) it is seen that for a given value the square wave roughness produced the biggest increase in the film thickness. This indicates that the larger the kurtosis of the waveform the larger its effect!

As an example, for $D\,d/h_{\text{rough}}$ equal to 5

at
$$h_{\text{sq}}/h_0 = 0.6 \quad h_{\text{sin}}/h_0 = 0.4 \quad \text{and} \quad h_{\text{saw}}/h_0 = 0.37$$
$$d/h_{\text{rough}} \to \infty \quad h_{\text{sq}}/h_0 \to 0.5 \quad h_{\text{sin}}/h_0 = h_{\text{saw}}/h_0 \to 0$$

The square gives the lowest friction and the sawtooth the highest.

The obvious ways to improve this treatment are by using the Reynolds equation and by improving the surface roughness model. The latter has been to some extent tackled by Nayak [254] who refers to earlier work by Carter [185] and Greenwood and Trip [9].

In an aside to one of his papers [256] Garnell attempts to estimate the time taken to expel any fluid film from between the moving members and comes to the interesting conclusion that the film would not have time to squash (even if it wanted to!) with the result that it would be difficult for asperities to contact through the film.

This observation is not always supported. Jackobson [257] has shown that once EHD has been subjected to sliding and shear beyond its limiting shear stress in one direction, then the fluid has no shear strength in any direction. Thus an elastically flattened asperity in the inlet zone can recover within the film by dispelling fluid sideways and, given enough time, can break through the film and make contact. In other words, the time to dispel fluid may be reduced by having lateral paths open to it to allow quicker decay. The asperity contacts are therefore more likely than thought possible.

In the absence of any fluid film the rolling friction is to a large extent the surface adhesion if the bodies are elastic [258]. Any surface roughness obviously reduces this adhesion. One adhesion parameter $\alpha$ is given by

$$\alpha = \left(\frac{4\sigma}{3}\right)\left(\frac{4E}{3\pi\beta^{1/2}\gamma}\right)^{2/3} \tag{7.215}$$

where Whitehouse's terminology for the radius of curvature of peaks $\beta$ and the RMS roughness $\sigma$ is used; $\gamma$ is the net change in surface energy per unit area on the formation of contact. Briscoe found that at large roughnesses the friction was lower, as expected, but for slightly rough surfaces and very elastic bodies such as perspex a small increase was found due to distortion in the balance between the real area of contact and viscoelastic effects. For practical cases the *roughness generally produces a lower friction,* although it should be said that the presence of a fluid film is usually regarded as vital in metal contacts.

A simple idea of the way in which the spalls and pits could be formed by the texture has been given [259]. It seems to be plausible to give quantitative information about damage if the Hertz pressure and the roughness ratio $D$ are known. Remember that $D$ is the ratio of film thickness to some roughness height parameters.

The basis for the argument is that the whole contact zone is responsible for spalls whereas the individual contacts between asperities are responsible for the smaller damage, sometimes called micropits. If some simple assumptions are made such that the asperities have a maximum height of $h$ and a peak radius of curvature, then each asperity generates a micro Hertzian contact when it collides.

Simply

$$\delta - h_{\text{r}} = \frac{2\pi^2}{E^{1/2}}p_\alpha^2\beta \tag{7.216}$$

where $\delta h_r = t$ is the difference between the asperity height and the film thickness. If $R_t$ is the peak-to-valley height and $p_\infty$ is asperity pressure

$$p\alpha = \frac{E^*}{\pi}\left|\frac{R_t - 2h_r}{4\beta}\right|^{1/2} \tag{7.217}$$

and the mean asperity contact pressure $\bar{p}_\infty$, is given by

$$\bar{P}_\alpha = \langle p_\alpha \rangle = \int_{h_r}^{R_t}(t/2)p_\alpha p(\delta)\,\mathrm{d}\delta \tag{7.218}$$

given a height distribution which is assumed to be a function similar to a beta function

$$p_n(\delta) = K_n(\delta_{\max}^2 - \delta^2)^n \tag{7.219}$$

which has been used elsewhere, gives the results from Berthe *et al* [202].

Table 7.6 shows that for common ground surfaces of $50 < \beta/\sigma < 100$, the influence of the roughness ratio here defined as $\sigma/h$ is large compared with the ratio $\beta/\sigma$ for the surface. In table 7.6 $h_r$ is the film thickness calculated for rough surfaces and $h_s$ is the film thickness calculated for smooth surfaces. The critical pressure above which damage or surface distress will occur is between $\sigma/h_s = 0.7$ and 1.1.

**Table 7.6**

| $\beta/\sigma$ | $\delta/h_r =$ 0.5 | 1 | 1.5 | 2 | 2.5 |
|---|---|---|---|---|---|
| 20 | 100 | 1800 | 3100 | 5400 | 7200 |
| 50 | 50 | 1100 | 2200 | 3300 | 4400 |
| 100 | 0 | 800 | 1600 | 2400 | 3200 |
| 200 | 0 | 600 | 1200 | 1800 | 2400 |
| $\delta/h_s =$ | 0.5 | 0.7 | 0.85 | 1 | 1.11 |

According to this, therefore, the critical pressure is not basically dependent on the surface asperity curvature $1/\beta$. The real damage is due to the presence or absence of contacts rather than the geometrical detail of the asperities.

The graphs in figures 7.107, 7.108 and 7.109 show some trends [260]. Figure 7.108 shows quite clearly that no micropits are formed below a ratio $\sigma/h_s$ of 1. In figure 7.107 the micropit formation rates vary in time. This agrees with figure 7.109 which shows that all micropits are formed early (i.e. $N\ 10^5$). It also suggests that the introduction of micropits accelerates the formation of spalls which is in agreement with the idea of debris damage.



**Figure 7.107** Fatigue life versus roughness ratio.

**Figure 7.108** Number of pits at $2 \times 10^6$ cycles versus roughness ratio.



**Figure 7.109** Nature of surface damage with cycles.

There have been a number of different ways of dealing with pitting damage where the surface acts as the source of noise [261] and the system the source [262]. Other variations exist in the behaviour and performance of rolling elements under normal loads. In particular, vibration of the system from external sources has an effect on pitting and fatigue life [263]. The situation becomes very much self-regenerative in the sense that the surface texture can initiate the vibration of the system, which under certain conditions usually involves low damping. This vibration acts in parallel to the normal effect of the roughness to reduce life considerably. The fact is that roughness can have more than one way of initiating damage.

Roughness can have three possible effects on pitting life:

(1) film breakdown and normal asperity contact;
(2) crack initiation due to either deep cracks or debris damage
(3) initiation of vibration in a system which produces squeeze film effects and hence more asperity interaction.

There are other cases in which rolling balls are expected to perform load-carrying duties other than in conventional ball bearings. These are in grease-packed bearings in which many miniature balls are secreted into the grease. There is no evidence that any film is produced. Indeed at the very low speeds which are usually involved it is highly unlikely [264] but the grease tends to keep the balls together by adhesion. Surface roughness tends to inhibit this grease adhesion, so the requirement on the balls is that they be very smooth.

### 7.5.6.3  *Rough surfaces and rolling*

The fact that roughness affects the life in rolling contact is undisputed. Also, as the presence of contacts through the fluid film is an important issue, the question arises as to the quantitative nature of the stresses produced by the roughness. This will be considered briefly in what follows.

Before this is examined, a review of the basic assumptions made in the static case will be given following Greenwood and Trip [24].

1. If the normal elastic compliance of the cylinder, say, is much larger than the RMS roughness then its contact width is given by

$$a = \left( 8 \frac{(1 - v^2)WR}{\pi E} \right),$$  (7.220)

This condition is satisfied when $2W(1-v^2) \geq \pi E \sigma$.
2. Under a normal load $W$ junctions form in the contact region.
3. The expected value of the density of junctions is approximately proportional to the local pressure $p$.
4. The average junction size only increases very slowly with $W$.

Nayak [254] made these assumptions (figure 7.110) to develop a qualitative theory of rolling in the presence of roughness. This will be outlined below, but it should be emphasized that the approach is very dependent on assumptions made about the number and the behaviour of junctions in the contact region.

It has been inferred earlier that in elastic rolling of a ball or cylinder the contact region is split into a 'locked' and 'slipped' region. This is a feature of smooth surface rolling contact. The slip velocity from above is $\zeta$

$$\zeta = \left( \frac{V_{10} - V_{20}}{V} \right)$$  (7.221)

where $V = (V_{10} + V_{20})/2$.

The mechanism for pitting has been investigated with respect to gears as well as balls. In fact, pitting is the commonest form of gear failure. One model [265] proposes a fracture mechanics model and simulates the gears by the contact of two cylinders. This approach seems overdone because gear design today uses empirical rules collated into the various national and international standards. It is interesting to notice that grain boundaries which block slip bands are attributed to the initiation rather than asperity contact. The worrying factor is that in order to apply finite element analysis smooth surfaces are assumed. It could well be that the initiation mechanism is being thrown out in order to get tractable simulations. Finite element analysis has a lot to answer for!



**Figure 7.110**  Nayak's treatment of lubrication.

The 'locked' region is located near to the leading edge in which the slip velocity is zero. This region extends over the range $a - 2a' < x < a$, where $a'$ is given by

$$a' = a(1 - T/\mu W)^{1/2} \tag{7.222}$$

where $T$ is the tangential load (if any) (this stops acceleration).

The frictional stress $\tau'$ for the smooth case is given by

$$\tau' = \mu p \left[ 1 - \left( \frac{a'^2 - (x - a + a')^2}{a^2} \right)^{1/2} \right] \tag{7.223}$$

where $p$ is given by Hertzian theory as an elliptical distribution

$$p = \frac{2W}{\pi a} \left( 1 - \frac{x^2}{a^2} \right)^{1/2}$$

and the creep

$$\xi' = [1 - (1 - T')^{1/2}] \frac{4(1 - v)}{\pi} \tag{7.224}$$

where $T' < 1 = T/\mu W$.

The range of the contact which is not 'locked' is a 'slipped' region in which the slip is non-zero. It extends over the range $-a \le x \le a - 2a'$ and the stress within it is $\mu p$.

Nayak uses a roughness index given by $D$ where

$$D = \frac{2 - v}{3(1 - v)\beta\eta} \left( \frac{\pi R E^*}{4W} \right)^{1/2} = \frac{0.062(2 - v)}{(1 - v)} \left( \frac{\sigma}{R} \right) \left( \frac{R}{\beta} \right)^{5/8} \left( \frac{E^*}{p(0)} \right)^{5/2} (\eta_A R^2)^{1/4} \tag{7.225}$$

where $\eta$ is the density of contact junctions, $\eta_A$ of the peaks and $\beta$ is the curvature of peaks (taken here to be the mean value). $R$ is the radius of the curved rolling member and $E^* = E/2(1 - v^2)$.

The creep as in (7.225) is shown as a function of $D$ in figure 7.111 [254]. $D$ has the physical significance that it is proportional to the ratio of the deformation necessary for a junction to develop a large shear stress over the maximum tangential elastic deformation in the contact region predicted by the smooth surface theory.

The conclusion reached for rough surfaces is as follows:

1. If the surface texture $\sigma$ is significantly less than the compliance ($\sim 1/10$) then smooth surface theory can be used freely. If not then the roughness becomes increasingly important.



**Figure 7.111** Creep as a function of $D$.

2. Qualitatively there are two distinct differences between the smooth surface and rough surface results:
  (a) The stress distributions along the $x$ axis for the rough surface case have continuous slopes along the $x$ axis, whereas the stress distributions given by the smooth theory have a distinct discontinuity at the junction of locked (or stuck) and slipped regions.
  (b) There is no such thing as a true distinction between the locked or stuck region and the slip region for rough surfaces. The slip is non-zero everywhere, from small at the leading edge to increasing monotonically at the trailing edge.

Applying this theory shows that for the case when $R$ is large the roughness effect is small and can be neglected (e.g. in railway wheels).



**Figure 7.112** Effect of roughness on rolling characteristic.

Figure 7.112 shows an example of the difference in shear stress across the contact region for the rough and smooth surface theory. Although this is qualitative it does show how the roughness is likely to affect the rolling characteristics.

Perhaps the important issue from the point of view of fatigue and pitting is the fact that on a continuous basis the shear stresses produced by the rolling motion are smaller when roughness is present (figure 7.111) than for the smooth case because of the degradation of the stick-slip region.

### 7.5.6.4  Pitting due to debris and subsequent surface effects

Another way in which roughness can affect pitting is its role in forming debris. This debris is often caused by initial wear or contaminated lubricant and can cause one of two effects: (i) a direct breakdown in the lubricant films which enables immediate asperity contact and accelerates the fatigue stress cycle failure; or (ii) the debris damages the surface which then affects the maintenance of the film and initiates the failure path. The debris can easily cause damage to the surface because, by the nature of rolling (i.e. with zero velocity at the contact), any debris within the film in the contact zone becomes trapped and pressed into the surface [194]. That the debris causes the indents can be checked by filtration. This raises the issue of two types of influence of surface geometry: one is the statistically uniform type dictated by the manufacturing process and the other is the flaw caused usually by debris, but not always. There has been a suggestion that the roughness therefore results in two types of failure:

1. Fatigue caused by repeated cyclic stress created by normal asperity interaction through the film. This causes subsurface failure resulting in micropits.
2. Fatigue caused by debris marks or deep valleys in the texture producing cracks which are incipient sources of propagation of strain. This sort of damage in the surface initiates and results in macropits or spalls.

Type 1 failure tends to occur after many cycles and represents typical fatigue failure, whereas type 2 represents short-life catastrophic failure. This type of failure would occur in much the same timescale as scuffing [267]. The latter type of failure is related to surface topography yet it can be reduced by filtering out those debris particles likely to cause damage. These are likely to be larger in size than the film thickness plus the surface roughness value of the peak-to-valley and are usually confined to debris small enough to get wedged in valleys, say about in the 40 $\mu$m to 2 $\mu$m range. There is evidence [267] that this is so. One point to note is that when debris from fatigue failure begins to occur it will proliferate damage very quickly [268]. It has been reported that filtering at 40 $\mu$m produces seven times more failure than when the lubricant is filtered at 3 $\mu$m, implying that debris less than 3 $\mu$m for a typical case simply passes through the contact within the fluid film.

### 7.5.7 Vibrating effects

Under this heading is vibration caused by and to the roughness of the workpiece, that is squeeze films and fretting.

### 7.5.7.1 Dynamic effects

One aspect of surface roughness not yet touched upon is the noise generation produced in roller bearings. Strictly this comes under the earlier heading of roundness or out-of-roundness. From a purely functional point of view it is well established in industry but not necessarily proven that the out-of-roundness greatly contributes to the noise and vibration in roller and ball bearing systems. Many companies, such as SKF, measure such vibration in assembled bearings and relate it to the texture. It is also true to say that a great many manufacturers, especially those making miniature precision ball bearings and taper bearings, prefer to relate the $dr/d\theta$ parameter of the circumferential roundness to the acceleration on the bearing. Usually the critical measurement is the maximum and average rate of change of radius $dr/d\theta$ (figure 7.113):

$$\overline{\frac{dr}{d\theta}} = \frac{1}{2\pi} \int_0^{2\pi} \left| \frac{dr}{d\theta} \right| d\theta. \tag{7.226}$$

The most important factor is the impulsive effect of an instantaneous change in the radius term. That is the impulsive effect is $m\dot{r}$—this causes the vibration and noise. If the angular velocity is constant within the time duration $\delta t$ then the impulsive effect is given by

$$m\dot{r} = m\frac{dr}{dt} = \frac{dr}{d\theta} \cdot \frac{d\theta}{dt} = mw\left( \frac{dr}{d\theta} \right) \tag{7.227}$$

Usually $\dfrac{dr}{d\theta}$ is measured over 5° which could and should represent the angular subtence of the Hertzian elastic zone of the ball in the track. The values of $dr$ correspond to the change in position of the centre of the ball imposed by the track on usually the inner race. This change is transmitted thereby to the outer ring and thence to the support case producing the vibration and noise.

There are a number of places where this rate of change parameter is measured, but the most usual and the most critical is the outer track of the inner ring—rather than the inner track of the outer ring. This imposes

**Figure 7.113** Effect of $\delta r/\delta\theta$ on variation of bearings.

most of the acceleration on the balls or roller and hence contributes most to the vibration and acoustic noise. This source of surface problem is rarely addressed in the technical literature. The usual conception is that it is the out-of-roundness of the rings which is important. This is not so, because the out-of-roundness of the rings, or at least the low-frequency content, is taken up by elastic give or compressive alignment.

The other source of vibration is the rollers or balls themselves. Some work has been done on this [261]. Basically the roughness of the rolling elements, say rollers or discs, simulating a wheel on a track, is considered to be the input to a dynamic system comprising the test apparatus (simulating the engineering problem). The output is the noise.

Whereas in most situations the flexible elements are in the actual system as above, it can be argued that the elastic deformation of the rolling bodies themselves in the region of contact constitutes the linear dynamic system. Natural resonances and damping can be attributed to the elastic bodies and elastic hysteresis.

Unfortunately the system in this form is far from linear as the force-compliance curve for elastic deformation shows. However, within small increments it is possible to use the Fokker–Planck equation [198] and involve linear systems analysis. By measuring the accelerations on the test rig and the spectra on the surface and carrying out a cross-spectral or cross-correlation exercise, it is possible to show that the surface roughness in the circumferential direction has a great influence on output noise. Note that this high correlation corresponds particularly to the high-frequency response. The lower, longer, wavelengths of higher power would have to be modified by the method of supporting the rollers.

Exercises involving cross-correlation show that the surface texture can functionally affect the operation of a dynamic system. On the other hand, exactly the opposite is possibly true, as has been seen in machine tools. This time chatter marks are set up on the workpiece because of the lack of stiffness of the machine tools. Damping is provided not only by the stiffness or damping in the machine joints but also in the cutting or machining mechanism. Evidence that exactly the same thing occurs in systems like railway lines does exist [199]. One thing that has emerged from these results is the very important role played by damping in the system. This is somewhat unfortunate as it is the least controllable.

### 7.5.8 Squeeze films and roughness

Squeeze effects are related to the dynamic character of systems in which vibration is present and which can indirectly be triggered by the roughness itself [263]. It is also related to the time it takes actually to squeeze a film out from between a contact and the actual time during which the contact exists. The two need not be the same. It has already been indicated that it takes quite a lot longer to evacuate a film than to maintain the contact if it is a rolling regime and there is a normal approach. Under these rolling conditions a squeeze effect can be produced.

In sliding also squeeze effects can be important. In general the effect of squeeze [270] is to increase the effect of roughness. This is due to the resistance to additional flow caused by the asperities (figure 7.114). However, the distinction between the extreme cases of rough surface sliding and rolling becomes less

**Figure 7.114** Squeeze film.

conspicuous with the increase of the squeeze flow parameter. The result is that the normal squeezing motion makes the stagnant fluid trapped in the roughness pockets, the cause of the difference due to roughness, move.

In terms of the Reynolds equation and using the normal convention of [270] and following Christensen, for transverse roughness with $p$ defined as

$$p = \frac{\overline{p}h_0}{6\eta\omega_s}(h_0/2R)^{3/2} \tag{7.228}$$

$$\overline{p} = \int_{-\infty}^{X} \frac{x\,\mathrm{d}x}{H^3(1 + R'^2_{max}/3H^2)} \tag{7.229}$$

and $H_T$ is the film thickness

$$W = \frac{\overline{\omega}h_0^2}{24\eta\omega_s R^2} = \int_{-\infty}^{\infty} p(x)\,\mathrm{d}x \tag{7.230}$$

then

$$\frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{\mathrm{d}p}{\mathrm{d}x}\frac{1}{E[H_T^{-3}]}\right) = -1. \tag{7.231}$$

For longitudinal roughness

$$\frac{\mathrm{d}}{\mathrm{d}x}\left(\frac{\mathrm{d}p}{\mathrm{d}x}E[H_T^3]\right) = 1. \tag{7.232}$$

The form of the equations, not necessarily their value, should be considered here because of the nature of the assumptions made to get them and the difficulty of ever applying them, other than on a computer. The value $R'_{max}$, for example, would be difficult to measure for a random waveform because it is a divergent variable related absolutely to the sample size. However, the form and general discussion is relevant. Tallian [203] gives a thorough background to the problems (also [272, 260]).

### 7.5.9 *Fretting and fretting fatigue*

This is the form of damage that occurs when two mating surfaces are in oscillatory motion tangential to each other. Strictly, this comes under three-body interactions, but because the source of the effect is vibration and is two body initially anyway, it will be introduced here. This type of fatigue is not rare although it may seem, at first impression, to be unlikely. It is almost always a direct result of vibration, usually at a joint, and so is likely to be very important in determining amongst other things the damping which occurs in, for example, machine tools. It is also obviously a source of wear in the joint.

One of the big problems in fretting is that the actual size of the debris and wear scars are very small. It is difficult to quantify. Weighing is one way and the normal approach is another.

The magnitude of the motion is usually very small. In practice it may be about *25 μm* and can be even lower at a thousand times less (*25 μm*). The essential point about fretting is that it occurs most usually

in an enclosed system, with the result (figure 7.115) that any debris formed by any means gets trapped. Any debris which is so trapped can then be instrumental in causing yet further damage. Thus it is an avalanching process.



**Figure 7.115** Fretting mechanism initiation.

The principal factors which cause fretting fatigue (which gives rise to the debris or initiates fatigue) are:

(1) local interactions between asperities resulting in local welding adhesion and general surface attrition;
(2) surface oxidation due to a rise in temperature produced by energy dissipation due to (1);
(3) abrasion caused mainly by oxide debris formed in (2).

With regard to the surface roughness it should be obvious that the asperities that interact to produce these welds are usually the major ones. These have a width in the scale of size of a few values of the correlation length of the surface. Here it is assumed that the two surfaces are similar (figure 7.116).



**Figure 7.116** Weld width in relation to asperity width.

This scale of size is the one which is the most relevant for debris size. The sort of mechanism is illustrated in figure 7.117.

This mechanism demonstrates [273] how the surface geometry is often an important factor, not only in sometimes being instrumental in the formation of debris but certainly always in trapping debris and making it a moving stress concentrator in much the same way as a ball bearing in pitting damage. The fact that the debris is usually an oxide and furthermore soon becomes much harder due to work hardening strengthens this analogy.

Pointing out the scale of size of the debris and the correlation length of the surface geometry is not accidental. Again there is a size link between manufacture and function (or in this case damage). The correlation length of the surface is naturally related to the manufacturing process parameters. Effective grit size, for example in grinding, is a determinant of the correlation length. It would seem that this in itself determines the scale of the major asperities. They are basically related to a reflection of the unit process event. In terms of

**Figure 7.117** Mechanism of fretting failure.

stress fields therefore, this scale of size is an important one and sets the scene for the scale of size of debris and, hence, fatigue cracking.

Waterhouse and Allery have tabulated a list of the ways in which fretting affects normal fatigue and, not surprisingly, they find that the fatigue curves are similar to the fatigue curves obtained in air but displaced to lower stress levels [274].

That the debris of 'fretting bridges' is responsible for the main effect has been shown [273, 274] simply by removing it and showing that, once removed, the fatigue life is unchanged.

Interestingly, above the minimum number of cycles needed to form the fretting fatigue cracks shown in figure 7.117(d) further fretting has no more deleterious impact than continuing the test without the fretting bridges. This is due to the fact that any further fretting cracks produced are less developed than those already in existence that will therefore produce the eventual failure.

So fretting is primarily related heavily to the roughness. It is concerned mainly with shortening the fatigue initiation stage and hence the life. Thus roughness is most important when workpiece life is mainly governed by initiation, that is when the first cracks form. Such cases are in the high-technology areas involving high-frequency low-stress levels with smooth or polished surfaces. Fretting is relevant to any situation in which the surface is constantly subject to vibrational loading, yet only occasionally used (i.e. a bearing or slideway).

If aggressive environments are present then corrosive effects come into effect, producing accelerated damage in a manner similar to ordinary fatigue crack initiation.

Quantitative prediction is still some way off, mainly because of the complexity of the behaviour.

Fretting failure is being identified under plastic conditions in which the roughness is not really involved and surface cracks and subsurface movement are severe [265,275].

Reduction of fretting and fatigue wear in general is a priority. Apart from the conventional introduction of residual stress by shot peening there is the possibility of using a process which helps reduce failure. One such attempt has been made by Schneider [210]. According to him his vibrorolling had many advantages over conventional processes including

  (i) Imparted compressive stresses.
 (ii) Improvement strength and hardness of the areas contiguous to the grooves, hence of the load carrying capacity.
(iii) Imparted independent features capable of multifunction optimization.
(iv) Reduced wear, corrosion, fretting, suffing.

The relative merits in running-in and wear are shown in figure 7.118.

According to the figure there is a two to one improvement over what presumably is diamond turning which is certainly significant.



**Figure 7.118** Running-in and fretting wear.

The aspect which is missed is the extent to which vibrorolling can modify size. Turning and grinding for example can remove as well as smooth material. Vibrorolling obviously moves material by large plastic strain. It seems incredible that subsurface damage is small.

The potential for this sort of process in which the roughness geometry and the physical parameters of subsurface can apparently be matched to function independently is high.

Like everything else, the best way to eliminate fretting failure is to remove the vibration in the first place and then the corrosive atmosphere! Failing that, as in normal fatigue move critical surfaces having edges or notches to regions of lower stress.

As in conventional fatigue the introduction of a compressive surface stress rather than leaving the tensile stress, which wants to open up, is a good stratagem. Shot peening, surface rolling and nitriding have all been found to increase the corrosion fatigue of steels. Surface coatings can also have compressive stresses, such as anodized layers, which are useful.

Note that many of the so-called inhibitors are designed to affect the surface itself. This is because of the central dominance of the surface. At the surface the crack itself can be initiated, but even if the crack is not initiated the necessary stresses to start the subsurface crack are communicated via the surfaces so that the surface geometry is still a very important factor.

### 7.6 One-body interactions

#### 7.6.1 General

This topic is taken to mean the interaction of a gas or a liquid or simply a physical phenomenon applied to a surface with a resulting functional effect.

It has been somewhat artificial putting fretting earlier, because of its obvious association with fatigue. However, the initial idea of what fretting is belongs to surface contact. Hence it has been lumped with pitting.

## 7.6.2 Fatigue

It is known that surfaces produced by manufacturing operations give rise to widely varying fatigue behaviour. Three factors are most significant:

1. Surface roughness: rougher surfaces are expected to encourage fatigue crack initiation especially for notch-sensitive materials and especially in high-cycle fatigue.
2. Residual stress: compressive or tensile residual stress patterns are set up either deliberately in the case of surface treatments such as shot peening, or incidentally as has been shown in chapter 6 in most machining operations.
3. Surface microstructure: since fatigue cracks spend the majority of their life in the near-surface layer, microstructural effects such as decarburization or grain orientation in castings will have strong influences.

It is well known that roughness is important [207], but how important?

Apart from the crude statements that smooth surfaces are good and rough surfaces are bad for fatigue resistance, what else is quantitatively known? The statements have been confirmed experimentally but furthering the scope of the field is still difficult. This is mainly due to the large variety of environments to which surfaces are subject. The type of loading is especially important. If fatigue failure is being caused by a small number of high loads surface roughness is considered to be less important than if the workpiece is subject to low loading with a high number of cycles.

Traditional formulae do exist which relate surface roughness to fatigue limit. The ordinary limit is multiplied by this $K$ factor (for the roughness in this case) which is always less than unity unless the specimen is perfectly smooth, that is $K$ (fatigue limit with surface roughness) $K_f$ (limit without roughness) $K_s$ (surface roughness factor), or

$$K_{fs} = K_f K_s. \qquad (7.233)$$

This empirical formula, usually called high-cycle fatigue because it relates to low loads and a high number of cycles, takes the form of a factor derived from surface roughness measurement which is then applied to the known smooth specimen fatigue limit of the metal to give a fatigue limit for components having that surface roughness.

One such formula [209] is

$$K_s = 1 - x(0.338 \times (S_u/(50-1) + 4x + (0.25 - x) \times (S_u/(50-3)^2 \qquad (7.234)$$

where $x$ is $\ln (s)/23.025\,85$. $K_s$ is the surface roughness factor at the endurance limit, $S$ is the surface roughness (presumed to be CLA, AA or $R_a$ in today's nomenclature) in microinches and $S_u$ is the tensile strength. Equation (7.234) allows for the factor that metals with a high tensile strength are more susceptible to surface finish induced damage than ones with a lower static strength.

Formulae like (7.234) are known to have severe limitations. The effect of a given surface intrusion will depend on how sharp it is, that is its curvature as well as its depth. Some reports do give modifications to the above formula which apply in particular cases. The effect of the depth of a valley is not linear; deeper ones can be less damaging than might be expected, probably due to subsurface effects. There are other modifications to fit other cases.

More general progress is to be expected from the use of fracture mechanics. This is concerned with predicting how fast a crack will grow under a given fatigue loading and with a given combination of notch geometry and crack length. It is widely used in the design of high-integrity structures such as air frames. Applied to the surface roughness problem the method requires detailed knowledge of the geometry of the surface irregularities, the loading applied and the fracture mechanical properties of the material from which the workpiece is made. A particular property is the threshold crack-tip stress intensity factor, which determines whether or not a crack will grow at all. This is well documented for many metals and the stress analysis methods underlying the use of fracture mechanics are well known. One of the difficulties concerns the surface roughness. Simple surface parameters are not comprehensive enough to be used and fatigue experts

are not usually aware of the modern refinements mentioned in chapters 2 and 3 of this book. The other main difficulty is that the fracture mechanics models for very short cracks are less reliable than those for cracks longer than about 1 mm, and it is specifically on the short-crack models that the research is needed. The short-crack problem is being investigated intensively because of its general importance and its specific applicability to the surface roughness problem [276–278].

Just how crude current information is, is shown by the version of the A 18 handbook referenced by the software handbook [209].

The available fatigue data are normally for a specified loading, specimen size and surface roughness. Thus the design life

$$N_i = K_1 K_d K_s N_1 \qquad (7.235)$$

where $K_1$ is the correction for loading, $K_d$ is for diameter and $K_s$ is for roughness.

Values of three factors for correcting standard Moore fatigue test data on steel are shown in figure 7.119. This figure is taken from the *American Metals Handbook* [276, 278]. It has to be accompanied by table 7.7.

These directives are only approximate but they do amount to what is available generally. In equation (7.235), for example, it is assumed that this relationship holds for other metals as well as steel because of the lack of other data.

The effect of corrosion in the figure is related to the surface roughness only if the fatigue does not take place within the corroding material, that is stress corrosion does not take place. Should this be the case, however, then the $s/N$ curves taken for the stress-corroding conditions should be used and the surface roughness of the part at initial conditions used. Here $s/N$ *is* the reversing stress versus number of cycles to failure—the fatigue endurance curve.



**Figure 7.119** Effect of roughness on fatigue.

**Table 7.7**

|  | Type of load | | |
|---|---|---|---|
| Factor | Bending | Torsion | Axial |
| $K_{1_1}$ | 1.0 | 0.58 | 0.9 |
| $K_d$ where $d < 0.4$ in | 1.0 | 1.0 | 1.0 |
| $0.4 < d < 2$ in | 0.9 | 0.9 | |
| $K_s$ | | from figure 7.119 | |

It has been noted earlier that the residual stress and the surface microstructure are relevant in addition to the surface roughness. It is strange that, although much work has been done to assess the influence of notches and cracks in high-cycle fatigue, no reliable quantitative models have been developed to determine the number and nature of 'notches or effective notches' which lie buried in the surface texture. It is presumed that fatigue cracks will be initiated from the greatest stress concentrators which would most likely be the deepest, sharpest grooves on the surface. Consequently it could be argued that the surface roughness parameter of most relevance to fatigue strength is that one which best describes the deepest and sharpest valley. This poses enormous problems because of the actual problem of finding it, let alone characterizing it. Fatigue in this sense is possibly the most difficult surface requirement of all. Deep valleys, especially deep, sharp valleys, bring out the worst in instrument resolution, worse even than the problems associated with finding the highest peaks. That the subject of surface roughness has only been rigorously investigated in this way in the last few years comes as no surprise [279].

Taylor and Clancy [279] did remove residual stress from the unknown by annealing steel specimens at 590$^{\circ}$C for 3 hours. This was confirmed by x-rays. For the surface specification first the $R_a$ value was used and then the $R_{\max}$, the deepest valley together with its width, polished surfaces (1 $\mu$m diamond paste) grinding, milling and shaping covering a wide if not very realistic surface range. The specimens were then fatigue tested at 100Hz constant amplitude.

The results are shown in figure 7.120 for the values of $R_{max}$. From this figure it can be seen that for a given stress the number of cycles to failure depends greatly on surface roughness.



**Figure 7.120** Influence of process on fatigue.

How is this sort of data likely to be modelled? One possible method is to assume that surface valleys are in effect small cracks and then to use the large literature on crack analysis in fatigue [280].

At first sight there is only intuitive evidence to suggest that surface valleys develop into cracks later. However, work conducted on machined notches and on long cracks with blunted tips suggests that notches may be crack-like providing that their root radii do not exceed a value dependent on the material and stress level, but typically about 10 $\mu$m. Whether these apply to surface valleys is not known, but it is a useful coincidence that curvatures of about 10 $\mu$m are the limit for geometric optics in the surface scatter of light before diffraction takes over. This fact could probably be used in making instruments for detecting fatigue cracks.

Making some sort of identification of valleys with short cracks allows the use of the Kitagawa plot [281]. Figure 7.121 is a Kitagawa plot for short-crack behaviour and figure 7.122 is a modified Kitagawa plot for surface roughness features (typical predictions for cracks and notches). The threshold stress range is the fatigue limit of material containing a given crack. The polished surface fatigue line appears as a horizontal line. The long-crack threshold ($\Delta K_{\text{th}}$) appears as a straight line of slope $= \dfrac{1}{2}$

**Figure 7.121** Kitagawa plot for short-crack behaviour.



**Figure 7.122** Relationship between roughness, crack length and stress.

The data found conforms to the lines at extreme values and the curve below the lines close to the intersection. The intersection point and the points of deviation from straight line behaviour can be defined by cracks of length $a_0$, $a_1$ and $a_2$ respectively.

It could be argued that the Kitagawa plot with roughness on the $x$ axis (as $R_{max}$) is a reasonable model.

From figures 7.120, 7.121 and 7.122 the rough surface model found in shaping does not fit too well with the short-crack model but does for the notch-based prediction, presumably because of the very curved surface profile in the valleys for the shaping processes.

Ground surfaces follow closely the short-crack model.

It seems that for relatively low roughness corresponding to grinding, polishing, etc, the short-crack data plot, with $R_{max}$ replacing the crack length abscissa is reasonable. As the roughness (and more importantly the type of cutting tool) increases, notched-based analysis takes over. Where these intersect at $a_0$ should give an upper bound to the short-crack approach [282].

Figures 7.120, 7.121 and 7.122 are based on inspired guesswork relating surface roughness to fatigue and they form an extraordinarily good basis for an understanding of the mechanism underlying crack formation. The question arises as to whether any or much theoretical work has been done. The problems are great, as expected. Nothing to do with surfaces is straightforward. Most works have failed to realize or at least to tackle the fact that the topography and the residual stress are both involved. Furthermore, the two have to be recognized as different and do not come under the same heading. Why this insistence on recognizing them? The answer stems mostly from the fact that manufacturing processes and functions have changed in the past few years. The older processes employing high-speed cutting tools and very shallow cuts left a high proportion of compressive stresses in the surface which are regarded as beneficial. The use of deeper cuts on the pretext of economy imparts more tension to the surface layer and this is bad for fatigue. The effects of laser cutting, ion beam milling, etc, usually are more stress free. One thing is clear, and this is that bulk properties

of materials are nowhere near as important as previously thought and surface properties and conditions are paramount. It is mainly these surface properties that control the fatigue strength of critical workpieces.

In theoretical terms it seems obvious that fracture mechanics should be central to the discussion. It assumes that, when it is used to select materials for structural fatigue strength cracks or flaws will be present in the structure of the surface at some point in its life.

Fracture mechanics methodology hopefully provides some design procedures which will minimize catastrophic service features.

The only work specifically to investigate such a correlation between surface fatigue and fracture mechanics [277,283] suffers from a very limited examination of the roughness, certainly nothing like enough to assess the relationship correctly. A general acceptance that the surface behaves in the same way as an edge crack under loading is a step forward. For very deep grooves (~ 50 $\mu$m and above) a theoretical linear elastic approach using fracture mechanics shows reasonable agreement with practical data, but otherwise the results tend to be not so conclusive. There is a really serious lack of work in this area!

There are, however, many indirect pointers to the importance of the roughness, which demonstrate that the fatigue life is shortened by the presence of roughness. What is difficult is linking the cracks initiated with the actual roughness itself. In reference [277] the comment is made that inspection with a profilometer is inadequate because of the inability to pick out the defects that cause failure. There seems to have been a considerable lack of appreciation of the fact that $R_a$ (CLA) does not measure the peak-to-valley parameters. Much of the lack of correlation between surface roughness and fatigue is due to poor measurement of the surfaces. Until this is corrected it is difficult to see how progress can be made.

A more comprehensive measurement routine has to be set up, perhaps using relocation methods and also using some of the latest surface-measuring instruments to locate and characterize potential crack initiators within the surface.

In the discussion so far, the emphasis has been on trying to quantify the influence of the surface roughness on the fatigue strength. There is, however, a completely different aspect to the relationship between the geometry and fatigue mechanisms. This is the use of surface topography to investigate fatigue and in particular fracture behaviour. The name given to this is FRASTA (fracture surface topography analysis), [284, 285]. Basically, this involves using the surface topography of the two surfaces resulting from a fracture being used to reconstruct the past history of the crack. This enables, at least in principle, the initiation site and its starting point to be found. Such information is obviously vital in giving true estimates of fatigue life.

Essentially the argument is as follows [284]. Consider the development of a microcrack in a material under stress. Initially, the material undergoes local plastic flow before failure begins, at some weak spot, in the form of a microcrack. The newly formed microfracture surfaces are free surfaces and consequently have zero stress with the result that the material underneath can undergo no more deformation. As a result, the applied load has to be supported by unfractured material in front of the crack tip. This will then deform, then crack, and so the mechanism proceeds.

Thus microfracture extension results from the sequential process of deformation, microfracture and redistribution of stress. This produces differences in the amount of plastic deformation experienced locally by the material as a function of distance from the microfracture nucleation site and also as a function of time after the initiation of microfracture nucleation. Thus this difference in the amount of plastic deformation experienced by the local material is a record of the history of the crack propagation.

It should be noted, however, that the path of the crack cannot be established from one surface. The crack front interacts with local microstructure features so that the total amount of plastic deformation developed at the crack is divided unevenly between the two fracture surfaces—usually called the conjugate fracture surfaces. Consequently, the total amount of plastic deformation occurring during fracture can only realistically be assessed by taking the mating surfaces into account.

Usually, the topographic information is explained in two ways (figure 7.123). One is a series of fractured–area projection plots, FAPPs, and the other a series of cross-sectional plots.

**Figure 7.123** FRASTA technique for crack history.

The area projection plots are taken parallel to the fracture plane at one instant during crack extension. This plot provides information on the location of the microcrack initiation sites and the projected area of the microcracks. By examining a series of the projection plots produced as a function of map separation (i.e. pulling the surfaces apart on a computer), a detailed picture of crack propagation can be obtained. Other information, albeit more tentative, can also be estimated from these plots. The average crack length as a function of map displacement can be found by computing the percentage of the fracture area and converting it to the actual area or crack length, assuming a rectangular window function. From this the growth rate of the crack can be assessed.

In the other form of display using the cross-sections perpendicular to the fracture surface, the way in which the two surfaces match each other can be found. This allows the plastic deformation in terms of over-lap to be found and also the amount of crack face opening displacement. The latter allows the stress intensity factor (or J integral) to be determined, from which the fracture toughness of the material can be assessed; also, it is possible to get some idea of the actual local loading conditions.

Using the topography therefore as a 'fingerprint' of the effects of crack propagation can allow (i) the fracture mechanisms to be characterized, (ii) the crack history to be characterized and (iii) the actual fracture mechanics parameters, such as the initial flow or even the toughness mechanisms, to be determined.

The basic problem with such investigations is that a great deal of background information on the mechanical properties of materials is needed before a realistic interpretation of the pictures can be given. However, the technique reinforces the statement that the surface roughness, or more properly the topography, is a valuable source of data to be exploited in many functional situations. How the data is obtained from the surfaces depends on the application, but a scanning technique obviously is needed so as to get the areal information. Optical 'follower' methods or laser scanners can be used effectively, as also can stylus techniques and the scanning electron microscope (SEM). The essential requirements are that there is good resolution (of a few micrometres) and that a wide scan area is made available. This should typically be of at least a few millimetres square. Current thinking favours the optical scan methods because of their lack of distortion and non-contact.

As corrosion fatigue is linked to fatigue, pitting and spalling it will be considered first.

### 7.6.3 Corrosion fatigue-general

Corrosion fatigue is caused by the combined action of a cycle stress and an aggressive environment. This means that all the usual factors known to influence fatigue behaviour, such as material composition, stress amplitude and spectrum, component geometry, surface finish and temperature apply equally to corrosion fatigue. The superposition of a range of corrosion conditions on top of these complicates the situation enormously.

As a result of this complexity it is obvious that there is hardly likely to be any universal theory of corrosion fatigue because the mechanisms of corrosion are so widely different in different environments and the fatigue characteristics over a range of stress systems are so dissimilar.

As is known, the fatigue failure of any surface can be considered in two stages: (i) crack initiation and early growth; (ii) crack propagation. The presence of an aggressive environment may affect either or both of the processes and an important aspect in the understanding of corrosion fatigue is to identify the controlling process, that is whether workpiece life is primarily governed by initiation or propagation.

A point to note is that under relatively inert environments it is generally found that for low-cycle fatigue life ($N < 10$) most of the workpiece life is governed by propagation (i.e. stage (ii)) and for high-cycle fatigue conditions ($N > 10$) about 90% of life is spent in initiation. However, if an aggressive environment is present the initiation stage may only account for 5~10% of the life and the major factor again becomes propagation.

Immediately this would indicate that surface roughness, which mostly affects initiation, would become less important. One of the problems in testing this and any other hypothesis in corrosion fatigue and corrosion is that of actually getting results of fatigue in a non-aggressive environment. Most people assume that air is non-aggressive, but this is not the case.

Materials have a variety of properties-physical, mechanical and chemical. Corrosion is most certainly considered to be concerned with the chemical properties, but corrosion fatigue is concerned with both the mechanical and chemical properties. Corrosion itself may be defined as follows: 'the reaction of a material with its environment with the consequent formation of a reaction product of the materials'. This does not refer to any material deterioration.

Generally corrosion has the following significant factors:

(1) rate per unit area
(2) extent
(3) form of attack
(4) location
(5) function of the workpiece subject to attack.

Obviously many of the tools of surface analysis are concerned with the spatial distribution of roughness. In so far as corrosion produces roughness, it could well be argued that roughness measurement methods could be used to classify the corroded area of a surface.

The influence of an aggressive atmosphere on a surface, however formed, whether by manufacture or use, accelerates bad functional effects as indicated in figure 7.123.

Figure 7.124 shows a schematic diagram of the possible role of surface roughness in the general area of corrosion and fatigue. It relates to the theory of roughness metrology as well as just the simple presence of roughness on the surface as a result of the manufacturing process.

Although roughness plays a direct part in the failure of the workpiece owing to crack initiation, it could play a part indirectly in the sense that much of the theory of metrology used in roughness could be used by both fatigue people and corrosion people with comparatively little effort. This would enable a better link of the theory between the two. One part in figure 7.125 shows that there is an element of roughness which is an outcome as well as an input. This is very much a situation which often arises in manufacture and function.

### 7.6.4   Corrosion

#### 7.6.4.1   General

Erosion is a two body form of wear in which particles impinge onto a surface. It is sometimes associated with the subject of corrosion—one body wear, but here is considered to be the mechanism whereby the surface is produced by bombardment by solid particles [286]. This is a wear mechanism which is intentional in

processes such as sand blasting, but more often found as a problem, as in aeronautical engineering [367]. The surface texture found initially on the surface does not in itself play an important role in erosion but it does indirectly because of the original process producing the surface. This determines subsurface stresses. In the case of ductile materials the erosion mechanism is mainly plastic flow, whereas for brittle surface the mechanism is more concerned with crack propagation from the surface.

The direct influence of the roughness is in the coefficient of adhesion $\gamma$ of particles to the bulk surface

Thus $\gamma = \dfrac{\gamma_{ad}}{HR_q}$ where $H$ is hardness of the material and $\gamma_{ad}$ *the* energy of adhesion. $R_q$ is the RMS roughness.

This adhesion wear is only marginally related to erosion as such.

This is a one-solid-body problem in which the effect on a solid surface by chemical agents is the function of interest. As indicated earlier, air is an aggressive environment because it allows the formation of oxide films. These usually change the propagation rate. The argument here is that, in a vacuum, oxides cannot form, but unless the workpiece is clean beforehand, surface films on the workpiece before the evacuation will remain!

Corrosion in general is related to the presence or the attack of water and chemicals contained in it. Obviously it is very closely related to the fretting and pitting fatigue problems. It is not the intention here to go into general corrosion. This is well documented (e.g. [287, 288]).

Corrosion is slightly different from other properties in that unlike mechanical and physical properties which can be defined in terms of environment-independent constants, the rate of corrosion of a material cannot. Because it is basically a chemical reaction the temperature, chemical composition, pressure, etc, are needed [289].

Just how complicated it is is shown in figure 7.126. Actual surface geometry only comes into the problem of corrosion:

(1) in its source of defects allowing crevice or crack corrosion;
(2) in its ability to assist or desist in the application of corrosion inhibitors in the sense that the geometry can better allow the deposition of protective films owing to their adhesive properties.

The surface defects, apart from the already-considered fretting corrosion effects that are caused by sources other than heterogeneities in the surface, are mainly from the machining. This is a more important problem than might be thought at first sight because of the modern tendency to grind brittle materials.



**Figure 7.124.** Causes of surfaces.

**Figure 7.125** Air in general increases the fatigue rate by a factor of 4 [215].



**Figure 7.126** Localized crevice corrosion.

Although this is not important in the case of ceramics, it certainly is in high-alloy metals because it can allow preferential attack, as seen in figure 7.127



**Figure 7.127** Attack in crevice-dependence on shape.

### 7.6.4.2 Localized attack—electromechanical

Localized attack can have many sources such as differential aeration, crevice corrosion, pitting, inter-granular attack, bimetallic corrosion, etc. But whatever the mechanism and source, corrosion is due to the thermodynamic instability of a material, usually metal, in its environment. The rate at which the reaction proceeds to equilibrium and the location will be affected by the heterogeneities which may be associated with the structure of the metal, the environment or the geometry.

Under these circumstances one or more areas of the surface will become predominantly anodic while other surrounding areas will become predominantly cathodic. The extent to which the corrosion is localized depends upon a number of factors.

These can be compared electrochemically in terms of the ratio of anodic current densities $i_a$ to cathodic densities $i_c$ The ratio for localized attack $R_{LA}$ is therefore given by

$$R_{LA} = \frac{\sum i_a}{\sum i_c} > 1. \tag{7.236}$$

In the case of a crack this can be very high. In pitting, for example, it can be 1000.

In equilibrium

$$I = \sum i_a S_a = \sum i_c S_c \tag{7.237}$$

where $i_c$, $i_a$ are current densities and $S_c$, $S_a$ are areas of cathode and anode activity respectively.

The critical statement about cracks is that $S_a$ is small while $I$ and $S_c$ in equation (7.237) are large.

### 7.6.4.3 Heterogeneities

These give rise to differential attack. They are shown below according to 'scale of size'.

1. Atomic (nanometric) effects:
   (a) Sites within a given surface layer-vary according to crystal plane.
   (b) Sites at edges of partially complete layers.
   (c) Point defects in the surface layer, vacancies (molecules running in the surface layer), kink sites (molecules missing at the edge of layer), molecules adsorbed on top of complete layer.
   (d) Disordered molecules at point of emergence of dislocations (screw or edge) in metal surface.
2. Microscopic (micrometre) effects:
   (a) Grain boundaries being more reactive than the grain interior.
   (b) Phases-metallic (single metals, solid solutions, intermetallic compounds) and non-metallic (metal compounds, impurities).

3. Macroscopic (millimetre and submillimetre) effects:
    (a) Grain boundaries.
    (b) Discontinuities on metallic surfaces—cut edge, scratches—discontinuities in oxide, failure or other chemical failure in applied metallic or non-metallic coatings.
    (c) Bimetallic couples or dissimilar metals.
    (d) Geometric factors, crevices, contact with non-metallic materials, etc.

Anodic areas are listed in table 7.8.

**Table 7.8**

| System/metal | Metal area which is predominantly anodic |
|---|---|
| Dissimilar metals in contact | Metal with the more negative corrosion potential in the environment under corrosion |
| Crevices, geometrical configuration which results in differences in the concentration of oxygen or other cathodic polarizers | Metal in contact with the lower concentration |
| Differences in metallurgical condition due to thermal or mechanical treatment | Cold-worked areas, anodic to annealed areas, metal subject to external stress |
| Discontinuities in films, oxide, applied metallic or non-metallic coatings | Exposed substrate-providing it is more electrochemically active than the coating |

### 7.6.4.4 Localized attack-electrochemical

Intense localized attack can and does occur in the vicinity of crevices and cracks caused by any of the heterogeneous reasons given earlier, but especially at the microscopic level. In particular cracks or pits caused by poor machining or work hardening or poor thermal history are potential problems.

In such a crack (figure 7.126) a small volume of stagnant solution can be entrapped causing intense activity. Take for example a pit or a crack in a metal immersed in sea water. Initially the metal outside the crack and within the crack will corrode at the same rate, but because of the difficulty of the renewal of oxygen by diffusion and convection the stagnant solution within the crack will rapidly become oxygen free and the corrosion of the metal within the crack will be entirely determined by oxygen reduction on the exterior surface which will become passive; hence it will form the cathode of the cell ((metal exterior to crack)/(metal within crack)). Chloride atoms will migrate from the water outside the crack to the interior. These ions favour active dissolution and work against passivation. The metal ions $M^+$ will now hydrolyse to form the metal hydroxide (according to [127]:

$$M^{2+} + 2H_2O \rightarrow (OH)_2 + 2H^+ \tag{7.238}$$

with the consequent generation of the $H^+$ ions and a decrease in pH. This results in a consequent increase in corrosion rate. The chlorine $Cl^-$ increases while the pH decreases. Charge rapidly transfers from the metal within the crack to metal outside. In other words, a really active cell has been set up caused mainly by diffusion difficulties at the boundaries set by the crack. Obviously under these circumstances cracks which are narrow and deep will be more active than those with low crack curvature in the valley.

Hence such cracks in the surface will be most active at high curvature points (figure 7.127(a)). This encourages further crack propagation in the presence of any local stress, which is why stress corrosion is so dangerous. Figure 7.127(b) shows a void which, if the metal were uniform, would not be active in the same way.

The fact that the overactivity of cracks and crevices can at least in part be attributed to the size and shape of the crack as well as to its position means that for a given solution or gaseous environment quantitative effects could be evaluated using the diffusion equations. In principle therefore, it should be possible

to form an atlas of corrosion cracks which sets out the key shapes and depths to ensure safety and quality. This could, in principle, lead to the prohibition of manufacturing processes which produced detrimental valleys in the surface. Although the shapes of fatigue cracks could never be predicted, the valleys of manufactured surfaces could be more predictable with reference to the 'unit of manufacture'.

It could never be argued that surface geometry is the prime consideration in corrosion. It is not. However, it is important in respect of cracks, as has been seen. It could also be important in the adsorption of hydrogen. This again is larger at cracks and it may be that hydrogen build-up and the associated pressure build-up in ferrous metals rely somewhat on the surface geometry. As in the case of fretting fatigue this possibility could be explored.

## 7.7 One body with radiation (optical). The effect of roughness on the scattering of electromagnetic and other radiation

### 7.7.1 Optical scatter-general

There has already been a great deal written elsewhere and in this book (section 4.3.2) on roughness and light scatter. Indeed it is impossible to consider interferometry, holography, etc, unless this is done. Much of the theory regarding scatter has already been covered therefore. Light in effect has been used as a tool, for example in measuring flatness, roughness, etc.

There is another way of looking at the subject and this is from the functional point of view. How does roughness affect light scatter?

In the next section the basic problem will be investigated from a purely engineering point of view. Can the surface be considered to be a basic component of an optical system? Has it got a numerical aperture or its equivalent?

It has already been mentioned that light can be used to estimate the surface roughness. The idea is that the surface roughness is estimated from the pattern of light scattered from the surface. Here the light is a tool for surface measurement. On the other hand there is a large body of applications in which the scatter itself is meaningful. In fact in some instances a controlled degree of scatter from a surface is preferred, for instance in the quality of photographic prints. Paper which is too glossy is not acceptable. Then there is the very important question of asking what parameters of the surface produce specific light-scattering effects and then what manufacturing process can be used to produce these surface parameters and hence indirectly produce the desired optical effect?

In an ideal situation there is a fixed relationship between the manufacture and the function of the part. This will be considered later on. To start with, the way in which the surface modifies or modulates the wavefront will be investigated.

Consider figure 7.128. This shows the possible ways of investigating the behaviour of waves incident on a surface.

As can be seen the options are considerable and have been covered in many places. One very important review has been given by Ogilvy [290]. The applications are numerous; only a few will be touched upon here.

### 7.7.1.1 Models

The general approach to surface behaviour which is true of contact, lubrication and optics is to try to break the problem down into tractable parts.

Instead of thinking of the surface as a complex, continuous, areal geometry it is often represented as a set of unit objects such as a set of hemispheres or cones or wedges. These units can be given different scales of size or different attributes such as curvature or slope. Then, the way in which these units are distributed in space is considered, usually in the $x$ and $y$ directions but sometimes also in the $z$ direction.

**Figure 7.128** Behaviour of waves incident on a surface.

From this range of simple models solutions of behaviour are obtained by three possibilities: (i) integral methods; (ii) perturbation methods; and (iii) iterative methods. These operations are used to estimate the average functional behaviour and are listed in order of complexity. A different approach is required for extreme-value situations such as occur in corrosion and fatigue.

Breaking the problem down like this makes it tractable. It also helps considerably to identify the essential physical behaviour. Finally it is usually accurate enough to work with. Unfortunately using simple models may be unrealistic. To imagine that a surface is comprised of a sprinkling of hemispheres can quickly be dispelled by looking at one. However, it does help as a first step in understanding what is going on and can point the way to more realistic investigations. Ideally the simple approaches should be followed immediately by simple experiments and simulations. Revision and refinement should then take place based on what has been discovered.

### 7.7.2 General optical

There are three basic ways of tackling the scatter problems of light from surfaces. Each has its advantages and disadvantages which are well documented: see, for example, Beckmann and Spizzichino [291], Bass and Fuchs [292] and Ogilvy [290]; a good review has been carried out by Teague *et al* [293]. These articles are excellent from an optical point of view but they fall down somewhat on their consideration of how to achieve the optical performance by means of manufacture.

Ideally it would be preferable to use the vector theory approach, for example as carried out by Le Maystre [294, 295], but the method is very complicated for the understanding of a purely engineering function and will not be dealt with in detail here, only in passing. Also, because considerable effects occur in polarization (the vector effect) owing to the presence on the surface of non-geometric features such as phase changes, dislocations, films, etc, it becomes difficult to relate the light scatter to the pure, geometrical effect and hence to simple manufacture. These non-geometric effects may be very important in some functions but are relatively uncontrollable and difficult to take into account. Much of the work of Crock and Prod'hamme [296], for example, has indicated that even in the presence of liquid the differences in light scatter due to polarization do not often exceed 10% and so may not even be worth taking into account. There are also other effects which will only be considered later as second order, and these include secondary scattering [295]. Also, double and multiple reflection and shadowing are generally considered secondary [264, 292], especially when the incident light is near to the normal and the surface has low slopes.

In this section scalar scatter will be considered mainly because it is used as the basis for most instruments and also because it is less sensitive to extraneous effects. The intensity is a function of phase and amplitude only and not of polarization. This allows a joint approach to be made with sound waves and light which is advantageous in terms of understanding the problems involved. It is also quite straightforward to show how the wavefront is affected by the surface.

Scalar theory therefore is very much a middle path of analysis between the vector theory and ray tracing or geometrical optics. However, it should be quite clear that it is only an approximation. The approaches are shown in table 7.9.

**Table 7.9**

| | |
|---|---|
| Vector | Phase and amplitude and polarization |
| Scalar | Phase and amplitude |
| Geometrical | Position and angle only $(\lambda = O)$ |

For engineering applications such as gloss and cosmetic appearance a mixture of the scalar and the geometric approach is ideal. It also has the considerable advantage that it allows the vast amount of literature which has been generated for the design of optical instruments to be utilized (see e.g. [227]). In effect the rough surface can be considered to be a component in an optical system even though it is flawed and can be examined for aberrations in the same way that optical components are. Scalar methods are more tractable, dealing with amplitude and phase only. Even then there are difficulties because of the complicated geometry found on surfaces, for instance the surface of the sea or the moon as well as engineering surfaces.

Because surface slopes made by conventional machining processes are generally low, the effects of multiple reflections and shadowing will be considered only as secondary effects. In the analysis which follows only one dimension will be evaluated for simplicity. The treatment can generally be expanded to the areal view without loss of generality.

Let a point on the wavefront scattered from the surface be $V(x)$—the simple Kirchhoff treatment

$$V(x) = R \exp[-jk\zeta(x)]. \tag{7.239}$$

In this equation it is assumed that the surface is modulating only the phase and not the amplitude which is kept at $R$ and is invariant to both $x$ and angle. This is an approximation, but for the purpose of this analysis it will be taken as representative of what happens to a light wave incident on a surface.

The practical arrangement for viewing is shown in figure 7.129. Collimated light is thrown onto the surface. The light scattered from the surface is viewed in the far field by means of a lens $L$ of focal length $f_L$. The intensity pattern of the light scattered from the surface is $I(\omega)$ where $\omega = 2\pi D / \gamma f_L$ (the angular frequency corresponding to a spacing of $D$ on the surface).

For practical convenience the incident light is assumed to be normal. Then the surface height $z(x)$ is related to the optical path length $\zeta(x)$ by the simple relationship (figure 7.129)

$$\zeta(x) = 2z(x). \tag{7.240}$$

In both the vector and scalar theory there is difficulty in taking into account the optics of the measuring system as well as the light scatter from the surface. The only practical way to achieve this is to use ray tracing.

Ray tracing takes no account of phase; the ray is in effect conveying energy flux per unit steradian per second. The wavelength of light can sometimes be regarded for simplicity in ray tracing as being very small. The beauty of ray tracing based on geometric optics is twofold. First, very large surface effects more commonly found in objects such as aspherics can be readily dealt with in addition to the surface roughness (which is treated as if it were an aberration). Secondly, much work [259] has been carried out in the past to

**Figure 7.129**

develop mathematics to optimize the design of optical systems. It seems a pity not to take advantage of such techniques when measuring surfaces.

Unfortunately there appears to be great confusion in the literature concerning the circumstances under which the light scatter can be considered to be geometrical or not. In the case where there is a wide spatial frequency separation between the long wavelength and the short wavelength on the surface, a piecemeal approach can be adopted, the ray tracing being used for the long-wavelength geometry and scalar theory for the shorter wavelengths. However, it is the purpose of this section to consider the surface roughness in the first instance to be random and then deterministic and to examine light scatter from it using scalar theory to extract the geometrical components. Finally, some manufacturing implications of the results are considered, although this aspect was discussed in more detail in chapter 6.

Obviously the angle of scattered light in the geometrical case is twice that of the surface as shown in figures 7.129 (*a*) and *(b)* in order to obey the laws of reflection. Spacings on the surface are the same as those on the wavefront.

The autocorrelation of the wavefront is given by

$$A_{\mathrm{w}}(\tau) = \left\langle V(x_1)V^*(x_2) \right\rangle = A_{\mathrm{w}}(x_2 - x_1) \tag{7.241}$$

where $\langle\ \rangle$ indicates an ensemble average over all statistical possibilities and the surface is assumed to be statistically stationary so that

$$A_{\mathrm{w}}(\tau) = A_{\mathrm{w}}(x_2 - x_1).$$

Equation (7.241) can be used in the light scatter received at the far field:

$$I(\omega) = R^2 \iint \left\langle V(x_1) V^*(x_2) \right\rangle \exp[-jk'\omega(x_2 - x_1)]\mathrm{d}x_2\,\mathrm{d}x_1$$

$$I(\omega) = R^2 \iint A_\mathrm{w}(x_2 - x_1)\exp[-jk'\omega(x_2 - x_1)]\mathrm{d}x_2\,\mathrm{d}x_1 \tag{7.242}$$

$$= \lim_{L \to \infty} \frac{1}{LJ} \int_{-L}^{+L} (L - |\tau|A_\mathrm{w}(\tau))\exp(-jk'\omega\tau)\mathrm{d}\tau$$

$$= \int A_\mathrm{w}(\tau)\exp(-jk'\omega\tau)\mathrm{d}\tau \tag{7.243}$$

where

$$k' = \frac{2\pi}{\lambda f_\mathrm{L}}. \tag{7.244}$$

Equation (7.242) is the well-known expression which demonstrates that the scattered intensity in the far field is the Fourier transform of the autocorrelation function of the wavefront emanating from the surface. Unfortunately, it does not refer to the actual surface itself but only to the wavefront scattered from it.

To see what form $A_w(r)$ takes in terms of the surface it is necessary to make some assumptions about the surface itself. Consider that the surface has a Gaussian height distribution. If this is so the slope distribution is also Gaussian. From equation (7.240) the slopes of the wavefront will also be Gaussian because they are made of the differences between Gaussian variates. Hence

$$V(x_1)V^*(x_2) = R^2\exp[-jk\zeta(x_1)]\exp[jk\zeta(x_2)]$$

$$= R^2\exp[-jk(\zeta(x_1) - \zeta(x_2))] \tag{7.245}$$

$$= G(t),\ \text{say}$$

and where

$$t = \zeta(x_1) - \zeta(x_2). \tag{7.246}$$

Therefore

$$\left\langle V(x_1)V^*(x_2) \right\rangle = \int G(t)p(G)\,\mathrm{d}G = \int G(t)p(t)\,\mathrm{d}t$$

but

$$p(t) = \frac{1}{\sigma_\mathrm{ws}\sqrt{2\pi}}\exp\left(\frac{-t^2}{2\sigma_\mathrm{ws}^2}\right). \tag{7.247}$$

Hence

$$\left\langle V(x_1)V^*(x_2) \right\rangle = \frac{R^2}{\sigma_\mathrm{ws}\sqrt{2\pi}}\int \exp(-jkt)\exp\left(\frac{-t^2}{2\sigma_\mathrm{ws}^2}\right)\mathrm{d}t$$

$$= \frac{R^2\exp}{\sigma_\mathrm{ws}\sqrt{2\pi}}\left(\frac{\sigma_\mathrm{ws}^2 k^2}{2}\right)\int_{-\infty}^{\infty}\exp\left[-\frac{1}{2\sigma_\mathrm{ws}^2}(t + j\sigma_\mathrm{ws}k)^2\right]\mathrm{d}t \tag{7.248}$$

$$= R^2\exp\left(-\frac{\sigma_\mathrm{ws}^2 k^2}{2}\right)$$

but

$$\sigma_\mathrm{ws}^2 = E[\zeta(x_1) - \zeta(x_2)]^2 = 2\sigma_\mathrm{w}^2[1 - A_\mathrm{s}(\tau)]. \tag{7.249}$$

In equation (7.249) $A_s(\tau)$ is the autocorrelation of the surface, $\sigma_w^2 = 4\sigma_s^2$, where $\sigma_s$ is the RMS value of the surface geometry, because for the normal incidence case (7.240) holds. Equation (7.249) is the structure function of the wavefront $S_w(r)$, in this case $\sigma_{ws}^2$

Hence

$$A_w(\tau) = R^2 \exp\{-4k^2\sigma_s^2[1 - A_s(\tau)]\} \tag{7.250}$$

so it can readily be seen that, even for a Gaussian surface, the wavefront statistics are not the same as those for the surface. All too often they are taken to be the same i.e., $A_w(\tau) \neq A_s(\tau)$.

Thus, assuming a Gaussian height distribution and using the Helmholtz–Kirchhoff integral the intensity pattern in the Fourier plane is

$$I(\omega) = R^2 \int \exp\{-k^2\sigma_w^2[1 - A_s(\tau)]\}\exp(-jk'\omega\tau)d\tau. \tag{7.251}$$

When $\sigma_s$ is small compared with $\lambda$ the surface is generally called a weak scatterer (i.e. $\sigma_s/\lambda$ (< $2\pi$. For $\sigma_s/\lambda > 2\pi$ it is a strong scatterer.



**Figure 7.130** Scatter for fine surfaces.

### 7.7.3 Smooth random surface

In this case the exponent can be approximated by its first two terms. Thus

$$I(\omega) = R^2 \int (1 - k^2\sigma_w^2)\exp(-jk'\omega\tau)d\tau + R^2 \int k^2\sigma_w^2 A_s(\tau)\exp(-jk'\omega\tau)d\tau \tag{7.252}$$

or

$$I(\omega =) \text{ aperture term } + C \times \text{ power spectrum of surface}$$

where $C$ is a constant independent of scattering angle. Equation (7.252) produces an intensity pattern as shown in figure 7.130. For the case of a smooth random surface therefore, the intensity pattern has two quite distinct components.

The aperture term contains no spacing terms related to the surface, only the optical configuration. Hence the aperture term is

$$R^2(1 - k^2\sigma_w^2)2\,\text{sinc}(k'\omega L). \tag{7.253}$$

This term always appears in scalar scattering but not in the geometric case.

To illustrate how the autocorrelation function of the surface affects the scatter, consider the case of $A(\tau) = \exp(-\tau^2/2\sigma_L^2)$ a Gaussian function having a correlation length $\sigma_L$. Ignoring the aperture term gives

$$I(\omega) = R^2 k^2 \sigma_w^2 \sqrt{2\pi} \exp\left[-\frac{1}{2}\left(\frac{\omega/f}{(\lambda/\sigma_L)/(1/2\pi)}\right)^2\right]$$ (7.254)

or

$$I(\omega) = C \exp\left(-\frac{1}{2}\frac{\alpha^2}{(\lambda/\sigma_L)^2(1/2\pi)^2}\right).$$ (7.255)

In equation (7.254) the diffraction angle is determined by the ratio $\lambda/\sigma_L$. The heights of the surface are not important except that they have to be small; here $\alpha = \omega/f$.

There is a small slope $I(\omega)$ brought about by the dependence between the two components of the spectrum. Thus

$$\int \omega^2 P(\omega)\,\mathrm{d}\omega = \sigma_w^2$$ (7.256)

where $\sigma_w^2$ is the variance of the wavefront slope. This is completely general for any type of surface.

Equation (7.252) relies for physical sense on the fact that the surface and hence the wavefront height standard deviation is small compared with the wavelength of light. It has been shown that the approximation to $P(\omega)$ is good up to $\sigma_s < \lambda/6$ for periodic surfaces and slightly better for random surfaces.

Up to now this treatment has used the scalar theory of physical optics. How does this relate to geometric optics ray tracing?

### 7.7.4 Geometric ray-tracing criterion—rough random surfaces

If the scattered light has a probability density directly proportional to the standard deviation of surface slopes then this constitutes the geometric condition for light reflection. This condition is given as

$$I(\omega) \propto p(\tan\alpha)$$ (7.257)

A similar condition could be applied to curvature and divergent light scatter.

The condition given by expression (7.257) is quite restrictive. It means that there is geometrical reflection scatter if there is a one-to-one relationship between the scattered light intensity and the slope distribution on the surface (and hence wavefront). It does not mean that the geometric condition is satisfied if one characteristic of the slope distribution of the surface is equal to one characteristic of the light intensity scatter. Just because the second moment of the scatter and the slope distribution yield $\sigma_w^2$ as in equation (7.253) does not mean that there is a geometric condition. In fact, because a Gaussian slope distribution is necessarily smooth and continuous it can never be equivalent to the intensity scatter given for smooth surfaces by equation (7.252) which is of a degenerate form.

From figure 7.129 the condition in equation (7.257) can be seen as $\alpha = \omega/f$. If the surface has a Gaussian distribution of slopes (and curvatures, both of which follow from Gaussian multivariate theory) the question arises as to whether equation (7.251) can be written in the form of (7.257).

Consider equation (7.251):

$$I(\omega) = R^2 \int \exp\{-k^2\sigma_w^2[1 - A_s(\tau)]\}\exp(-\mathrm{j}k'\omega\tau)d\tau.$$ (7.258)

This equation can only yield a Gaussian term outside the integral sign, thereby representing the geometric slope distribution, by being able to complete squares in the exponential bracket in a way which does not contain $\tau$ and yet does involve $\sigma_\alpha$. This is equivalent to saying that for Gaussian statistics $A_w$ and $p(\alpha)$ are transform pairs.

Expanding $A_s(\tau)$ by Taylor's theorem in even powers gives

$$A_s(\tau) = A(0)\left(1 + \frac{\tau^2}{2} A''(0)/A(0) + \dots\right) \qquad (7.259)$$

where $A(0) = \sigma_w^2$

Expanding $A_s(\tau)$ in even powers is necessary because $A_s(\tau)$ is an even function. For the surface to be well behaved and not fractal the autocorrelation $A_s(\tau)$ should also have well-defined differentials at the origin; that is, $A''(0)$, $A^{iv}(0)$, etc, should exist and be finite. Restricting equation (7.259) to two terms only, to avoid introducing terms in $\tau^4$ gives $A_s(\tau)$ truncated to $A_{st}(\tau)$:

$$A_{st}(\tau) = A(0)\left(1 + \frac{\tau^2}{2} \frac{A''(0)}{A(0)}\right). \qquad (7.260)$$

This equation is equivalent to approximating the autocorrelation function by a sagittal drop term $\tau^2 C/2$ from $A(0)$, the autocorrelation function at the origin where $C = A''(0)$.

This expression in equation (7.260) does not necessarily imply an approximation to a Gaussian auto-correlation function $A_s(\tau)$ or a Lorentzian function or even a cosine function. Another form based on Gaussian statistics is allowed for in equation (7.260) by letting $A''(0)/A(0) = -\sigma_s^2/\sigma_w^2$ where $\sigma_\alpha^2$ is the variance of surface slope and $\sigma_w^2$ is the variance of heights of the wavefront.

Hence

$$A_{st} = \sigma^2 \left[1 - \left(\frac{\tau^2 \sigma_\alpha^2}{2\sigma^2}\right)\right] \qquad (7.261)$$

and thus

$$I(\omega) = R^2 \int \exp\left[-\frac{k^2\sigma^2}{2}\left(\frac{\sigma_\alpha^2}{\sigma^2}\right)\left(\tau + j\frac{k'\omega\sigma}{\sigma_\alpha}\right)^2\right]d\tau \exp\left[-\left(\frac{k^2\omega^2}{2\sigma^2}\frac{1}{k^2}\right)\right] \qquad (7.262)$$

from which

$$I(\omega) = \frac{R^2\lambda}{\sigma_\alpha\sqrt{2\pi}} \exp\left(-\frac{(\omega/f)^2}{2\sigma_\alpha^2}\right). \qquad (7.263)$$

$I(\omega)$ represents the desired Gaussian condition if $\omega/f = \tan\alpha \approx \alpha$; it becomes the Fourier transform pair of $A_w(\tau)$. Hence

$$I(\omega) \propto (\tan\alpha) = \frac{R^2\lambda}{\sigma_\alpha\sqrt{2\pi}} \exp\left(-\frac{\alpha^2}{2\sigma_\alpha^2}\right). \qquad (7.264)$$

Therefore, for the surface autocorrelation function given by equation (2.261) the surface behaves as a geometric reflector on a point-by-point basis.

Equation (7.261) is obtained only by making a suitable choice for the surface autocorrelation function to make the autocorrelation of the wavefront Gaussian. It should be pointed out that (7.261) is in a divergent form. It therefore violates the definition of an autocorrelation function, that is $|A(\tau)|(r) \to \infty$ as $\tau \to \infty$.

Another obvious point is that the phase terms all vanish in the integration of equation (7.262) thereby allowing a geometric rather than a scalar solution. Also, if the surface has an exact Gaussian autocorrelation function then the geometric condition cannot be obeyed.

### 7.7.4.1 Effect of surface curvature

The surface has been considered to be made up of random tilted facets of undetermined length. This type of profile produces a wavefront having $z$ heights and a complex field given by $V(x)$.

In practice the situation is more complicated because of the presence of many other surface features. One such complication is the presence of curvature on the facets rather than flats. This modifies the angles at which light rays are scattered, as seen in figure 7.131.

**Figure 7.131** Effect of curvature on scatter angle.

Curvature has the effect of reducing the light scattered in the specular direction. Hence $p(\alpha)$ is reduced for the same input and the effective angle of scatter becomes $\alpha$, not $f'(x)$, due to local curvature:

$$\alpha = f'(x) + xf''(x). \tag{7.265}$$

In the figure $f$ is the average size of facet, and $x$ is taken from the mid-point of the facet.

The variance of rays from the surface will obviously increase as a result of the term $xf''(x)$. If $\sigma'_\alpha$ denotes the standard deviation of the reflected ray at $\alpha$ then

$$\sigma'^2_\alpha = E[f'(x) + xf''(x)]^2 = \sigma^2_{sc} + E[x''f]^2 = \sigma^2_{sc} + E[x^2]E[f''^2] = \sigma^2_{sc} + \frac{l^2}{12}\sigma^2_{sc} \tag{7.266}$$

where it is assumed that $f'(x)$ is independent of $f''(x)$ and $x$ is independent of either. Furthermore it is assumed that $p(x)$ is uniform from $l/2$ to $-l/2$. The variances $\sigma^2_\alpha$ and $\sigma^2_c$ are those of profile slope without curvature and with curvature respectively, $\sigma^2_{sc}$ is the remnant of the variance in the specular direction.

From equation (7.266) it is seem that the effective value of scatter has been increased by $(l^2/12)\,\sigma^2_c$ to take account of curvature. Hence $\sigma_{sc} = \sigma_\alpha - (l/12)\sigma_c$, where $\sigma_{sc}$ replaces $\sigma_\alpha$.

The geometric condition taking account of both slopes and curvature is therefore given in this treatment by a modification of equation (7.264) to take account of (7.266). Thus

$$I(\omega) = \frac{R^2\lambda}{\sigma_{sc}\sqrt{2\pi}}\exp\left(-\frac{\alpha^2}{2\sigma^2_{sc}}\right) \tag{7.267}$$

or

$$I(\omega) = \frac{R^2\lambda}{\sqrt{2\pi}[\sigma^2_\alpha - (l^2/12)\sigma^2_c]^{1/2}}\exp\left(-\frac{(\omega/f)^2}{2[\sigma^2_\alpha - (l^2/12)\sigma^2_c]}\right) \tag{7.268}$$

where $\sigma_{sc}$ is the standard deviation of the local specular reflection with curvature taken into account.

When this extra requirement is transmitted back to equation (7.264) it means that the autocorrelation function of the surface, $A_{st}(\tau)$, has to take a modified form of

$$A_{st}(\tau) = \sigma^2\left[1 - \frac{\tau^2}{2\sigma^2}\left(\sigma^2_\alpha - \frac{l^2}{12}\sigma^2_c\right)\right]. \tag{7.269}$$

Note that because this is a simplified treatment it does not require an extra term including $\tau^4$ in the expansion to cater for curvature, just an additional term in $\tau^2$—which is not surprising because the curvature term here simply modifies the local facet slope.

Equation (7.269) is useful because it can be compared with the case when the correlation function is allowed to take an extra term. Thus $A_{st}(\tau)$ can be taken to an extra term in $\tau^4$ giving

$$A_{st}(\tau) = \sigma^2\left(1 + \frac{\tau^2}{2\sigma^2}A''(0) + \frac{\tau^4 A^{iv}(0)}{24\sigma^2}\right). \tag{7.270}$$

How does equation (7.270) compare with (7.269), which takes account of curvature? Using random process analysis $A''(0) = -\sigma_\alpha^2/\sigma^2$ and $A^{iv}(0) = -\sigma_c^2/\sigma_\alpha^2$, yielding

$$A_{st}(\tau) = \sigma^2\left[1 - \frac{\tau^2}{2\sigma^2}\left(\sigma_\alpha^2 - \frac{\tau^2}{12}\sigma_c^2\right)\right]. \tag{7.271}$$

This shows the interesting result that, if $\tau$ is taken to be equal to $l$, then equation (7.269) is identical to (7.271). The question now arises as to the value of $l$.

### 7.7.4.2   Estimation of I, the facet length

There are a number of ways of evaluating this. The most obvious, yet not the easiest, way is to resort to random process analysis [297,298].



**Figure 7.132** Inflection points on waveform.

It could be argued that the end points of facets are marked by points of inflection as shown in figure (7.132) and the average length of a facet is therefore the reciprocal of the mean distance between inflection points. Thus if $\bar{c}$, $\bar{m}$ and $\bar{n}$ represent the average distance between crossings, peaks and inflection points of a random waveform respectively (in this instance in the one-dimensional case), it works out that

$$\bar{n} = l = \pi[A(0)(\pi\bar{c})^2(2\pi\bar{m})^2]^{1/2}. \tag{7.272}$$

Unfortunately, this definition of $l$ presupposes well-behaved differentials of $A_s(\tau)$ up to and including the sixth differential. This is not possible to assure because of the very restricted series allowed for $A_s(\tau)$ in equation (7.270).

### 7.7.5   Scatter from deterministic surfaces

As before, the criterion for scatter is determined from the wavefront and in particular the expected or average phase difference produced by the surface.

Hence if the surface is represented by a sine wave

$$z(x) = \tilde{R}\sin x$$

$$\begin{aligned}
E[V_1, V_2] &= \frac{1}{2\pi}\int \exp\{-jk\tilde{R}\,[\sin x - \sin(x + \tau)]\}\,dx \\
&= I_0(j2k\tilde{R}\,\sin\tau/2) \quad or \quad J_0(-2k\tilde{R}\,\sin\tau/2) \\
&= J_0(2k\tilde{R}\,\sin\tau/2)
\end{aligned} \tag{7.273}$$

because $J_0(-)$ is even. Hence the intensity becomes

$$I(\omega) = R^2 \int_{\infty}^{-\infty} J_0(2\widetilde{R}\sin\tau/2)\exp(-jk'\omega\tau)\ \mathrm{d}\tau. \tag{7.274a}$$

If the wavelength of the signal is $\lambda_L$, equation (7.274$a$) can be solved to give the general solution

$$I(2\gamma) = \lambda_L\cos(\pi\gamma)J_\gamma^2(2\pi\widetilde{R}/\lambda)(-1)^\gamma \tag{7.274b}$$

where $(\lambda_L/\lambda)\ \alpha = \gamma$ and $\alpha = \omega/f$.

### 7.7.6   Smooth deterministic signal

Consider equation (7.274) when $\widetilde{R}$ is small. Under these circumstances the argument for the Bessel function of the first kind, zero order, is itself small. Therefore it can be approximated closely by the first two terms of its ascending power series:

$$\begin{aligned} J_0(2\widetilde{R}\ k\ \sin\tau/2) &\approx 1 - 4\widetilde{R}^2\ k^2\sin^2\tau/2 \\ &= \left(1 - \frac{k^2\widetilde{R}^2}{2}\right) + \frac{k^2\widetilde{R}^2}{2}\cos\left(\frac{2\pi}{\lambda_L}\tau\right). \end{aligned} \tag{7.275}$$

Hence

$$I(\omega) = R^2\int\left(1 - \frac{k^2\widetilde{R}^2}{2}\right)\exp(-jk\omega\tau)\ \mathrm{d}\tau + R^2\int\frac{k^2\widetilde{R}^2}{2}\cos\left(\frac{2\pi\tau}{\lambda_L}\right)\exp(-jk'\omega\tau)\ \mathrm{d}\tau \tag{7.276}$$

$$I(\omega) = R^2\left(1 - \frac{k^2\widetilde{R}^2}{2}\right)\left[2\ \mathrm{sinc}\left(\frac{k'\omega L}{2}\right)\right] + R^2\ \frac{k^2\widetilde{R}^2}{2}\int\cos\left(\frac{2\pi\tau}{\lambda_L}\right)\exp(-jk'\omega\tau)\ \mathrm{d}\tau. \tag{7.277}$$

Compare equation (7.277) with (7.252). The aperture term in both cases is a function of the optical system and the mean square signal.

### 7.7.7   Geometric ray condition, rough deterministic surfaces

The general case is given by equation (7.277) subject only to the assumptions made in this derivation. In equation (7.274) the actual size of the surface is determined by the Bessel argument $R/\lambda$.

The geometric condition can be obtained in exactly the same way for deterministic surfaces as for the random case. In this the argument of the wavefront autocorrelation function $A_w\ (\tau)$ is limited in its power series expansion. Thus $2\widetilde{R}k\ \sin\tau/2$ becomes simply $k\widetilde{R}\tau$ which reflects the geometric ray's dependence on linear facets. Thus

$$J_0(2k\widetilde{R}\sin\tau/2) \simeq J_0(k\widetilde{R}\tau) \tag{7.278}$$

and equation (7.274$b$) becomes $I(\omega) = I(2\gamma)$ where

$$I(2\gamma) = \gamma_L\left(\frac{1}{\sqrt{\alpha_{\max} - \alpha^2}}\right). \tag{7.279}$$

In this equation $\alpha_{\max}$ is the maximum slope for the waveform, which for a sine wave is $R\omega_0$. The general slope $\alpha$ is given by $\alpha = \omega_0 z$.

Equation (7.279) represents exactly the slope distribution to be obtained from a deterministic sinusoidal component. Hence the intensity distribution is the slope distribution of the surface so that the geometric reflection criterion is achieved.

### 7.7.8 *Summary of results, scalar and geometrical*

A number of points concerning symmetry can be made. These are shown in tabular form in table 7.10.
From this table some particular comments can be made:

1. If the surface is rough and random and the autocorrelation function of the surface takes a specified form which is not quite Gaussian or Lorentzian, the surface can behave as a scattering object according to geometric laws.
2. Surface curvature and slopes can both be taken into account using the specified form for $A_s(\tau)$.
3. If the surface is fine, despite having a Gaussian autocovariance for the wavefront, it will not behave as a geometric scatterer, only as a diffraction scatterer because the intensity distribution is degenerate—it produces separate aperture and diffraction terms.
4. Features of slope, height, curvature and spatial frequency can all be components of light scattered from a surface. Furthermore the geometric effects and diffraction effects can coexist if a rough component and fine component occur together on the surface.
5. Surfaces having autocorrelation functions which are not well behaved in terms of their differentials cannot be made to take on geometric properties. This category includes surfaces such as fractals, which are based on Poissonian rather than Gaussian statistics.

$$I(\omega) = \int \left\langle V_1 V_2^* \right\rangle \exp(-jk'\omega\tau)d\tau \tag{7.243}$$

**Table 7.10**

In conclusion, the amplitude of the surface determines the mode of scatter whether physical or geometric. For the former the spacings on the surface determine the shape of the intensity distribution; for the latter the slopes on the surface determine the shape of the intensity distribution.

This general conclusion is meant to be a guide to the way in which surfaces react to light. As such it is conditional on the assumptions made for scalar scattering theory. It is, however, useful in the sense that it enables some characterization of surfaces by their light-scattering properties. Furthermore, it pinpoints those geometric parameters of the surface which are important in scatter and the way in which they are important.

One obvious conclusion is that the parameters are not dominated by peak (summit) properties as they tend to be in two-body interactions. The parameters take in the whole of the surface.

Extreme properties tend to be important only in secondary effects such as multiple reflections and shadowing which will be briefly discussed in section 7.7.10.

### 7.7.9 Mixtures of two random components

$$z_T(x) = z_1(x) + z_2(x). \tag{7.280}$$

It is assumed that the two components are independent and that one component, say $z_1(x)$, representing the roughness, is wideband yet of high spatial frequency and the other component $z_2(x)$ is narrow band and of low spatial frequency. This $z_2(x)$ is what is in fact typical of waviness (figure 7.133).



**Figure 7.133** Scatter from surface with different scales of size.

Using the same procedure as before and assuming normal statistics

$$A_w(\tau) = \langle V_1^*(x_1)V_2^*(x_2)\rangle = \exp[-k^2\sigma_1^2 + \sigma_2^2 - \sigma_1^2 A_1(\tau) - \sigma_2^2 A_1(\tau)] \tag{7.281}$$

If the two components are smooth, $\sigma_1, \sigma_2 \ll \lambda$. Using the same reduction in the exponential series

$$I(\omega) = K_1 P_1(\omega) + K_1(\omega) \tag{7.282}$$

ignoring the aperture term. The spectra are simply additive.

If one component is large and the other small, $\sigma_1 < \lambda$, $\sigma_2 > \lambda$, then

$$A_w(\tau) = \exp\{-k^2\sigma_1^2[1 - A_1(\tau)]\}\exp\{-k^2\sigma_2^2[1 - A_2(\tau)]\}$$

$$= [(1 - k^2\sigma_1^2) + k^2\sigma_1^2 A(\tau)]\exp\left(-\frac{k^2\sigma_\alpha^2\tau^2}{2}\right)$$

where $\sigma_\alpha$ is the standard deviation of slope of the rough component, giving

$$I(\omega) = K_1(\sigma_1)\exp\left(-\frac{\alpha^2}{2\alpha_\alpha^2}\right) + k_2(\sigma_1\sigma_L)\exp\left(-\frac{\alpha^2}{2\alpha_{\alpha L}^2}\right) \tag{7.283}$$

where

$$\alpha_{\alpha L}^2 = \left[ \sigma_\alpha^2 + \left( \frac{\lambda}{\sigma_L} \right)^2 \left( \frac{1}{2\pi} \right)^2 \right] \qquad (7.284)$$

$$= k_1(\sigma_1) p(\alpha) + k_2(\sigma_1 \sigma_L) p(\alpha') \qquad (7.285)$$

where

$$\alpha' = \frac{\alpha}{[1 + (\lambda^2/\sigma_L \sigma_\alpha)^2 (1/2\pi)^2]^{1/2}}. \qquad (7.286)$$

Hence the scatter term of purely geometric nature relating to the large component is retained and in addition there is a mixed term which incorporates both spacing and angle components. The diffusion therefore is increased with the presence of the small component, which is to be expected.

In other words, the fine surface acts as a diffraction scatter on top of the geometric scatter, as seen in figure 7.133 so that the geometrical and physical scattering modes can exist side by side. It may be that this explains why on rough surfaces the scatter of light is sometimes much. wider than expected due to the term $\lambda/\sigma_L$ in $\sigma_{\alpha L}$. Another minor point is that the aperture term is small because $\sigma_1 > \lambda$.

The same breakdown holds for when one or both signals are deterministic.

One requirement in surface metrology is to obtain quantitative data about the geometry as quickly as possible. This can be done using fast scanning tactile methods but still not quickly enough for in-process measurement. Optical methods can be used as is seen in chapter 4 but often the optical method is more responsive to surface slope than to the actual height. Nomarsky and light scattering methods come into this category. It has been suggested that it should be possible to get at least a good estimate of the surface profile, say, given the values of the slope by invoking Taylor's expansion. Because the slope information is noisy by virtue of its high frequency amplification a low pass filter has to be incorporated [299]. In this method it is not adequate to derive all derivatives via the next lower one i.e. $f'$ in terms of $f''$ and $f'''$ in terms of $f'$ etc. because all estimates are ultimately dependent on $f'$ which means that the noise gets embedded in the signal.

Any method of reconstruction via the differentials requires some assumptions to be made [124]. The noise is one problem and somewhere there is an arbitrary constant of integration which has to be calibrated out. Consequently any method can be made to work over a limited range in amplitude and for a restricted waveform but not generally. So these methods are confined usually to specific processes such as diamond turning. The technique by Mansfield is particularly intriguing.

### 7.7.10   Other considerations on light scatter

Many excellent books are now available on the subject which expand considerably on the simple light reflection considerations discussed above [290].

Some bring in second-order effects which occur when the source and detector are perhaps in difficult positions. An example is when the incident light is at an extreme angle. Under these conditions the following can occur:

(1)  the wavelength of the light is effectively increased due to the non-normal angle;
(2)  multiple scattering;
(3)  shadowing of the surface.

These will be considered in what follows.

### 7.7.10.1   Angle effects

In figure 7.134 the path difference is $2(z_1 - z_2)$, whereas from the geometry in case (*b*) it is $2\sin\beta(z_1 - z_2)$.

Hence the phase difference is

$$(a) \quad \frac{4\pi}{\lambda}(z_1 - z_2)$$

$$(b) \quad \frac{4\pi(z_1 - z_2)}{\lambda/\sin\beta}.$$

(7.287)



**Figure 7.134** Angular (or obliquity) effect on scatter.

Consequently the effective wavelength of the incident light is $\lambda/\sin\beta$ rather than $\lambda$. From this it can be seen that diffraction or wide-angled scatter is more likely than even the glancing angle of the incident light would suggest. This oblique-angled illumination is sometimes used to increase the effective wavelength. Apart from the disadvantage that it foreshortens detail, it has the advantage that it enables rough surfaces to be measured using the effective longer wavelength and still satisfies diffraction criteria without using quartz lenses and special infrared detectors.

The angles generally used in the literature are shown in figure 7.135. In this case equation (7.287) becomes

$$\frac{4\pi(z1 - z_2)}{\lambda/(\cos\theta_1 + \cos\theta_2)}$$

(7.288)



**Figure 7.135** Angle convention.

and all equations involving $2\pi/\lambda$ become $2\pi(\cos\theta_1 + \cos\theta_2)\lambda$. The reason for this difference in angle nomenclature is because surface angles are conventionally measured from the horizontal, whereas in optics they are measured from the normal—yet another source of confusion!

The viewing of surfaces is usually done by holding the surface at right angles to the line of sight of the eye and also at various angles. In getting an idea of roughness a picture is built up by the succession of images obtained at various angles. Because the surface appearance changes with angle, such angular considerations are important. What parameters of roughness influence the angular picture? Obviously the local slope of the surface and 'lay,' corresponding in the simplest terms to 'altitude' and 'azimuth' effects.

### 7.7.10.2   Multiple reflections

In questions of the functional importance of surfaces in optical (or other) scatter the issues are basically as follows. What factors of the geometry are important in determining the properties? What property of the geometry influences gloss, lustre, contrast, directionality, etc? Often the description of the property is subjective. The reason is because these properties have been described by people who only need to identify quality in general terms and not numerically. Up to now the scatter has been restricted to single reflection. Obviously multiple scattering is possible and so also is shadowing if the light is incident at an acute glancing angle. The complexity of the analysis for such cases is quite high and yet the importance functionally is low. In real situations it is necessary to apply Pareto's rule concerning priorities. Concentrate on what is the most important issue rather than muddy the picture by less important considerations. Multiple scattering and shadowing fall into this latter category.

The same argument follows for the use of scalar and vector methods in determining scatter. While the latter are undoubtedly more accurate there are considerably more complicated and often more restrictive. Whether or not this accuracy is worth the effort is questionable for engineering considerations. For a thorough comparison of methods see Ogilvy [290].

In what has been considered so far, light is directed at the surface and scattered back. It is not contained within the surface asperities except from the limited effect of the reflectivity by the surface. In multiple scattering, however, the light is trapped for a time within the geometry. It interacts with more than one point of the surface before finally being released (figure 7.136).



**Figure 7.136**  Multiple scatter of beam.

This type of behaviour is not probable for surfaces which have been generated by typical machined surfaces because the surface slopes are not very great. However, there is a distinct possibility that on nanoroughness and quasifractal types of surfaces the probability is higher because of the higher slopes. When slopes are of the order of 45° multiple reflections are very possible. Unfortunately, under these circumstances the slopes cannot be measured by conventional surface instruments, which often have an angle cut-off of 45° (e.g. 90° stylus) and so will not register. As a result, multiple scattering might happen where it has been dismissed because of the low slopes measured. The order of scattering is shown in figure 7.137. Notice that figure 7.136 and 7.137 have been drawn to exaggerate the problem.



**Figure 7.137**  Complicated direction of multiple scattered rays.

Earlier attempts to estimate the probability of multiple reflections for acoustic waves [299] and for electromagnetic waves [300] amongst others have been made. Some use the concept of an array of proturberances or 'bosses' scattered on an otherwise flat surface. These bosses have defined shapes such as hemispheres or even non-spherical shapes (figure 7.139). Rayleigh scattering or the treatment of the rough surface as an array of dipole sources [301] has been used for sound waves but the principles are exactly the same (figure 7.138).



**Figure 7.138** Distribution of unit scatterers.



**Figure 7.139** Rayleigh scattering with 'bosses'.

It is interesting to note that theories of contact developed with 'bosses' on surfaces in a parallel way at the same time.

Boss methods have only been partially useful [290] and then only at high angles of incidence.

Several people have used Kirchhoff theory to provide a scattering coefficient for single scattering and then to sum this coefficient over one or more further surface interactions. Lynch and Wagner [302] have considered second-order multiple scattering by this method and in addition have included geometric optics by taking $\lambda \to 0$. This leads to the probability of a ray intersecting a surface twice or more. Because of this geometric correlations between the scattering points are ignored.

In general it is found, and this is intuitively obvious, that the presence of multiple scattering increases the scatter in all directions.

Very often in fact in the case of multiple scattering the theory is complicated enough to justify consideration of only specular points, as seen via geometric optics; the diffuse terms would tend to be neglected. In fact in many attempts to predict the scattering a surprisingly simplistic model of the surface is used as a basis, from which very sophisticated mathematics is developed. The question arises: are the mathematics correct and the physical model wrong?

To see how the simple approach is tackled consider a situation as drawn in figure 7.138. If a plane wave is incident on the surface, usually taken to be perfectly conductive, and the scalar scattering is considered, the intensity $U(P)$ at $P$ can be considered to be

$$U(P) = \exp(-jkz) + \exp(jkz) + \sum_{i=\infty}^{\infty} U_i(P) \tag{7.289}$$

where $k$ is the usual wavenumber. ($z$ is the coordinate in the direction of the incident wave (and opposite sense) in reference [301]). The first term describes the incident wave, the second term scatter from an ideally smooth reflecting plane at $z = 0$. $U_i(P)$ describes the contribution of the $i$th hemisphere. Using this approach therefore the problem has been split up into two distinct parts:

1. The behaviour of the 'unit scatterer' (i.e. the hemispheres).
2. The distribution of the unit scatterers in the horizontal plane.

What many investigators forget is the distribution of scatterers in the vertical plane, although it turns out that this is not quite as important here as might be expected. It should be noted that there is a close similarity between the model here and the model of a manufactured surface involving unit manufactured events distributed randomly on a plane.

The usual assumptions are:

1. All scatterers behave the same.
2. The field around any scatterer is centrosymmetrical about a normal at $O_i$, that is independent of the local $\varphi_i$.

Hence $U_i(P)$ is given by centrosymmetrical functions centered at the local $O_i$. Thus

$$U_1(P) = \sum_{n=\infty}^{\infty} a_n h_n^{(1)}(kr_i) P_n(\cos\theta_i) \tag{7.290}$$

where $h_n^{(1)}$ denotes the spherical Hankel function of order $n$. $P_n$ is the Legendre polynomial of order $n$ and $r_i$, $\theta_i$, $\varphi_i$ are the spherical coordinates of $P$ in the 'local' reference system centred at $O_i$ as in figure 7.138

This field at $O_i$, can be looked at from the origin of the system $O_0$ and can be considered to be made up of a superposition of spherical waves of orders $v$ and $\mu$ centred at $O$. Thus

$$h_n^{(1)}(kr_i) P_n(\cos\theta_i) = \sum_{v=0}^{\infty} \sum_{\mu=-v}^{\infty} F_{v,\mu}^{(n,i)}(kr) P_v^{\mu}(\cos\theta)\exp(j\mu\varphi). \tag{7.291}$$

where $P_v^{\mu}$ denote Legendre fuctions of orders $v,\, \mu$. If the distance $O_0 - O_i \,\rho_i$, then

$$F_{v,\mu}^{(n,i)}(kr) = g_{v\mu}^{(n)}(k\rho_i)\, j_v(kr) \qquad \text{for } r < \rho_i$$

$$F_{v,\mu}^{(n,i)}(kr) = f_{v\mu}^{(n)}(k\rho_i)\, h_v^{(i)}(kr) \qquad \text{for } r > \rho_i \tag{7.292}$$

where $j_v$, denotes the spherical Bessel function of order $v$. Equation (7.293) has to be put in (7.292) and then (7.294).

So it is seen that even with this very simple start the actual calculation is becoming complex. Some approximation can be made. The way in which this type of analysis works is to consider the field at $P$ due to the sum of the spherical waves from the individual 'bosses'. This corresponds to the case for first-order scattering. For higher-order scattering the interaction between the hemispheres has to be taken into account.

Carrying this calculation further [260] an iterative method enables estimates of any order of scattering

to be found. Thus

$$U_i(P) = \sum_{n=0}^{\infty} a_n W_n^i(P) \qquad (1 \neq 0)$$

where

(7.293)

$$W_n^i(P) = \sum_{v=\infty}^{\infty} F_{v,0}^{(n,i)}(kr) P_v(\cos\theta).$$

The boundary conditions on a metal surface are such that the normal derivative of the total field $U(P)$ vanishes. This implies that odd terms are zero.

After some calculation the following is the result:

$$a_n h_n^{(1)\prime}(kR) = -2(2n+1)(-1)^{n/2} j_n'(kR) - j'(kR) \sum_{v=0}^{\infty} A_{vn} a_v$$

(7.294)

where

$$A_{vn} = \sum_{i=1}^{\infty} g_n^{(v)}(k\rho_i) - N \int_{2R}^{\infty} g_n^{(v)}(k\rho_i) \, P(\rho_i) \; \mathrm{d}\rho_i$$

(7.295)

where $N$ denotes the number of hemispheres per unit surface and $P(\rho_i)$ denotes the probability density of occurrence of the centre of a hemisphere at distance $\rho_i$ from the origin O.

Both equations (7.297) and (7.296) are handy for evaluating the multiple scattering, that is the field scattered by any hemisphere illuminated by the field of $j$th order scattered by all other hemispheres.

Starting with first-order scattering, $a_n$ is evaluated by neglecting all interactions between hemispheres, that is $A_{vn} = 0$ in the right-hand side of (7.296). Hence

$$a_n = a_n^{(1)} = -2(2n+1)(-1)^{n/2} S_n(kR)$$

(7.296)

where

$$S_n(kR) = j_n^1(kR) h_n^{(1)\prime}(kR).$$

(7.297)

For the second order $a_n = a_n^{(1)} + a_n^{(2)}$ and $a_v$ is replaced with $a_v^{(1)}$ on the right-hand side of (7.296). Hence

$$a_n^{(2)} = -S_n(kR) \sum_{v=0}^{\infty} A_{vn} a_v^{(1)}.$$

In general by putting

(7.298)

$$a_n = a_n^{(1)} + a_n^{(2)} + a_n^{(3)} + \dots + a_n^{(j)}$$

for $j \geqslant 1$ the scattering number is

$$a_n^{(j+1)} = -S_n(kR) \sum_{v=0}^{\infty} A_{vn} a_v^{(j)}.$$

(7.299)

This equation is an iterative method for getting an idea of the extent of multiple scattering. It is meant to illustrate the way in which the far field at a point $P$ changes as interactions are allowed between the fields emanating from surrounding hemispheres. It shows that the field converges for a relatively small amount of multiple scatterings so that high-order multiple scattering is rare. The calculations, however, are very messy and the question arises as to whether it is actually justified in terms of the reality of the physical model. The only real benefit is that it does allow a consideration to be given to correlation between scatterers. Obviously for first-order scattering the correlation is zero and no interactions are allowed, but as the higher order of scattering is allowed the correlation coefficients between scatters rise. In other words,

instead of a scatterer being regarded simply as a 'specular point or bright point' for single scattering, it becomes more of a 'diffuse point' as multiscattering is allowed.

### 7.7.10.3 Shadowing

When light is incident on a surface it is usually ensured that all the surface is illuminated if the incident light is uniform. This is not so because shadowing can take place (figure 7.140). Parts of the surface will be illuminated and parts will not. Also parts will be in semi-shadow.



**Figure 7.140** Shadowing and multiple hits.

Early attempts to predict shadowing assumed that the scattered light could be explained by a shadow function $S(\theta)$:

$$\langle I_{\mathrm{sh}} \rangle = \langle I \rangle S(\theta_1) \tag{7.300}$$

where $0 \leq S(\theta_1) \leq 1$. One point that is dear is that shadowing is fundamentally a slope condition and not primarily a function of surface height. Also it is clear from figure 7.140 that shadowing and multiple scattering are related effects; multiple scattering has to be taken into account when shadowing is being considered. The problem is simply put in reference [290].

Just a look at a scatter picture like figure 7.140 shows how complicated the situation is. Sometimes there is a direct reflection from the source via the surface to the receiver (1). This causes a shadow. Ray (2) reflects direct to the receiver, ray (3) hits the surface twice, whereas ray (4) hits the surface twice, once when in shadow.

Thus the scattering average for field intensity becomes

$$\langle I_{\mathrm{sh}} \rangle = \langle I_0 P_1(K_{\mathrm{inc}} \mid \alpha, \beta) \rangle \tag{7.301}$$

where $\langle I_0 \rangle$ is the unshadowed result and $P_1$ is simply the probability of an arbitrary surface point of 2D gradients $\alpha, \beta$ being illuminated when the incident ray is along the vector $K_{\mathrm{inc}}$. For bistatic scattering, that is when multiple (i.e. two) reflections can take place, equation (7.303) is replaced by

$$\langle I_{\mathrm{sh}} \rangle = \langle I_0 P_2(K_{\mathrm{inc}}, K_{\mathrm{sc}} \mid \alpha, \beta) \rangle \tag{7.302}$$

where $P_2(K_{\mathrm{inc}}, K_{\mathrm{sc}} | \alpha, \beta)$ is the conditional probability that an arbitrary surface point at $x$ and $y$, gradients $\alpha$ and $\beta$, is illuminated and directly in sight of the receiver when the incident wave direction is along $K_{\mathrm{inc}}$ and the scattered wave is along $K_{\mathrm{sc}}$. If it is assumed that the averages of the scattered intensity and the shadowing probabilities are independent, then

$$\begin{aligned} \langle I_0 P_1 \rangle &= \langle I_0 \rangle \langle P_1 \rangle \\ \langle I_0 P_2 \rangle &= \langle I_0 \rangle \langle P_2 \rangle \end{aligned} \tag{7.303}$$

where the shadowing functions are $S(\theta_1)$ for $P_1$ and $S(\theta_1, \theta_2)$ for $P_2$. Equation (7.305) is an approximation when the intensity $I_0$ itself depends on the angles $\alpha$ and $\beta$. However, it seems reasonable physically to assume that the shadowing functions are independent of the scattered field.

The surface shadowing functions are equal to the probability that any arbitrary surface area is illuminated.

The conventional way to start this assessment of probability is to assume that the wave is plane and incident in the $xy$ plane. The probability of any point being shadowed by any other point will be a function of $x$ only. Let $\tau$ be a measure of distance along $x$ pointing to the source. Consider an arbitrary point on the surface $\tau = 0$. If $S(\theta_1, \tau)$ is defined as the probability that $z(0)$ is not shadowed by any $z$ up to $z(\tau)$ then the required shadowing function is

$$S(\theta_1) = \lim_{\tau \to \infty} S(\theta, \tau) \tag{7.304}$$

here $\theta$ is height above the mean line

$$S(\theta_1, \tau + \Delta\tau) = S(\theta_1, \tau)[1 - g(\tau, \theta_1)\Delta\tau] \tag{7.305}$$

where $g(\tau, \theta_1)$ is the probability that $z(0)$ is shadowed by $z$ in the interval $(\tau, \tau + \Delta\tau)$ given that it is not shaded in the interval $0, \tau$ for incident and scattered values of $\theta$.

Performing a Taylor expansion about $\tau$ and retaining terms to $\theta(\tau)$ leads to the differential equation

$$\frac{dS(\theta, \tau)}{d\tau} = -g(\tau, \theta_1)S(\theta, \tau) \tag{7.306}$$

which gives, by means of (7.304),

$$S(\theta_1) = S(0)\exp\left(-\int_0^\infty g(\tau_1, \theta_1)\ d\tau\right). \tag{7.307}$$

Equation (7.306) is the key equation and in itself is complete. However, the real problem arises when evaluating $S(\theta)$ and this requires a knowledge of $g(\tau, \theta)$.

A number of people have considered it. Wagner [236] derives it for a Gaussian height distribution

$$S(\theta_1) = \frac{[1 + \text{erf}(v)][-\exp(-2b)]}{4b}$$

where $\tag{7.308}$

$$b = \frac{\exp(-v^2) - \sqrt{\pi}v\ \text{erfc}(v)}{4\sqrt{\pi}v} \quad \text{and} \quad v = \frac{1}{\sigma_\tau \sqrt{2}\ \tan\theta_1}$$

where $\sigma_\tau$ is the RMS surface gradient in the $|\tau|$ direction. Hence it is asserted that shadowing is a function of all surface geometry, not just height. This can be interpreted as follows.

For a given angle of incidence the probability of one point being shadowed by a nearby point depends primarily on the ratio of their height difference to their separation in the $x$ direction, that is $(z_1 - z_2)/(x_2 - x_1)$. Wagner derives a shadowing function when correlations between shadowing and shadowed surface points are taken into account. It is similar to equation (7.309) only $b$ is a more complicated form of $v$. The parameters $\sigma_\tau$ and $\theta_1$ remain the only parameters of the surface and the wave to affect shadowing. Bistatic shadow functions are also derived which take into account the situation where $\tau = 0$ is shadowed by a second point some distance $\tau_2$ away in the direction of the scattered wave, in the consideration of $g(\tau_1, \theta_1)$ from equation (7.309) This has been approximated by Beckmann and Spizzichino [291] by the probability that $z$ will interrupt the ray connecting shadowing and shadowed points in the interval $\tau, \tau + d\tau$ if $dz/d\tau > \cot\theta$. From figure

7.141 this approximation is reasonable only if one surface point is shadowed by one other surface point and no more. It will therefore break down for some angle of incidence close to grazing when the shadowing function itself becomes close to zero.



**Figure 7.141**

This discussion on shadowing is very similar to the 'motif' method of identifying long- and short-wavelength features on surfaces used in the automotive industry (i.e. the $R$ and $W$ parameters) and also the parameters of the 'rainfall count' used to work out equivalent cumulative stress in fatigue analysis.

Other methods of looking at shadowing have been attempted. One by Hardin [237] allows the source to be a finite height above the surface. This produces the sensible result that the shadowing function increases as the source height increases for fixed surface parameters, or alternatively the shadowing function increases as the correlation length of the surface increases (for a given source height). In the limit when the source is at infinity the results agree with Wagner [304].

All these results demonstrate the importance of slope. Wagner shows that for angles of incidence $\theta_1$ smaller than twice the RMS surface slope Kirchhoff theory is generally not reliable. Geometric shadowing corrections can increase the error in forward scattering as multiple scattering effects are then important. In fact Thorsos [306] indicates that the most important parameter in Kirchhoff theory is the ratio of the correlation length of the surface to the wavelength of light, a comment which is echoed throughout this book! This is followed by the ratio of $\tau_c$ to $\sigma$.

One more way of looking at shadows uses random process analysis and in particular overshoot theory.

Overshoot theory has been used extensively in engineering by investigators such as Crandall [307] in order to find the times of first passage of the displacement or force over a fixed, usually safe, limit.

Consider a random function $f(x)$ and its relationship to another function $\varphi(x)$, as shown in figure 7.142. If, after intersection, $f(x)$ continues upwards, then $\partial f / \partial x > \partial \varphi / \partial x$, the overshoot is positive and vice versa.



**Figure 7.142** Overshoot theory.

Shadowing is an overshoot problem in which $\varphi(x)$ *is* a straight inclined line (of the incident ray). The theory has been used to evaluate bistatic shadowing functions involving two reflections at $\theta_1$ and $\theta_2$.

Bass and Fuchs [292] arrive at a shadowing function

$$S(\theta_1, \theta_2) = \frac{1}{1 + F(s_1) + F(s_2)} \tag{7.309}$$

where $s_1 = 1/\sigma_\alpha \tan\theta_1$, $s_2 = 1/\sigma_\alpha \tan\theta_2$ and $\sigma_\alpha$ is the RMS gradient for surface profiles within the azimuth plane (i.e. along $\tau$). $F$ is a function depending on a single parameter:

$$F\left(\frac{1}{\sigma_\alpha \tan\theta_j}\right) = \tan(\theta_j) \int_{\cot\theta_j}^{\infty} (s - \cot\theta_\alpha) p(s, \sigma_\tau) \, \mathrm{d}s \qquad j = 1, 2 \tag{7.310}$$

where $p(s, \sigma_\tau)$ is the probability density along the direction $\tau$.

The proportion of surface area in semi-shadow is made to be small compared with the full shadowed area or the full illuminated area—this effectively rules out the effect of diffraction in the scattering.

Imposing these conditions leads to the inequality

$$\frac{\lambda}{\tau_c} \ll \frac{1}{\sigma_\alpha} (1 + \sigma_\alpha)^{3/2}. \tag{7.311}$$

This inequality relating the surface spatial parameter (the correlation length) to the wavelength of light again acts as a yardstick for diffraction versus geometric optics, because as $\lambda \to 0$ the geometric optics case results.

In the above, shadowing has been seen in the role of an error or deviation from straightforward scatter. It is, in the sense that it spoils the calculations, but it can be used to estimate the spectrum of the surface. This in effect again uses overshoot theory. The higher peaks throw the shadows as seen in figure 7.143.

The shadows can be used as an approximate way of getting quick estimates of the dominant wavelengths and lay patterns on the surface simply by sweeping the incident beam through an angle of $\pi/2$, preferably in two direction at right angles (figures 7.144 and 7.145).



**Figure 7.143**

**Figure 7.144** Power spectrum estimation by shadowing (normalized).



θ large     θ small

**Figure 7.145** Lay estimation by shadowing.

The average shadow length $\langle l_s \rangle$ is a measure of the presence of this wavelength on the surface. It is usually expressed as a ratio:

$$\langle l_s \rangle_\theta \cdot \sum_{j=1}^{N} l_{sj} \Big/ L \tag{7.312}$$

where $L$ is the length of the surface and $N$ is the number of shadows. Another way of getting this estimate is to cross-correlate the shadow pattern with a projected bar pattern on the surface (figure 7.146).



**Figure 7.146** Cross-correlation by shadowing.

The detector (via suitable optics) integrates the shadow pattern after it has been modulated by the scanning bar pattern, thereby giving an estimate of $P(\omega)$. Moving the pattern across the shadows improves the estimate and allows the variance to be estimated. A big advantage of these methods is that convolution and correlation can be achieved optically.

The shadowing technique when used at different angles can also be used as a quick measure either to determine the optimum sampling length for measurement of the surface or to decide the appropriate motif length for graphic or computer evaluation.

In general the shadowing function $S(\theta)$ can be regarded, for a given $\theta$, as an intrinsic property of the surface because the surface itself produces the shadow by modulating the intensity pattern of the incident light in the same way as vee blocks are used in roundness to provide an intrinsic coordinate datum. Consequently, in no way can shadowing be regarded as a hindrance to understanding; it can often be used as a positive help [308].

### 7.7.11 Scattering from non-Gaussian surfaces

So far no one has completely solved this problem. Efforts have been concentrated on trying to estimate the behaviour of surfaces having multiple scales (e.g. fractal surfaces).

Also, such fields may be generated when there are only a few 'bright spots' contained within the illuminating area. In this case the central limit theorem does not apply because the various scatterers making up the scattering pattern are correlated. This is because the illumination (or resolution of the detector) is only a small fraction of the correlation length of the surface (e.g. [309]).

The non-Gaussian statistics arise when:

1. The dimensions of illumination are comparable with the correlation length.
2. Scattering is dominated by geometric-type reflections (i.e. the specular points).
3. Scattering is in the near field.

Although the complete solution is still unsolved, phase-changing screen models have been used to try to simulate the behaviour. Using such methods Jakeman and colleagues [309, 310] have come to the conclusion that some scattered fields (i.e. from the sea) have an amplitude which is not Gaussian but obeys the so-called semi-empirical $K$ distribution:

$$p(\rho) = \frac{2b}{\Gamma(\alpha)} \left(\frac{b\rho}{2}\right)^\alpha K_{\alpha-1}(b\rho) \tag{7.313}$$

where $\rho$ is amplitude, $\alpha$ and $b$ are constants, $K_{\alpha-1}$ is a modified Bessel function of the second kind and $\Gamma$ is the gamma function. It is surmised that this sort of probability density arises as a result of fractal-type surfaces. However, this is not yet proven.

It is possible that these so-called $K$ surfaces are of the fractal or quasifractal type. In particular, when the surface is of the form designated 'subfractal' by Jakeman, equation (7.313) can be simplified somewhat to give a formula (equation (7.314)) which is the probability density of intensity $p(I)$, where $p(I)$ is

$$p(I) = 2K_0(2\sqrt{I}) \tag{7.314}$$

and $K_0$ is a Bessel function of zero order.

The reason for this type of distribution is not known but the form can be shown to result when the number of steps fluctuates according to a negative binomial distribution. It seems that this type of distribution would also follow for a Poissonian distribution, in which case it would be readily explained in terms of the 'unit event' of surface generation and as such would relate to the 'unit event of scatter—the specular point'.

### 7.7.11.1 Fresnel scattering

Whereas Fraunhofer scattering from a surface corresponds to the scatter behaviour of a plane incident wave when the scattered light is picked up a long way away, Fresnel scattering relates to a spherical wavefront hitting the surface and the scattered spherical wavefronts being picked up close to the surface. Fraunhofer

scattering corresponds to the linear form of scatter (and diffraction) whereas Fresnel scattering corresponds to the quadratic form. This makes the calculations in the Fresnel case more difficult than in the Fraunhofer.

This difference reveals itself in the formulae in one dimension and for uniform illumination:

$$u = C \int_{\text{aperture}} \exp[-j(2\pi/\lambda(Y/x))] \; dy \qquad (7.315)$$

from Fraunhofer, and

$$u = C \int_{r_1}^{r_2} \exp(-j\pi r^2/2) \; dr \qquad (7.316)$$

for Fresnel, where

$$r = (Y - y)\sqrt{\frac{2}{\lambda X}} \qquad r_1 = Y + \frac{d}{2}\sqrt{\frac{2}{\lambda X}} \qquad r_2 = Y - \frac{d}{2}\sqrt{\frac{2}{\lambda X}} \qquad (7.317)$$

and $d$ is the slit width. In general the Fourier case is only an approximation of the Fresnel case.

It follows that in Fraunhofer conditions the scatter is dependent on the illuminated area and the angle of scatter $\theta$, whereas in Fresnel conditions the scatter is dependent more on the distance of the detector.

This spherical scatter shows itself when the observer is close to an illuminated surface.

The actual light as seen by a detector placed near to a random phase screen suffers speckle if the detector is close enough. This is due to the presence of local caustic or spherical waveforms generated by the surface roughness. As the detector is moved away from the screen the contrast rises until it reaches a maximum, after which the contrast assumes Gaussian statistics as shown, especially if the source is poly-chromatic (figure 7.147).



**Figure 7.147** Speckle contrast.



**Figure 7.148** A is focus of marginal rays, B is focus of paraxial rays.

### 7.7.11.2 Caustics

The caustic effect is shown above in figure 7.148 at large scale and assumes geometric optics within the simple facet although not outside.

Light reflected or transmitted through an optical element, in this case a curved facet of many sides making up the surface, does not focus at a single point if the curve is not a true parabola. The cone of rays resulting is called the caustic surface. It is a curved surface which is a function of the curved facet. It is usually centrosymmetric.

Because the effect considered is close to the random phase screen in figure 7.148, at the surface it is governed by the Huygens–Fresnel diffraction integral whose simple form is given in equation (7.313).

### 7.7.11.3 Fractal surfaces

It has been observed that many surfaces having a highly polished appearance have a power spectrum which is characterized by its very simple form in terms of an inverse power law:

$$P(\omega) = k\omega^{-\nu} \tag{7.318}$$

where $k$ and $n$ are constants for the surface. This form has been slightly refined by Mulvaney *et al* [269] into the form

$$P(\omega) = \frac{k}{1 + (\omega/\omega_c)^2} \tag{7.319}$$

or similar. To start with, it is obvious that the intrinsic parameters associated with such surfaces are $k$ and $\nu$ taken from the power spectrum. This is different from the height and sparing criteria used for ordinary surfaces, that is $\sigma$ and $\tau_c$, where $\sigma$ is the RMS (or $R_a$ or $R_t$) acting as the height parameter and $\tau_c$ the correlation length (or average wavelength or $S_m$) is acting as the spacing ordinate.

The main difference between the representation of equations (7.318) and (7.319) is that the value of '$\nu$' is not necessarily integral. The other difference is that when plotted on a log scale the equation (7.318) is a straight line of slope determined by $\nu$. On the other hand the form of equation (7.319) has a corner when $\left(\dfrac{\omega}{\omega_0}\right) \sim 1$.

The critical point here is that a fractal surface is naturally measured using $k$ and $\nu$. Conventional surfaces have parameters associated with them whose reliability is directly proportional to the number of degrees of freedom present in the data. This is taken to be $L/\tau_c$. According to the standard this would normally be in the region of 30 or more in order to get about 5% reliability in the value. Unfortunately fractal surfaces tend to have values of $\sigma$ which are divergent and not convergent either in value or in accuracy. In other words, the value of $\sigma$ will depend on the length of record as well as its variability. This is not the case for conventional surfaces, so values of $\sigma$ on fractal surfaces would necessarily exhibit severe statistical fluctuations. These would tend to obscure any dependence of the measured roughness parameters on the measurement process.

Before considering the scatter it is interesting to note that whereas equation (7.319) represents a fractal-type surface, it only does this in one dimension. In two dimensions and for an isotropic surface the form is quite unusual:

$$P(\omega) = \frac{\Gamma((n+1)/2)}{2\Gamma(\frac{1}{2})\Gamma(n/2)} \frac{K_n}{\omega^{n+1}} \tag{7.320}$$

where $\Gamma$ is the gamma function. It is this quantity which, strictly speaking, should appear in the analysis of the surface.

Another point is that the mathematical analysis of fractals makes use of two different parameters: the

Hansdorff–Besicovitch dimension $D$

$$D = \tfrac{1}{2}(5 - v) \tag{7.321}$$

and the so-called length topothesy $L$ where

$$L^{v-n} = -\frac{1}{2}\,\frac{(2\pi)^n}{\Gamma(n)\,\cos(n\pi/2)}\,K_n. \tag{7.322}$$

Physically this is the length over which the chord connecting two points has an RMS slope of unity:

$$
\begin{aligned}
&\text{for } n = 1, \; D = 2 \quad \text{the fractal is called `extreme'} \\
&\text{for } n = 2, \; D = 1.5 \quad \text{the fractal is `Brownian'} \\
&\text{for } n = 3, \; D = 1 \quad \text{it is a `marginal' fractal}.
\end{aligned} \tag{7.323}
$$

In the case of $n = 2$ the isotropic Brownian $P(\omega) = K_2/4f^3\,\mu\text{m}^4$ and $L = 2\pi 2K_2\mu\text{m}$.

Fractal profiles have the interesting feature that when they are stretched along the $x$ axis by a factor $M$ (see [270] and along the $z$ axis by $M^{2-D} = M^{(n-1)/2}$, the result has the same statistical properties as the original profile. It is self-similar if both magnifications are equal, but self-affine if not.

This apparent absence of internal length scales seems plausible for the roughness of highly polished surfaces since the presence of such scales would be an unwanted signature of the polishing process. At the same time it makes the characterizing of such surfaces difficult in terms of intrinsic height and length parameters.

The idea of 'long-range' order and infinite correlation lengths mean that definitions of peaks and even RMS in the height as well as all length parameters are relatively meaningless. Thus, not only are there considerable problems in digital sampling techniques, as shown earlier in chapter 3, but there are considerable conceptual problems because of the nature of the surfaces themselves.

On the one hand, there is the difficulty of characterization of such fractal surfaces using traditional means yet, on the other hand there is one small benefit of fractal surfaces which is that when magnified equally in both directions they have the same properties for all scales of size. This means that all scales of size can be dealt with using the same characterization procedure, unlike the waviness, roughness, error of form split which has tended to develop over the years because investigators have mixed up the reasons why they were separated in the first case.

Luckily no surfaces have been found that are fractal over more than three orders of magnitude. This all relates to methods of generation or manufacture whether erosive, abrasive or by growth. One reason is instrumental: the measuring instrument will always provide the short-wavelength cut-off because of resolution limits and a long cut-off because of range limitation, so it will always only ever be possible to characterize what can be measured.

Even then fractal characteristics do not follow across many scales of size. Franklin and Schneider [312] have found that, in fracture surfaces and stress corrosion cracks, fractal behaviour exists over only a limited range and that the deviations outside this range are great. Similar results have been found by Underwood and Banerji [313] and others [314, 315]. Modified or subfractal surfaces have been advocated as being more realistic. These have the fractal behaviour but over a severely limited range. It seems that the whole concept of fractal characterization becomes somewhat meaningless when this subdivision has to be used. It constitutes ways of forcing theories onto real surfaces almost for the sake of it! Even when a surface has been designed as fractal, how to measure its dimension $D$ has been found to be difficult [316]. Different results are being found depending on whether or not box-counting techniques or Fourier analysis techniques are being used. It seems that, on balance, the spectral method is more reliable. The reason for the enthusiasm about fractal surfaces is based on their 'self-similar' behaviour through different scales of size.

But there is a conceptual problem here which is concerned with 'scale invariance'. Fractals are scale invariant but is this a useful property? In one respect it is decidedly not. This is because many functional

properties, and details of manufacture are highly dependent on the scale of size. Such properties would therefore look invisible to fractal analysis. So to attempt to use fractal analysis to characterize any dynamic use of the surface would be a waste of time.

A very good example is the behaviour of surfaces when irradiated with light. The behaviour is very scale dependent: it depends on the wavelength of light. It has just been shown how the type of light scatter depends on the ratio of the roughness $R_q$ to the wavelength $\lambda$ of the light. Scale invariance can only happen above $\lambda$ or below $\lambda$ but not through it!

Dynamic behaviour of surfaces also has pronounced scale dependence. Consider Figure 7.149.



**Figure 7.149** Force/energy balance.

This shows the well-known change in effectiveness of the different components of force with scale of size. Mass effects reduce in value by $l^{-3}$ damping by $l^{-2}$ and elastic by $l^{-1}$ as the scale of size $l$ is reduced. Rotating forces concerned with inertia fall off even more rapidly as $l^{-5}$. As a consequence of this the dynamic balance between forces changes dramatically with size; fractal parameters have no use in such situations. The same effect occurs in laminar fluid flow but not necessarily in turbulent flow.

If interest is in atomic surfaces and nanosurfaces and their physical properties (for example properties of adsorption or fracture), then if the surfaces were fractal in nature it should be possible to take measurements in a more convenient, measureable, range (i.e. micrometres) for which there are many available instruments and then to scale the findings back down to the atomic level; the fractal behaviour enabling a way to 'abseil' down the scales of size! This property would be useful in the same way that the Wigner function has been suggested here as a mathematical tool to allow a progression up the scale of size from the atomic scale to the basic engineering scale. In both cases the basic philosophy is that of linking together the different scales of size.

Even if the correlation function is infinite in a true fractal surface, it is possible to have some sort of measure which is spectrally related and finite, such as the structure function first used in surface analysis by Whitehouse in 1970 [297] and subsequently in detail by Sayles [317]. This has been defined elsewhere in this book together with its properties.

However, as

$$S(\tau) = \langle z(x) - z(x + \tau) \rangle^2 \tag{7.324}$$

this exists for a fractal surface and takes the form of

$$L^{3-\nu} \, | \, \tau \, |^{\nu-1} \tag{7.325}$$

where $L$ is the topothesy for scattered waves off a fractal surface. Berry [318] gives this a name—diffractals—which expresses the behaviour rather well.

There are other ways of describing fractals based on a variant of Mandelbrot's work [319]. In fact the Mandelbrot ideal fractal $f(x)$ can be incorporated into the Weierstrass function

$$f(x) = \sum_{n=1}^{\infty} \frac{\cos 2\pi\gamma^n x}{\gamma^{(2-D)n}} \qquad (7.326)$$

in which the wave numbers are 1, $\gamma$, $\gamma^2$. This can be scaled to give the Weierstrass-Mandelbrot (W-M) function.

$$f(x) = \sum_{n=-\infty}^{\infty} \frac{(1-\cos n\pi\gamma^n x)}{\gamma^{(2-D)n}}. \qquad (7.327)$$

Berry and Lewis [254] have plotted W–M functions calculated using $\gamma = 1.5$ and different values of $D$. They also show a very interesting result, which is that the discrete frequencies of the W–M function approximate to a continuous spectral density in frequency $\omega$:

$$G(\omega) = \frac{1}{2\ln\gamma}\ \frac{1}{\omega^{(5-2D)}} \qquad (7.328)$$

which suggests a link between the spectral and fractal approach. It also suggests that similar curves might be obtained by taking

$$f(x) = \sum_{n=-1}^{\infty} \frac{1}{n^{(2.5-D)}} \cos(2\pi nx + \varepsilon_n) \qquad (7.329)$$

provided that $\varepsilon_n$ is random. This suggestion comes from Greenwood, who argues that this brings us back to the Sayles and Thomas $G(\omega) \simeq 1/\omega^2$ again, confirming that the link between spectral and fractal methods is in fact two views of the same thing and that the two approaches are anything but independent. The major difference is that the index for fractals need not be integer.

Going back to structure functions and light scatter, it is sometimes simpler to define the structure function as

$$S(\tau) = |\tau|^{\nu-2}/L^{\nu-2} \qquad \text{with } 2 < \nu < 4 \qquad (7.330)$$

in which the surface is differentiable and has a continuous slope. This is important because it is always assumed for surfaces that investigation of the scatter in the far field is some measure of the proportion of slopes facing the appropriate direction. To use this concept is difficult when the applicability of rays and geometric optics is being questioned. The intensity pattern $\langle I(\theta)\rangle$ is in fact a reasonably stable form when

$$\langle I(\theta)\rangle = P_0([kL]^{1-2/\nu}\sin\theta) \qquad (7.331)$$

where

$$\int P_0\exp(-j\lambda x) = \exp(-A\lambda^\nu). \qquad (7.332)$$

For a smoothly varying surface

$$\langle I(\theta)\rangle \simeq p(m)(\sin\theta) \qquad (7.333)$$

where $p(m)$ is the slope distribution [310].

The scattered distribution has a tail which falls off as $\sin^{\nu+1}\theta$. According to Jakeman $\nu$ and $L$ can be worked out from this distribution. According to Whitehouse [368] there are much more elegant ways.

In equation 7.330 the fractal parameters are given in terms of the structure function. A logarithmic plot of structure function against $\tau$ has been used to determine the fractal dimension. The plot is shown in figure 7.150.

**Figure 7.150**

In this the dimension $D$ is found from the slope and then the topothesy $L$ from the intercept. The problem is that the plot involving a number of values is required.

There are other ways which are much faster. One of these is given below [368].

From the approximation that $S(\tau) \sim \tau^2 A''(0)$ the further derivations of the autocorrelation $A(\tau)$ can be found e.g. $A^{iv}(0)$. Then using the well-known expression that the peak count density on the surface profile '$m$'

is given by $m = \dfrac{1}{2\pi} \sqrt{\dfrac{A^{iv}(0)}{-A''(0)}}$ peak density can be related directly to '$m$' giving

$$m = \frac{1}{\pi\tau} \sqrt{(D-1)\left(D - \frac{1}{2}\right)} \tag{7.334}$$

the topothesy can be found either given a value of the structure function or a value of the high spot density $n$

$$L = \left(n\pi R_q \tau^{D-1}\right)^{\frac{2}{2-1}} \tag{7.335}$$

So

Where $R_q$ is the root RMS surface roughness. However care should be taken if the surface is generated by grinding when $D$ is close to 2 when the $D$ value would be overestimated.

The strategy outlined above links the fractal parameters to a measurable unit—in this case peak count via the structure function making the approximation that $\tau''(0 \sim S(\tau))$.

An alternative method which bypasses this approximation uses the 3 point discrete value of peak density given in chapter 2.

$$\text{Thus} \quad m = \frac{1}{\pi\tau} \tan^{-1}\left(\frac{3 - 4p_1 + 2p_2}{1 - p_2}\right)^{\frac{1}{2}} \tag{7.336}$$

Where $p_1 = A(0)A(\tau)$ and $p_2 - A(0)A(2\tau)$

(7.338) can be written with the structure function instead of the autocorrelation function.

$$\text{Thus} \quad m = \frac{1}{\pi\tau} \tan^{-1}\left(\frac{wS(\tau) - S(2\tau)}{s(2\tau)}\right)^{\frac{1}{2}} \tag{7.337}$$

$$\text{giving} \quad m = \frac{1}{\pi\tau} \tan^{-1}\left(1 - 2^{2(D-1)}\right) \tag{7.338}$$

### 7.7.11.4 Fractal slopes: Subfractal model

Unlike the fractal height, if the class of surface called subfractal here is used, some use can be made of geometric optics. As previously said, this type of surface has a structure function $S(\tau) = |\tau|^{\nu-2}/L^{\nu-2}$. As in the case of smooth, continuous surfaces the angular distribution of intensity is determined by the surface slope distribution, unlike the fractal height surface which is not. In the case of the Fraunhofer and Fresnel regions the behaviour of light scatter of the second moment $\langle I^2 \rangle / \langle I \rangle^2$ points to the absence of any geometric optics effects. In neither of these regions is there enhancement of the Gaussian intensity distribution

$$\langle I^2 \rangle / \langle I \rangle^2 = n! = 2. \tag{7.339}$$

Only single functions dependent on products of function $kL$ and $\omega$ or $z$ are obtained. Normal continuous smooth surfaces have families of curves which are obtained as a function of surface height. For example, the caustics described earlier should always be observed on conventional surfaces (these are in effect geometric singularities in the intensity pattern) but not in the case of fractal surfaces where everything in the intensity pattern is diffuse.

Usually the higher-order moments of the intensity distribution are of interest in terms of illustrating the multiscale nature of surfaces. For example, in the far field a geometrical optical regime should exist which is dominated by the slope distribution (when the diffraction regime is not operable, as when $k^2\sigma^2 \ll 1$). Here an incident beam is reflected through an angle $\Omega_2$ whose RMS slope is determined by the slope structure function at the outer scale size, $\sqrt{S(\tau)}$, and it spreads out as a function through angle $\Omega_2$ which is determined by the slope structure function corresponding to the aperture $W$, that is $\Omega_2 = \sqrt{S(W)}$.

If a rectangular light distribution incident on the surface is assumed, then

$$\frac{\langle I^2 \rangle}{\langle I \rangle^2} \simeq \left( \frac{\zeta^2}{W^2} \right)^{(\nu-2)/2} \tag{7.340}$$

where $\zeta$ is the largest scale of size of the phase fluctuations.

The power law dependence on $W$ is characteristic of the model. It differs from the behaviour predicted for smooth surfaces, which gives a deviation from Gaussian statistics that is inversely proportional to the illuminated area (i.e. proportional to $W^{-2}$).

If no outer scale value is included then the amplitude fluctuations saturate at large distances from the surface at a value greater than for Gaussian speckle. This is because ray density fluctuations do not average



**Figure 7.151** Speckle contrast for surface without large-scale detail: (a) one-scale surface; (b) multiscaled surface.

out; the surface slope is correlated over an infinite range (figure 7.151). If an outer range is introduced then the intensity fluctuations tend to the values expected for Gaussian speckle as uncorrelated regions of the surface begin to contribute to the scattered field, thereby smoothing it out somewhat The intensity in this model is substantially similar to the $K$ distribution mentioned earlier.

The most interesting feature of the slope fractal model is that although surface slope is well defined, its curvature is not, so caustics or focusing cannot occur.

Jakeman [242a] notes that although most surfaces seem to have some fractal behaviour somewhere within their geometric range, fluids do not. They are intrinsically smoother and scatter in much the same way as continuous media. However, the possibility that some 'subfractal behaviour' might be attributed to turbulent systems has not been ruled out! But, in engineering, Markov properties dominate.

### 7.7.12 Aspherics

There is one area in optics where the geometry of the surface is of paramount importance. This is in aspheric optics.

The utilization of short wavelength sources for ultra precise microlithography capable of producing a submicron line has brought some unwelcome problems. As the wavelength of the source gets shorter and the power of the light sources increases problems such as loss of light surface damage and poor signal to noise ratios due to scatter have stretched optical design techniques. In order to improve the $S/N$ ratio a surface roughness of less than 0.2 nm $R_q$ is required. Since low absorption materials in the short wavelengths are not available, the number of lenses has to be reduced. Also the demand for lighter and smaller goods leads to fewer optical elements. Aberrations therefore such as spherical aberration have to be corrected with fewer lenses. All this points to using more aspheric lenses.—with very low surface finish [320].

Correction of Seidal aberrations is one of the few examples where the geometrical form of the surface represents the function and the surface roughness is an irritant.

The need to correct for spherical aberration stems from the central limit theorem in statistics. Generation of a lens or mirror surface automatically results in a spherical surface. Maudesley understood this and counteracted this tendency by using three glass or metal blocks, rubbing sequentially against each other, to make engineers' flats in the 1700s in Manchester.

If the manufacturing process is sufficiently random, which is most often the case in lapping, the statistics of generating is Gaussian i.e. in two dimensions $x$ and $y$; the generator is the circle $R^2 = x^2 + y^2$ from the statistics $\exp\left(-\left(\dfrac{x^2 + y^2}{2}\right)\right)$ imposed by the central limit theorem and hence, by extension, to the sphere. As a consequence lenses were spherical. Unfortunately the requirement for a point focus from parallel light is a parabolic wavefront. The spherical surface as therefore to be modified optically, (e.g. using a Schmidt plate), or by ensuring that the correction is made mechanically.

The general equation of an aspheric is given below. It comprises a basic conic section onto which an axial symmetric deviation is superimposed.

$$z = \frac{\mathrm{sgn}.x^2}{R = \sqrt{r^2 - (1 + K)}} + a_1\left|X\right| + a_2 X^2 \, \ldots \ldots \, a_{12} X^{12} \tag{7.341}$$

where $x$ is the horizontal distance from the aspheric axis and $z$ is the corresponding vertical distance.

Sometimes the power series extends to $A_{24} x^{24}$. These extra terms are rarely needed and usually result because of the occasional need to circumvent proprietary shapes.

In equation (7.341) sgn is $+1$ for concave profiles and $-1$ for convex. $R$ is the base radius of curvature from which the deviations are to be made.

K is the conic constant  
K<−1     hyperboloid  
K=−l     paraboloid  
1<K<0     ellipsoid  
K = 0     sphere  
K>0     oblate ellipsoid  

See Figure 7.153

**Figure 7.152** Conic sections.

**Figure 7.153** Conic shapes and deviations.

The general conic set is shown in Figure 7.152 as generated from the standard cone shape by intersection of planes.

## 7.8 Scattering by different sorts of waves

### 7.8.1 General

Consider first the Kirchhoff type of scattering. This has been examined in some detail in chapter 4. The differences which arise between the different wave types have not, however, been specifically pointed out. In what follows an explanation of the different types of scatter will be given. The basic incident waves involved are electromagnetic, acoustic or elastic. Each of these different types of scatter have special cases. For instance, within EM scattering could be included x-ray scattering, which has problems of its own owing to the high absorption. Also Raman scattering should be considered. Figure. 7.154 shows a breakdown of the different scattering modes in terms of elastic and non-elastic properties.



**Figure 7.154** Radiation regimes scattered from surface.

### 7.8.1.1 Electromagnetic waves

In the general case of vector theory the surface modulates the elastic field vector, which is at right angles to the direction of propagation. This field vector for any specific value in the propagation path can be an ellipse, a circle or a line, the general case being an ellipse.

Because any wave can be considered to be the superposition of two orthogonal linearly polarized waves, it is often most convenient when considering the depolarizing effects of surfaces to confine the incident ray to being simply linearly polarized. The polarization can be made either vertical or horizontal, but it is usual to consider those directions which are going to 'couple" with the direction of the lay of the surface roughness. In figure 7.155, for example, the horizontally polarized rays will be little affected by the surface, whereas the vertical ones will.



**Figure 7.155** Polarized light incident on corrugated surface.

### 7.8.1.2  Elastic, ultrasonic and acoustic scattering

Elastic and acoustic waves are in principle no different, and elastic waves can in many instances be regarded as similar to EM waves. Elastic waves have a compression component in the direction of the incident ray, like acoustic waves, but also two transverse shear waves orthogonal to each other and perpendicular to the direction of incidence, making three nodes in all. They are termed elastic because the radiation wavelength is not changed by the reflection. Elastic scattering is a vector field phenomenon and can be complicated because there can be an interaction between the surface roughness and all three wave modes, as opposed to two for EM waves and one for acoustic waves. This three-mode interaction makes the calculation of the reflection coefficients complicated. However, it has been accomplished for surfaces which are corrugated in one dimension. This corrugation is usually chosen to be perpendicular to the plane containing the incident wave; hence horizontal transverse coupling is zero [321].

Without going into the theory, which is based on Kirchhoff's equations as before, in electromagnetic radiation some useful facts emerge. These are that the *surface heights* influence the *phase* of the scattered wave and the *surface gradient* influences the *amplitude.*

### 7.8.2  Scattering from particles and the influence of roughness

Scattering of waves from particles is a large subject and is very important because it basically determines what happens not only in gases but also in solids. Essentially the material properties can no longer be ignored.

### 7.8.2.1  Rayleigh scattering

In Rayleigh scattering the wave is incident on particles or a collection of particles. The mode of scattering depends on whether the wavelength is large or small compared with the diameter of the particle. This is definitely a case where dimensional metrology is important. Also important is the shape of the particle, that is are the particles angular or truly spherical? This is more an application of surface metrology because it is shape dependent as in roundness. It also depends on the distribution of particles, whether random or periodic, and their average separation. Rayleigh scattering is the case when the wavelength of the light is usually much greater than the particle size. Hence only a collection of particles is considered.

If the wavelength is small and the particles regular, this can be considered to be a case of Bragg scattering.

### 7.8.3  Bragg scattering

Bragg scattering is usually referred to in the scattering of EM radiation, usually x-rays in a solid crystal, resulting in x-ray diffraction. Here the wavelength may be of the order of 0.1–3 Å and comparable with the crystal lattice spacing. This reduction of wavelength over the visible spectrum by $10^3$ or more provides a metrology unit which is comparable with the atoms themselves and, as such, x-rays are beginning to provide just the metrology tool needed to enable precision engineering to enter into the nanometric and atomic domains.

Interference maxima are obtained using the well-known Bragg condition

$$2a\sin\theta = m\lambda \tag{7.342}$$

where $m$ is an integer, $\theta$ the scattered angle and $a$ the lattice spacing. The calculation assumes that a secondary wave is created at each atom, which is slightly different from conventional diffraction because there are no Huygens waves to consider. The secondary wave created at each atom is the induced dipole emission of the incident wave and not a fictitious wave required by Huygens' principle. Nevertheless, it is always referred

to as x-ray diffraction. Because of the absorption problems x-rays have to be incident at a very small glancing angle, say 1° or less, if a reflection is to be achieved at the surface (figure 7.156).



**Figure 7.156** X-ray scattering off rough and smooth surfaces.

The surface of the crystal face is usually smooth, but there is now evidence that the roughness of it is important and can impair the performance of x-ray microscopes and telescopes which have to work with components having very acute glancing angles.

### 7.8.3.1 Non-elastic scattering

The type of scattering looked at so far could be considered to be elastic in the sense that the wavelength of the incident wave is not changed by the interaction, only its direction. An elastically scattered photon is one with the same energy $hv$ and a change in direction, whereas an inelastic photon changes both direction and frequency upon interaction with the solid.

An extreme case of the inelastic scattering of light is Compton scattering by electrons and ions.

There are two ways in which linear chains of ions can vibrate: in the optical mode and in the acoustic mode. These are shown very simplistically in figure 7.157 for the optical mode of vibration by the ion chain and also the acoustic mode. In the latter, neighbouring ions move in the same direction.



**Figure 7.157** Non-elastic scattering—simplified model for ion chains.

Light incident on the solid which interacts by means of the acoustic mode of vibration of the chain of ions is called Brillouin scattering. This produces a weak scatter pattern in direction and also in frequency. The frequency of the oscillation, $\Omega$, is given approximately by

$$\Omega = 2V_s 2\pi / \lambda \tag{7.343}$$

where $V_s$ is the velocity of the phonon in the crystal.

At very small angles of incidence x-rays are totally reflected from solid surfaces. Above some critical angle the x-ray reflectivity falls very rapidly. In a perfect scheme, with non-absorbing material and a perfectly sharp and flat, smooth interface, the fall in transmission begins abruptly at $\theta_c$, the critical angle, and drops at $(2\theta)^{-4}$. In practice there is some absorption; $\theta_c$ is a fraction of a degree.

Roughness is easily seen as an effect because, even for a perfect silicon substrate, surface roughness causes the drop-off to occur more rapidly, as in figure 7.158(b). The difference is clear and assignable to the roughness. It can therefore be used to measure the roughness. This consequently becomes the most sensitive method yet of measuring roughness. It is potentially better even than the AFM and STM because it has no natural lower-wavelength cut-off unlike the scanning microscopes in which an integrating stylus is present. This technique applies to specular reflection; diffuse scatter will be considered later. For the acoustic mode the oscillation frequency becomes smaller and smaller with longer and longer wavelengths. The optical



**Figure 7.158** X-ray spectrum for smooth and rough crystal: (*a*) without and (*b*) with roughness.

mode scattering, called Raman scattering, which is like the Brillouin scattering, is inelastic. The change in frequency is given

$$\Omega = \left[ 2f \left( \frac{1}{m_+} + \frac{1}{m_-} \right) \right]^{1/2}$$

(7.344)

where $m$ is the ion mass and $f$ is the force between them. The effect is to produce a spectral line above and below that of the incident wavelength which is much smaller but measurable.

For molecules the frequency spectra caused are very complicated but understandable.

Just how the roughness influences the scatter in inelastic conditions is a difficult subject. It obviously inhibits the ion vibration at or near the surface interface on the atomic scale. This has not previously been a factor, but now that x-ray diffraction and lithography are becoming more of an engineering tool in the semiconductor industry, it is. Also, it could be that because the friction and general dynamic performance of miniature electromagnetic systems (MEMs) are governed by the surface texture, the Raman scattering could even be used to measure it.

### 7.8.3.2 Influence of roughness—thin-film measurement

The interface between films and surfaces is in general rough, not smooth. Geometric roughness and variations with the material near to the surface—the subsurface—are also important. For example, changes in electron density can change the angular scatter from x-rays [322]. This is especially true in the measurement of thin films. x-rays have an advantage over optics because the influence of many chemicals on the surface is simply ignored and only the material itself is considered to be significant. The specular scattering in the region just above the critical angle contains information not only of the electron density and the thickness of the film but also of the surface and interface roughness on the atomic scale.

The roughness effect has to be calculated for the thin-film layers, for example, with some difficulty and involves the work of Parratt [248]. This takes the following abbreviated form.

The reflectivity amplitude $R_{j-1,j}$ at the ideal interface between layers $j-1$ and $j$ is related to $R_{j,j+1}$ by the recursive relationship

$$R_{j-1,j} = a_{j-1}^4 \left( \frac{R_{j,j-1} + r_{j-1,j}}{R_{j,j+1}\, r_{j-1,j} + 1} \right) \tag{7.345}$$

where

$$R_{j,j+1} = a_j^2 \left( \frac{E^r_{j,j+1} + 1}{E_{j,j+1}} \right) \tag{7.346}$$

where $E^r_{j,j+1}$ and $E_{j,j+1}$ are the reflected and incident amplitudes on the interface between layers $j$ and $j+1$ or an abrupt ideal interface. $r_{j-1,j}$ is given by

$$r_{j-1,j} = \left( \frac{f_{j-1} - f_j}{f_{j-1} + f_j} \right) \tag{7.347}$$

where $f_j = (\theta^2 - 2\delta_j - 2\sqrt{-1}\beta_j)^{1/2}$ and the refractive index $nj$ for layer $n = n_j = 1 - \delta_j - \sqrt{-1}\beta_j$. $a_j$ is the phase factor corresponding to half the layer thickness $t_j$:

$$a_j = \exp\left( -\sqrt{-1}\frac{\pi}{\lambda} f_j t_j \right). \tag{7.348}$$

Equation (7.345) is effectively a Fresnel equation and is successfully applied to each interface in turn working up from the substrate to the film surface. The intensity reflectivity is simply the square of the amplitude reflectance at the surface. The calculation is carried out for each angular setting of the incident angle.

If the reflectivity is small the Born approximation is made. Thus

$$r = r_F \int_{-\infty}^{\infty} \left( \frac{-1\mathrm{d}\rho(z)}{\rho_{-\infty}\mathrm{d}z} \right) \exp(-jQz)\ \mathrm{d}z \tag{7.349}$$

where $\rho_{-\infty}$ is the electron density in the substrate, $r_F$ is the Fresnel reflection amplitude and $\rho(z)$ is the longitudinal density distribution; $Q = k_s - k_i$ (scattered and incident waveform).

For a rough interface whose interfacial position is described by a Gaussian distribution the Fourier integral (7.351) gives the familiar form of the reflectivity amplitude for a rough interface with root mean square width $\sigma_{j-1,j}$.

$$r = r_F \exp[-\tfrac{1}{2}(\sigma_{j-1,j}Q_j)^2]. \tag{7.350}$$

For reflectances near unity a distorted-wave Born approximation (DWBA) is used. Equation (7.347) modified to give

$$r_{j-1,j} = \left( \frac{f_{j-1} - f_j}{f_{j-1} + f_j} \right) \exp[-\tfrac{1}{2}(\sigma^2_{j-1,j}QQ_{j-1})^{-1}]. \tag{7.351}$$

This equation gives the reflection amplitude from the interface between the $j$th and $(j-1)$th layers. It is this modification to the theory (i.e. from (7.345) to (7.351)) which is used to incorporate the effects of surface and interfacial roughness. The advantage of incorporating roughness in this way is that it is taken care of in the theoretical analysis and so does not slow down computation by adding it into the simulation of scatter.

## 7.9 System function

### 7.9.1 Surface geometry and tolerances and fits

For a long time the relationship between tolerances of size and surface roughness was not realized. Upper and lower limits of size had to be specified and met. Also upper and lower limits on surface roughness were specified. In fact the sets of limits are not independent and to be able to fit components together satisfactorily requires that the two are linked together. A typical situation is shown in figure 7.159.



**Figure 7.159** Surface geometry and tolerance.

### 7.9.2 Tolerances

Tolerance can be specified in two ways, either as a deviation in one direction only from the precise or basic size, or as an allowable deviation in either direction from the nominal. Generally a dimension to a surface which is not contacting anything is given by tolerance of size on either side of the nominal. This is known usually as a bilateral tolerance. Dimensions to surfaces which have a fit fundamental to the function of the part (e.g. shafts in holes as in figure 7.160) are most likely to have a tolerance in one direction only. This is referred to as a unilateral tolerance. Unilateral tolerances are usually arranged so that nominal size dimensions give the tightest fit. This automatically gives a situation in which the nominal size is the 'maximum material' condition; that is, further material could be removed and still allow the dimension to be within the required tolerance. This is useful because it encourages manufacture to be undertaken which is least likely to result in scrap. The whole philosophy behind this is that it is much easier to remove material than it is to add it!



**Figure 7.160** Limits on surface texture using material ratio values.

There was an attempt [324] to link roughness and tolerance early in the 1960s using the material (then the Abbott-Firestone) ratio curve and the mean lines $L_{m1}$ and $L_{m2}$ of the two surfaces (figure 7.160) The assumption for the surface roughness is that the probabilities $P_1$ and $P_2$ are both normal (Gaussian). Taking the case where $P_1$ and $P_2$ are independent

$$\langle l \rangle = \langle P_1 \rangle + \langle a \rangle + \langle P_2 \rangle$$
$$\sigma_e^2 = \sigma_{P_1}^2 + \sigma_{P_2}^2 + \sigma_a^2$$

(7.352)

when both sides are equal as in the case of a shaft shown ($P_1 = P_2$):

$$\mu_l = \mu_a + 2\mu_P$$
$$\sigma_l^2 = \sigma_a^2 + 2\sigma_P^2 \tag{7.353}$$

where $\mu_l$ is the mean of $l$ and so on. If $l_{max}$ is the maximum allowable dimension and $l_{min}$ the minimum allowable dimension then, when the normal graph of errors $N(\mu_l, \sigma_l^2)$ is plotted as a function of $l$ as in figure 7.161, sufficient area should be enclosed between $N(\mu_l, \sigma_l^2)$ and $l_{max} - l_{min}$ to give an acceptable probability of occurrence (i.e. a large enough area). The surface roughness (as determined by $\sigma_l$ and $\mu_l$) is usually adequate for the tolerance $l_{max} - l_{min}$.



**Figure 7.161** Tolerance value with to tolerance distribution.

It is interesting to see the way in which size tolerances change as a function of size (figure 7.162).



**Figure 7.162** Tolerance as a function of component size.

Section 1.6 of BS 4500 suggests that geometry, form and surface texture should be included in the specification of tolerances. It could be argued that the situation is precisely the other way round. The principal reason why tolerances in dimension are needed is because of the presence of surface texture, form, etc. The reason why the practical tolerance curve of figure 7.137 rises with dimension is simply because, as the size increases, the potential number of surface geometric errors increases [325]. The basic question therefore is whether a size tolerance is no more than a measure of the added texture and form errors—or is it more? The answer is that although the texture and form errors undoubtedly are large components of the allowed size tolerance, they can only be regarded as a part. Other factors include positioning errors in the machine tool,

thermal distortion produced by a change in the environment, and so on. In fact in so far that the surface texture is process controlled it is true to say that size tolerance and surface texture are related (but not uniquely) (see [326] on this). Vectorial ways of working out composite tolerances in complex systems will reduce confusion.

If the nominal size of a part increases from 1 to 500 mm the allowed value for the total roughness height increases by a factor of 6.3 for the same tolerance grade. If it increases from 1 to 10 000 mm the roughness is permitted to increase by a factor of 40. Since the numerical values correspond to the preferred number series R5 there is an increase of 60% from one tolerance grade to the next nominal diameter step. Therefore the allowed total height increases to the hundredfold value between tolerance grades 4 and 14. The upper $R_t$ limit is 250 $\mu$m for all basic sizes and tolerance grades.

The roughness grade numbers are defined in ISO 1302 or ONORM M.1115.

Another way of looking at tolerances and surface geometry is to avoid unnecessarily high requirements for surface roughness. This in itself is worth doing because it cuts down processes and reduces work stages [303]. Quite often the problem of tolerance versus roughness is more one of lack of communication than that of a technical criterion. If roughness is left off a drawing it is usually added as a precaution or put on at a later stage by a department other than development (e.g. the inspection department).

A relationship between roughness and tolerance has been based on the work of Kohlhage and Opitz and VDI3219. One practical empirical result is shown in table 7.11. It can be seen that the allowed roughness increases with both the size and the ISO tolerance grade.

The issue is that roughness is sometimes not specified or specified too tightly. Occasionally the roughness value which can be included is subject to other constraints which do not allow the preferred values to be kept to. However, the tables have been compiled from thousands of practical cases by Osanna and Totewa and therefore represent usable realistic values. One other constraint which has to be included is that of position (i.e. geometric deviations of position) [327]. The problem of the interaction between surface roughness and form also affects the tolerance of size by itself. A nomogram showing the relationship has been proposed for ease of use and is based on ONORM M.1116 [260], also [328, 329,].

It can be shown [260] that as a guide a working scheme would be that for workpieces with dimensions within tolerance; form deviations and roughness values (measured as profile height) can be allowed to have approximately equal magnitude unless higher or smaller numerical values are specified because of functional requirements.

To specify and collect the relevant data for compiling nomograms it has been necessary to wait for the skilled use of CMMs (coordinate-measuring machines). Lack of suitable data has held back the compilation of guides for many years, but now such data correlating roughnesses, form and size tolerances is being incorporated into computer-aided design systems.

As a rough guide as to how this has been done, pairs of data for roughness and form, $R$ and $F$, were collected on parts over a large number of plants and machines. This gives samples $(R_1, F_1)$, $(R_2, F_2)$,..., $(R_{n,Fn})$ of a two-dimensional population for analysis purposes with $n \geqslant 5$. On the basis of the sample a linear regression is tried:

$$e = \sum_{j=1}^{n} (F_j - bR_j - k)^2 = (n-1)(s_1^2 - b^2 s_2^2) \tag{7.354}$$

where $e$ is made a minimum and the regression gives $b = S_{RF}/s_1^2$, and where $s_1$ and $s_2$ are the sample standard deviations and $S_{RF}$ is the covariance of the sample (correlation).

From the regression coefficient a confidence interval for the regression coefficient $\beta$ of the parent population is determined.

**Table 7.11** Maximum allowable total height $R_t$ ($R_{max}$) in micrometres.

| Nominal diameter steps (mm) | ISO tolerance grades | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| 1–3 | 1 | 1.6 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | | |
| 3–6 | 1 | 1.6 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | | |
| 6–10 | 1.6 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | | |
| 10–18 | 1.6 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | | |
| 18–30 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | | |
| 30–50 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | | |
| 50–80 | 2.5 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | | |
| 80–120 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | | |
| 120–180 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | | |
| 180–250 | 4 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | | |
| 250–315 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | | |
| 315–400 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | | |
| 400–500 | 6.3 | 10 | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | | |
| 500–630 | | | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 |
| 630–800 | | | 16 | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 |
| 800–1000 | | | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 |
| 1000–1250 | | | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 |
| 1250–1600 | | | 25 | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 |
| 1600–2000 | | | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 2000–2500 | | | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 2500–3150 | | | 40 | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 3150–4000 | | | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 4000–5000 | | | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 5000–6300 | | | 63 | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 6300–8000 | | | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |
| 8000–10000 | | | 100 | 160 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 | 250 |

Assuming Gaussian random statistics, the mean

$$\mu\langle R_j \rangle = \beta R_j + K \tag{7.355}$$

and the variance $\sigma^2$ do not depend on $R$. It has been found that $\log R$ and $\log F$ are most likely [260], although it turns out that this change does not make any signficant change to the result.

The per cent confidence interval is thereby worked out:

$$b - J \leqslant \beta \leqslant b + J. \tag{7.356}$$

There has to be a significance value $\gamma$ for this expression. It is most often taken as being $\gamma = 99\%$. $J$ is given by

$$J = ce^{1/2} / (s_1(n-1)(n-2))^{1/2} \tag{7.357}$$

with values of $c$ quoted from tables of the $t$ distribution because of the relationship

$$F(c) = (\gamma + 1)/2. \qquad (7.358)$$

Knowing $c$ and hence $J$ from equation (7.358), and having evaluated $b$ from the data, means that the confidence intervals in (7.357) can be set.

The advantage of such an approach is that it relies upon the use of known existing data and does not require the use of mathematical models such as those using the beta function.

## 7.10.  Discussion

### 7.10.1   Profile parameters

In this chapter an attempt has been made to collect together the evidence for the functional importance of surface texture. It has emerged that there is very little quantitative information but much indirect evidence.

This is despite the availability of a multiplicity of parameters on modern instruments.

Consider two cases which purport to match parameters with function.

In table 7.12 on the left hand side a number of typical engineering uses are listed. Despite the large amount of work which has been done compiling the table, the actual amount of information it reveals is disappointingly small. It seems for example that $R_a$, $R_{sk}$, $R_{\Delta q}$ and lay are all important for most of the applications listed. Presumably any parameter will do! There is simply no real discrimination. The list conveys the feeling that the surface is important, but no more.

Table 7.13 shows a similar list of applications this time using surface parameters based on the French Motif system as reported in ISO Handbook 33 1997. Some attempt has been made to categorize the function in a way similar to the function map concept developed here.

**Table 7.12**  Typical match of parameter to function key.

| Function | Heights | Distribution and Shape | Slopes and Curvature | Lengths and Peak Spacing | Lay and Lead |
|---|---|---|---|---|---|
| Typical Parameters | $R_s$ $R_q$ $R_t$ | $R_{sk}$ $R_{ku}$ | $R_{\Delta q}$ | $R_{max}$ HSC | Std Sal |
| Bearings | √√ | √√ | √ | √ | √√ |
| Seals | √√ | √√ | √√ | √ | √√ |
| Friction | √√ | √√ | √√ | √√ | √√ |
| Joint stiffness | √√ | √√ | √ | √ | √ |
| Slideways | √√ | √√ | √ | √√ | √√ |
| Contacts (elec/therm) | √√ | √√ | √√ | √√ | |
| Wear | √√ | √√ | √√ | √√ | √√ |
| Galling | √√ | √ | √√ | | |
| Adhesion & bonding | √√ | √√ | √ | √ | √ |
| Plating & painting | √√ | √ | √ | √ | |
| Forming & drawing | √√ | √ | √ | √√ | √ |
| Fatigue | √√ | √ | x | x | √√ |
| Stress & fracture | √√ | x | √ | | √√ |
| Reflectivity | √√ | | √√ | √√ | √√ |
| Hygene | √√ | √ | √ | | |

Key: √√ - much evidence    √ - some evidence    x - little or circumstantial evidence

**Table 7.13** Relation between motif parameters and function of surfaces.

Most important parameter: specify at least one of them.

Secondary parameter: to be specified it necessary according to the part functions.

The indication (0.8). For example, means that if the symbol FG indicated on the drawing, and W not otherwise specified, the upper tolerance on W is equal to the upper tolerance on R multiplied by 0.8.

The symbols (FG, ect.) are acronyms of French designations.

| Surface | | Functions applied to the surface | | Parameters | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | Roughness profile | | | Waviness profile | | | | Primary profile | |
| | | Designations | Symbol | $R$ | $Rz$ | $AR$ | $W$ | $Wx$ | $Wk$ | $AW$ | $P1$ | $Pk$ |
| two parts in contact | with relative displacement | Slipping (lubricated) | FG | ● | | | ≤0.8R | | | O | | ● |
| | | Dry frictions | FS | ● | | O | | ● | | O | | |
| | | Rolling | FR | ● | | | ≤0.3R | ● | | O | | O |
| | | Resistance to hammering | RM | O | | O | O | | | O | | ● |
| | | Fluid friction | FF | ● | | O | | | | O | | |
| | | Dynamic sealing — with gasket | ED | ● | O | O | ≤0.6R | ● | | O | | |
| | | Dynamic sealing — without gasket | | O | ● | | ≤0.6R | | | | | ● |
| | without displacement | Stastic sealing — with gasket | ES | O | ● | | ≤R | | O | O | | |
| | | Stastic sealing — without gasket | | O | ● | | ≤R | | ● | | | |
| | | Adjustment without displacement with stress | AC | O | | | | | | | | ● |
| | | Adherence (bonding) | AD | ● | | | | | | | O | |
| indepen-dent surface | with stress | Tools (cutting surface) | OC | O | | O | ● | | | ● | | |
| | | Fatigue strengths | EA | O | ● | O | | | | | | O |
| | without stress | Corrosion resistance | RC | ● | ● | | | | | | | |
| | | Paint coating | RE | | | O | | | | O | | |
| | | Electrolytic coating | DE | ● | ≤2R | ● | | | | | | |
| | | Measures | ME | ● | | | ≤R | | | | | |
| | | Appearance (aspect) | AS | ● | | O | O | | | O | | |

Although stress is considered, it and the other surface properties are treated as attributes, and not as variables which again pulls back from quantifying the link between texture and function. Notice that even though two bodies are mentioned either directly in table 7.13 or indirectly in table 7.12 no attempt is made to bring them into parameters.

The situation up until now has not been satisfactory for a number of reasons:

1. There are too many similar profile parameters such as $R_a$, $R_q$: $P_{tm}$; $R_z$, $S_m$, HSC etc.
2. There are too many arbitrary areal (3D) parameters.
3. Not enough thought has been given to characterizing the effective texture of the system (usually two) of surfaces within the framework of the function map.
4. The function itself has not been adequately characterized.

### 7.10.2 *Functional types of parameter*

To expect slightly different definitions of parameters to be functionally significant is unrealistic. There is no evidence that a functional situation exists in which the $R_a$ value is critical but the $R_q$ value is not! Similarly with $R_{tm}$ and $R_z$ or $R_{3z}$.

The simplest and most effective approach seems to be to sort the parameters into types. Such a typology could be as follows.

(a) Average estimate of distribution obtained from a profile e.g. $R_a$ or $R_q$.
(b) Extreme estimate of distribution such as $R_y$.
(c) Average extreme estimate such as $R_z$.
(d) Areal (3D) structure such as lay from surface statistics.
(e) Positional variation in statistics of areal distribution.
(f) Freak behaviour not associated with process or machine tool (e.g. a deep scratch) and usually few in number.

Add to this plots of the distributions such as the material ratio curve or the slope distribution and geometrical aspects of unit functional agents such as the radius of curvature of asperities in a contact situation. Of these probably a minimum set is (a) (b) (d) and (f).

The simplest question is whether or not the above typology is good for performance or bad.

This has already been addressed in section 7.1 but is repeated here as Table 7.14.

**Table 7.14** Dual roles of surface metrology and relationship to the statistics.



Splitting parameters into type can considerably reduce the number of parameters. However, usefulness even of the 'types' depends on how closely they can be matched to a function. For example peak curvature is definitely important in contact. Material ratio important in load carrying capacity. This realization is tending to redirect parameter specification.

It is because of the difficulty of getting functional specifications that has supported or even extended the 'parameter rash', originally this term was used for profile parameters but now unfortunately, can be used to describe the growth of areal (3D) parameters.

A classic functional parameter described in detail in chapter 2 is the material ratio curve and parameters derived from it. These parameters are shown in Figure 7.162.

The first thing to notice is the amount of construction which is necessary to get the values of the parameters. As described in chapter 2 the idea is that the curve can be segmented into three parts; the top for peaks the middle for core or kernel and the bottom for valleys. Every step is arbitrary! In fact if the 'sleeping ess' curve is plotted on probability paper rather than linear the curve disappears and results in a straight negatively tilted line.

If the layered format is accepted then the three parameters which result, $R_{kp}$ $R_{kv}$ and $R_k$ are deemed to be functionally useful. $R_{kp}$, is the amount which is removed quickly if the surface is rubbed. $R_k$ determines life of the part and $R_{kv}$ indicates the oil retention and debris capacity of the surface. A surface made by plateau honing (i.e. cylinder bores) is measured extensively using this method.

### 7.10.3  *Areal (3D) parameters*

These are split into two groups. The one group is an extension of the profile parameters and the other group is called 'functional parameters'. Here only one or two of each will be given. More information is in chapter 2.

**Figure 7.163**

*Comments on areal parameters*

Despite its slow start there have been some attempts to parameterize the whole surface. This effort has been led by Stout who has been active in European Standards committees and has produced some useful if interim guidelines. Basically his approach has been to use as many areal equivalent profile parameters to the areal equivalent. Obviously this helps to reduce confusion but some less obvious parameters he has suggested called 'functional parameters' are based on common sense rather than proven usefulness. Some of these will be outlined below to illustrate the problems involved. A first important step has been to agree to designate areal parameters by the letter $S$ rather than $R$ for profile parameters, so that $S_a$ is the equivalent of $R_a$ and so one. The 'set' comprises:

(1) Amplitude parameters
(2) Spatial parameters        } as in the profile
(3) Hybrid
(4) Area and volume parameters
(5) Functional parameters

Table 7.15 gives the symbols for the 'conventional' parameters.

**Table 7.15** Conventional parameters—Eur 15178EN provisional.

RMS of surface $S_q$
Ten point height $S_z$
Skew $S_{sk}$
Kurtosis $S_{ku}$
Density of summits $S_{ds}$
Texture aspect ratio $S_n$
Texture direction $S_d$
Shortest correlation length $S_{ai}$
RMS slope $S_{\Delta q}$
Mean summit curvature $S_{sc}$
Material surface area ration $S_{ar}$

These are defined below.
*Amplitude parameters*

1. Arithmetic average $S_a$ corresponding to $R_a$

$$S_a = \frac{1}{L_1 L_2} \int_0^{L_1} \int_0^{L_2} \left| f(x,y) - \bar{f} \right| \, dx \, dy \qquad (7.359)$$

where $\bar{f}$ is the height of the mean plane and $L_1$ and $L_2$ are the extent of the sample. $f(x, y)$ is the surface height at $x,y$. Also $|\ |$ indicates that the sign is ignored.

2. Root mean square ($R_q$) value in 3D areally is $S_q$

$$S_q = \sqrt{\frac{1}{L_1 L_2} \int_0^{L_1} \int_0^{L_2} \left( f(x,y) - \bar{f} \right)^2 \, dx \, dy} \qquad (7.360)$$

3. Skew $S_{sk}$

$$S_{sk} = \frac{1}{L_1 L_2 S_q^3} \int_0^{L_1} \int_0^{L_2} \left( f(x,y) - \bar{f} \right)^3 \, dx \, dy \qquad (7.361)$$

4. Kurtosis $S_{ku}$

$$S_{ku} = \frac{1}{L_1 L_2 S_q^4} \int_0^{L_1} \int_0^{L_2} \left( f(x,y) - \bar{f} \right)^4 \; dx \; dy \tag{7.362}$$

5. Ten point height $S_z$

$$S_z = \left( \sum_{L=1}^{5} \left| P_{i\,max} \right| + \sum_{L=1}^{5} \left| V_{i\,max} \right| \right) \Big/ 5 \tag{7.363}$$

*Justification for functional parameters*

Some advantages of using functional parameters are given below. They are subjective and so need to be validated. They have not been adopted as yet:

  (i)  S notation is usually used for stressed surfaces.
 (ii)  Functional parameters can be specifically matched to a given use.
(iii)  Functional software is easier to understand?

Item (iii) above is stated but not quantified.

Table 7.16 gives some of the functional parameters suggested. They are loosely tied to the material ratio curve which makes sense because it links directly back to the profile.

**Table 7.16** Functional parameters.

Surface material ratio $S_q$
Void volume ratio $S_{vr}$
Surface bearing index $S_{bi}$
Core fluid retention $S_{ci}$
Valley fluid retention $S_{vi}$

**Table 7.17** Additional possibilities some functional parameters.

$S_{bc} = S_q / Z_{0.015}$
Where $Z_{0.015}$ is the height of the surface at 5% material ratio.
Core fluid retention index
$S_{ci} = (V_v(h = 0.8)) / S_q$ (unit area) where $V$ is valley
If $S_{ci}$ is large then there is good fluid retention.
Relationship is $0 < S_c < 0.95 - (h_{0.05} - h_{0.8})$
Valley fluid retention index
$(S_{vi} (V_v(h) = 0.8)) / Sq$ (unit area)
$0 < S_{vi} < 0.2 - (h_{9.8} h_{0.05})$

**Table 7.18** lsotropy.

This can be found by:
$\underline{Longest\ bandwidth} = 1 / isotropy\ ratio$
Shortest bandwidth
$\underline{Shortest\ correlation\ length} = isotropy\ ratio$
Longest correlation length
$\underline{Shortest\ average\ wavelength} = isotropy\ ratio$
Longest average wavelength
$\underline{Shortest\ S_m} = isotropy\ ratio$
Longest $S_m$
The isotropy wave property of the surface has also been called the texture aspect ratio.

Table 7.15 lists a few of the parameters and Table 7.16 gives some of the mathematical formulae for such parameters.

Table 7.16 gives the names of some of the functional parameters suggested and Table 7.17 some of the derivations.

Table 7.18 shows four possible ways of defining isotropy—lay or the lack of it. Many other attempts have been made to define the areal properties. Peklenik and Kubo, Longuet–Higgins etc are but a few. The problems encountered in areal situations are many times more difficult than for the profile. Differences in definitions are possible even for a simple parameter like isotropy. All of these can be related if the surfaces are random. It could be argued that any of these definitions would be sufficient in practice.

To be consistent each parameter should itself be areal (i.e. a plot of any of these isotropic possibilities).

Figure 7.164 and 7.165 give some of the graphical background to a conventional parameter (material ratio) and a suggested functional parameter (void volume ratio $S_{vr}(h)$)

Table 7.19 shows some of the parameters suggested up to now. This table shows a conservative set of parameters and even this has 17 in number.

Stout has, with some difficulty, reduced this to a basic set of 14 possibilities but even this number poses problems.

In figure 7.164 the equivalent of the material ratio is considered and figure 7.165 is an example of a function parameter.

$$\left.\begin{array}{l} \textit{Void volume ratio} \ = \dfrac{V_v(h)}{V_v(h_{\max})} = S_{vr}(h) \\[2mm] \textit{Similarly } V_m = \textit{Volume (material)} = S_{mr} \end{array}\right\}$$

Notice that these are not independent!



**Figure 7.164** Bearing ratio (material ratio).

In sections 7.10.1 and 7.10.3 some profile parameters and areal (3D) parameters have been examined to see their relevance to function. The ideas was to link key parameters to function. What has emerged has not been very satisfactory. Most of the evidence has been indirect. It seems that the most convincing results have occurred when the derivation of the parameter has taken some account of the function. Examples of this are peak curvature for contact or material ratio for rubbing.

The biggest problem has always been to determine the objective in terms of function. To help this the concept of the 'function map' has been introduced. Figure 7.1 and other derivatives figures in the chapter (e.g. 7.163, 7.164). Although simplistic in form it does allow some form of targeting of the behaviour to occur.

**Figure 7.165** Functional parameter.

**Table 7.19** This is a list of some proposed areal (3D) parameters EUR 15178EN from which a minimum set (shown ticked) is indicated

|   | Amplitude parameters | Possible Symbol | In minimum set? |
|---|---|---|---|
| 1 | Root mean square | $S_q$ | √ |
| 2 | Ten point height | $S_z$ | √ |
| 3 | Skew | $S_{sk}$ | √ |
| 4 | Kurtosis | $S_{ku}$ | √ |
|   | Spatial parameters |  |  |
| 1 | Density of summits | $S_{ds}$ | √ |
| 2 | Texture aspect ratio | $S_{tr}$ | √ |
| 3 | Texture direction | $S_d$ | √ |
| 4 | Shortest correlation length | $S_{ae}$ | √ |
|   | Hybrid parameters |  |  |
| 1 | Root mean square slope | $S_{\Delta°v}$ | √ |
| 2 | Mean summit curvature | $S_{sc}$ | √ |
| 3 | Developed surface area ratio | $S_{ar}$ | √ |
|   | Some functional parameters |  |  |
| 1 | Surface bearing area ratio | $\phi$ | — |
| 2 | Void volume ratio | $S_{vr}$ | — |
| 3 | Material volume ratio | $S_{mr}$ | — |
| 4 | Surface bearing index | $S_{li}$ | √ |
| 5 | Core fluid retention | $S_{ci}$ | √ |
| 6 | Valley fluid retention | $S_{vi}$ | √ |

Note: This is not a standard. It is just one attempt to put together a set of parameters to be measured over on cover of the surface. The idea is to select from the minimum set of 14 perhaps one or two from each group.

### 7.10.4 *Function maps and surface parameters*

Simply drawing two axes representing some functional situations is not enough. The map has to be localized into fundamental phenomena. Figure 7.167 breaks up the map into two flow regimes, one normal and one lateral. Surface finish is critical in the vertical flow of, say, electricity or heat because it determines the number and shape of contacts. This, in turn is mainly a manufacturing process problem as has been explained in chapter 6. The texture is not quite as important in the transverse direction because flow can occur despite the texture. However, there is a quantifiable relationship between the areal (3D) pattern and both the load carrying capacity and friction.

**Figure 7.166** Metrology map → function map.

This force regime can be broken down as shown in figure 7.167. On this lines can be drawn representing range boundaries within the constraints of being a function of a spacing and an independent velocity, as in Reynolds's equation for pressure gradient.

The axes can be made dimensionless by referring the gap to $\sqrt{\sigma_\sigma^2 + \sigma_2^2}$ and the velocity relative to the speed of a reference length per second on either or both surfaces. Alternatively the functions could be made dimensionally homogeneous. Figure 7.168 comprises non-dimensional contours of the various regimes. This graph shows the general idea and is not definitive. However, the contours are joint boundaries of the various solutions to the differential equations describing the dynamics of the surface system.

Another possibility is to use energy equations. It may be that the use of scalars could simplify the characterization. One example could be to have the abscissa as the Lagrangian and the ordinate as the Hamiltonian, suitably modified to non-particle behaviour. In this the ordinate would comprise first order equations whilst the abscissa would be second order to encompass damping losses. It remains to be seen which of the various methods of characterizing the function map will be the most useful. The objective once the boundaries are established is to determine those surface parameters which fit into the newly defined regions.



**Figure 7.167** Function map flow regimes.

The concept of pinpointing the function in the ways outlined above represents a shift of emphasis from manufacture towards function or performance. To some extent the manufacturing requirements of the surface has masked the more relevant functional parameters. Ideally, in order to be of use, the definition and operation should mimic the function. Parameters such as $R_a$ are simply descriptions of the waveform. Another departure from the traditional way to characterize the surface, in what follows, will be made in which the parameter includes the functional variable, gap and relative velocity.

**Figure 7.168** Force boundary domain.

### 7.10.5  System approach

The chapter has been trying to tie in surface parameters to performance. In this the investigation has been usually confined to single surfaces. In chapter 2 on characterization it has been suggested that perhaps the characterization of surfaces should be revised to include the system of surfaces. It has already been mentioned many times in this book that the metrology should wherever possible follow the function; the measurement should anticipate performance. So what is needed is an operation on the data to produce the parameter which



**Figure 7.169** Force regime map.

should mimic the actual function of the surface. In the simple method advocated here of using function maps it means devising operations which include the normal gap, the relative velocity of the surfaces and the surface data or some acceptable approximate to it e.g. a surface parameter. Can the system exemplified in the function map form the basis for surface system parameters rather than surface parameters? In what follows a possibility will be given based on convolution methods.

According to the systems philosophy associated with function maps, a functional surface parameter should embody the gap, the relative lateral velocity between the surfaces and some parameter associated with the system of surfaces. Ideally this system parameter should involve an operator to image the actual function.

**Figure 7.170** Function map.



**Figure 7.171** General surface relevance.



**Figure 7.172** Surface parameter map.

For the case of sliding with or without oil or air film the basic equation has been developed to take into account all of the above requirements.

The starting point for such a description is to define a system equivalent to the surface profile or contour. For the normal separation, the addition of the surface geometries is an acceptable system profile.

$$\text{Thus} \qquad z_s = z_1 + z_2 + g \qquad (7.364)$$

Where the gap is measured from the two mean lines $m_1$ and $m_2$. Once accepted, the following parameters for the gap axis, developed from conventional profiles, are given below. The suffix $s$ indicates system.

So

(i) $R_{qs}^2 = R_{q1}^2 + R_{q2}^2$ (Root mean value) (7.365)

(ii) $R_{sks} = R_{sk1}R_{q1}^3 + R_{sk2}R_{q2}^3$ (System skew)

If $R_{q1} = R_{q2} = 1$

$R_{sks} = R_{sk1} + R_{sk2}$ (7.366)

(iii) $R_{kus} = R_{ku1} \cdot R_{q1}^4 + 6R_{q1}^2 R_{q2}^2 + R_{ku2}R_{q2}^4$ (System kurtosis)

If $R_{q1} = R_{q2} = 1$ and the surfaces are Gaussian

$R_{kus} = R_{ku1} + R_{ku2}$ (7.367)

(iv) Parameters involving peaks or valleys in general are not additive because there is no guarantee that peaks on one surface have the same spatial address as those on the other surface. So $R_2$, $R_{tm}$, $R_y$ etc are not acceptable parameters.

(v) The material ratio $MR_s$ can be found readily because the system profile $z_s = z_1 + z_2 + g$ which corresponds to the convolution process linking the material ratio values of $MR_1$ and $MR_2$.

$$MR_s(z) = \int_{-\infty}^{\infty} MR1(s).MR2(z - s) \ ds \qquad (7.368)$$

$$MR_s = MR1 * MR2 \qquad (7.369)$$

where the symbol * denotes convolution. $MR_s$ is in effect a composite material ratio of the system. Terms involving individual $MR_1$ and $MR_2$ are left out because they do not contribute to the interaction. See chapter 2.

The above parameters possess the critical height information of the system.

Wavelength information can also be developed for the system in static mode.

Thus

$$A_s(\tau) = A_1(\tau) + A_2(\tau) \qquad (7.370)$$

where $A(0)$ is the autocorrelation function and $\tau$ is the shift. When normalized with respect to $A_s(\omega)$ the origin at $\tau = 0$ is unity, otherwise it is as $A_0(0) = R_q^2$.

There is a similar linear relationship between the power spectral density $P_s(\omega) = P_1(w) + P_2(\omega)$. Each spectrum or correlation function is insensitive to the spatial position $x$ and $y$.

The term $MR_s(z)$ is particularly important because it is a measure of the interaction which is allowed between the surfaces in the vertical ($z$) direction. This situation has been called 'static' in the function map because there is no lateral movement. However, the system need not be static in the $z$ direction. The $z$ value

can be a function of time to take into account functions which involve a fluctuating vertical position such as squeeze film operation and also the rolling condition. In this, although the contact lateral velocity is zero, the approach of the two surfaces is cyclic.

The system approach so far has dwelt on the behaviour of $z_1(x) + z_2(x)$ in terms of the probability density or the material ratio and in terms of normal average properties. There has been no dynamic aspect. This can be taken into account by investigating $z$ values having different $x$ values.

The question arises as to which are the most useful dynamic parameters. These must be the dynamic separation and the dynamic slope or wedge. It is plausible therefore to investigate the properties of $g-z(x)$ i.e. $z_1(x) + z_2(x_2)$ and $z_1(x_1) + z_2(x_2)$ by their mean square values.

$$\text{Thus} \quad \left. \begin{array}{l} E(z_1(x_1) + z_2(x_2))^2 = E(z_2)^2 + E(z_2)^2 - 2E(z_1(x_1) \cdot z_2(x_2)) \\ E(z_1(x_1) - z_2(x_2))^2 = E(z_2)^2 + E(z_2)^2 - 2E(z_1(x_1) \cdot z_2(x_2)) \end{array} \right\}. \tag{7.371}$$

The interactive component of both estimates is $2E(z_1(x_1) \cdot z_2(x_2))$ over all $x_1$ and $x_2$.

In order to mimic dynamic behaviour $x_1$ and $x_2$ have to be related. Perhaps the obvious would be the arguement of the cross correlation $x$ and $(x+\tau)$. In this case as the shift $\tau$ is advanced there would be relative positioning of $z_1$ and $z_2$. However, this does not readily translate into a dynamic relative velocity suitable to act as the abscissa of the function map. A much more promising possibility is to use the folding (faltung) property of convolution rather than correlation. This implies using $x$ and $(\tau-x)$.

So the interaction becomes $2E(z_1(x_1) \cdot z_2(\tau-x))$ rather than $2E(z_1(x_1) \cdot z_2(\tau+x))$ \hfill (7.372)

It should be pointed out here that $\tau$ for convolution is not the same as the $\tau$ for correlation. Also some liberty has been taken here by treating the convolution process as an average like the correlation rather than a simple integral for convolution.

Thus cross convolution $CCon(\tau)$ is given by

$$CCon(\tau) = \frac{1}{L} \int_0^L z_1(x) \cdot z_2(\tau - x)\, dx. \tag{7.373}$$

Equation (7.373) provides the necessary scan of $z_1$ across $z_2$ for each $\tau$ i.e. $z_1 * z_2$.
In the Fourier transform plane $CCon(\tau)$ becomes

$$F_{12}(w) = F_1(w) \cdot F_2(w) = A_1(w)A_2(w)\exp(j(\varphi_1(w) + \varphi_2(w))). \tag{7.374}$$

Compare this with the cross correlation transform $C_{12}(w)$.

$$C_{12}(w) = F_1(w) \cdot F_2^*(w) = A_1(w)A_2(w)\exp(j(\varphi_1(w) - \varphi_2(w))). \tag{7.375}$$

Taking $A_1 \cdot A_2$ as the power spectrum $P_{12}(w)$ it can be seen that the effective difference between the two approaches is that cross convolution preserves spatial phase and cross correlation does not! As both contain $P_{12}(w)$ it is clear that the essential surface–system interaction is the phase term. This is usually regarded as a source of variation and removed by dealing with correlations. In this form of system approach the opposite is true.

The actual variation of $z(x_1).z(x_2)$ corresponds to the system profile and parameters can be invented to describe these variations such as $R_a$ (system) in the same way as a surface profile is treated.

By representing $x$ by $vt$ where $v$ is the relative velocity of the two surfaces it is now possible to represent the system in terms of gap '$g$' and velocity '$v$' —the real functional parameters — as well as topographic properties such as $R_a$(system). Because the static gap value influences the value of the dynamic effect, the

system characterization can be expressed neatly by the product of the two convolutions.

$$Syst(g.v) = (MR_1 * MR_2) . (z_1 * z_2) \tag{7.376}$$

or really by $(MR_1(z)*MR_2(z)).(z_1(x)*z_2(w)).(z_1(y)*z_2(y))$ (7.377)

shown symbolically in figure 7.173.

Equation (7.377) represents an operator working on the data of both surfaces $z_1$ and $z_2$ of the system and including the functional parameters $g$ and $v$, both of which have a temporal equivalent.

So the interesting point is that the system involving the surfaces and functional parameters can be described by 'across' convolution operations in $x, y$ and $z$.

The operation can be represented conveniently by the across diagram (figure 7.173) where the centre of the drawing is the convolution symbol.

Most of the tribological regimes of the function map can be accommodated by this operator.



(a) Movement in $x$ direction     (b) Movement in $x$ and $y$ direction

**Figure 7.173** A system approach to a surface function operation.

## 7.11 Conclusions

The chapter is in two parts. The first part deals with the relationships between the conventional parameters and function (performance). It becomes obvious that one reason for the disappointing correlation is that the expectation is too high. For example many height parameters have just the same effect on performance. It turns out that the best approach is to classify the surface parameters into 'types' such as the average value of the distribution of the surface or the extremes of the distribution.

This condensed set gives a much more realistic correlation than any individual parameter. Which parameter to use i.e. $R_a$ or $R_q$ reduces to a matter of choice or convenience: providing it is within the 'type' it is usually satisfactory.

In order to focus properly on function a classification of function—here called a function map—has been introduced. This is based on two simple axes: one is the gap between surfaces and the other is their relative lateral velocity. Before classifying the surface parameter, the function must be clarified. It turns out that many tribological and other applications fit comfortably within this concept.

In an effort to reduce the 'parameter rash', which now runs into a minimum of 14 parameters per surface, a closer look between the surface and the application has been given.

This has resulted in a new approach, in which the parameter becomes more of an operator which produces an output value linked to the two surfaces: the operation embodies the functional mechanism.

This gives a parameter value dependent not only on the two surface roughnesses but also on the functional parameters 'gap' and velocity. It therefore becomes a true surface system parameter in which the number of parameters is reduced, from dozens for each of the surfaces to a few for the system. The mechanism identified to satisfy the surface system approach is cross convolution. This is called here 'across convolution', which can

embody the system parameters i.e. the gap and relative velocity as well as the two surfaces. Using this mathematical operation should increase the fidelity of the parameters as well as drastically reducing the number of parameters. The only real problem is that of having access to both surfaces for measurement rather than just one.

# References

[1] Whitehouse D J 2001 Function maps and the role of surfaces *Int. J. Mech. Tools Manuf.* **41** 1847–61
[2] Thomas T R 1999 Rough surfaces (*Imperial College Press*)
[3] Williamson J B P 1967/68 Microtopography of solid surfaces *Proc. IMechE* **180** 21–30
[4] Love A E H 1944 *A Treatise on the Mathematical Theory of Elasticity* (1892) (Cambridge: Cambridge University Press)
[5] Prescott J 1961 *Applied Elasticity* (1924) (London: Longmans Green)
[6] Hertz H 1881 On the contact of elastic bodies *J. reine Angew. Math.* **92** 156
[7] Kellog O D 1929 *Foundations of Potential Theory* (New York: Murray)
[8] Puttock M J and Thwaite E O 1969 *National Standards Lab. Tech. Paper No* 25, (Melbourne: CSIRO)
[9] Pashley M D 1984 Further considerations of the DMT model for elastic contact *Colloids Surf* **12** 69–77
[10] Greenwood J A and Williamson J B P 1964 The contact of nominally flat surfaces *Burndy Res. Rep.* No. 15
[11] Whitehouse D J and Archard J F 1970 *Proc. R. Soc.* **A316** 97
[12] Nayak P R 1973 *Wear* **26** 305
[13] Majumdar A and Tien C L 1990 Fractal characterization and simulation of rough surfaces *Wear* **136** 313–327
[14] Majundar A and Bhushan B 1991 Fractal model of elastoplastic contact between rough surfaces *Trans. ASME Journal of Tribology* **113** 1–11
[15] Sayles R S 1978 Computer simulation of the contact of rough surfaces. *Wear* **49** 273–296
[16] Archard J F 1957 Elastic deformation and the laws of friction *Proc. R. Soc.* **A243** 190–205
[17] Dyson J and Hirst W 1954 The true contact area between solids *Proc. Phys. Soc.* **B67** 309–12
[18] Clavarella M and Demelio G 2001 Elastic multiscale contact of rough surfaces: Archord's model revisited and comparisons with modern fractal models *J. Appl. Mechanical Trans ASME* **68** 496–498
[19] Clavarella M, Demelio G Borber Jr and Jang V H Linear elastic contact of the weterstross profile *Proc. Roy Soc land* A. **456** 387–405
[20] Greenwood J A and Tripp J H 1967 The elastic contact of rough spheres. *ASME J. App Mech* **34** 53
[21] Majunder A and Bhushan B 1945 Characterization and modelling of surface roughness and contact mechanic *Handbook of Micro/nano tribology* (Boca Raton: CRC Press) 109–165
[22] Whitehouse D J 2001 Fractal or Fiction *Wear* **249** 345–353
[23] Greenwood J A 1967 The area of contact between rough surfaces and flats *Trans. ASME, J. Lubr. Technol.* **89F** 81–9
[24] Greenwood J A and Trip J H The elastic contact of rough surfaces *Trans. ASME, J. Appl. Mech.* **34E** 153–9
[25] Greenwood J A *Tribology Rough Surfaces* ed T R Thomas (London: Longman) p 191
[26] Greenwood J A, Johnson K L and Matsubara E 1984 A surface roughness parameter in Hertz contact *Wear* 47–57
[27] Mikic B B and Roca R T 1974 *J. Heat Mass Transfer* **17** 205
[28] Yip F C and Venart J E S 1971 An elastic analysis of the deformation of rough spheres, rough cylinders and rough annuli in contact *J. Phys. D: Appl. Phys.* **4** 1470–80
[29] Rossmanith H P 1976 Stochastic analysis of non top top–contacting rough surfaces *Wear* **37** 201–8
[30] Zhuravlev V A 1940 On the physical basis of the Amontons Coulomb law of friction *J Tech Phys (USSR)* **10** 1447
[31] Greenwood J A and Williamson J B P 1964 The contact of nominally flat surfaces *Burndy Res. Rep. No.* 15
[32] Bickel E 1963 Some fundamental problems in the measurement of surfaces *Proc. mt. Conf. Prod. Eng. Res., Pittsburgh*
[33] Tallian T E, Chui Y P, Huttenlocher D F, Kamanshma J A, Sibley L B and Sidlinger N E 1964 Lubricant films in rolling contact of rough surfaces *ASLE Trans* **7(2)** 109–26
[34] Pesante M 1964 Determination of surface roughness typology by means of amplitude density curves *Ann CIRP* **12** 61
[35] Lipg F 1958 On asperity distributions of metallic surfaces *J. Appl. Phys.* **29** 1168
[36] Bowden F P and Tabor D 1954 *The Friction and Lubrication of Solids* (Oxford: AppI. Phys. University Press)
[37] Block H 1952 *Proc. R. Soc.* **A212** 480
[38] Halliday J 1955 *Proc. IMechE* **16a** 177
[39] Greenwood, J A and Morales–Espejel G E 1997 The amplitude of the complementary function for wavy EHL contacts in elastohydrodynamic' 96, Proceedings of the 23rd Leeds–Lyon Symposium in Tribology, Eds D Dowson *et al*
[40] Whitehouse D J and Archard J F 1970 *Proc. R. Soc.* **A316** 97
[41] Onions R A and Archard J F 1973 *J Phys. D*: *Appl. Phys.* **6** 289
[42] Gupta P K and Cook N H 1972 *J. Lubr. Technol.* **94** 19
[43] Nayak P R 1973 *Wear* **26** 305

[44] Francis H A 1977 Application of spherical indentation mechanics to reversible and irreversible contact between rough surfaces *Wear* **45** 221–9
[45] Wu Chengwel and Zheng Linqing 1988 A general expression for the plasticity index *Wear* **121** 161–72
[46] Bush A N, Gibson RD and Keogh G P 1978 Strongly anisotropic rough surfaces *ASME,* paper 78, Lubr 16
[47] Bush A N, Gibson RD and Keogh G P 1980 The limit of elastic deformation in the contact of rough surfaces *Mech. Res. Commun.* **3** 169–74
[48] Pullen J and Wilhamson J B 1972 On the plastic contact of rough surfaces *Proc. R. Soc.* **A327** 159–73
[49] Seabra J and Berthe D 1987 Influence of surface waviness and roughness on the normal pressure distribution in the Hertzian contact *Trans ASME* **109** 462–5
[50] Johnson N I 1949 System of frequency curves generated by methods of translation *Biometrika* **36** 149–167
[51] Chilamakuri S R and Bushan B 1998 Contact analysis of non–gaussian random surfaces *Pro. Inst. Mech. Engrs.* **212** 19
[52] Whitehouse D 1978 Beta functions for surface typology *Ann. CIRP* **27** 441
[53] McCool J I 1992 Non–Gaussian effects in micro contact *Int. Journal Mech. Tools Manuf.* **32** 115–123
[54] Adler R J and Firman D 1981 A non Gaussian model for random surfaces *Phil. Trans. Roy. Soc. A* **303** 433
[55] Oden P I, Maumdar A, Bhushan B, Padmanabhan A and Graham J J 1992 AFM imaging, roughness analysis and contact mechanics of magnetic tapes and head surfaces *Trans. ASME J. of Trib.* **114** p 666
[56] Russ J C 1994 *Fractal Surfaces* (New York: Plenum Press )
[57] Hasegawa M, Lui J Okuda K and Nunobiki M 1996 Calculations for the fractal dimensions of machined surface profiles *Wear* **193** 40–45
[58] He L and Zhu J 1997 The fractal character of processed metal surfaces *Wear* **208** 17–24
[59] Zhou G Y, Leu M C and Blackmore D 1993 Fractal model for wear prediction *Wear* **170** 1–14
[60] Gupta P K and Walowit J A 1974 Contact stresses between an elastic cylinder and a layered elastic solid *Trans. ASME J. Lubric. Tech.* **96** 250
[61] Volver A 1997 Correlation factors for the 2D coated Hertzian contact problem *Wear* **212** 265
[62] Cole S J and Sayles R S 1991 A numerical method for the contact of layered elastic bodies with real rough surfaces *J. Tribol.* **114** 729
[63] Chang W R An elastic plastic contact model for a rough surface with an ion plated soft metallic contact.
[64] Johnson K L 1985 *Contact Mechanics* (Cambridge: Cambridge University Press)
[65] Thomas T R and Sayles R 1977 Stiffness of machine tool joints—a random process approach *Trans. ASME, J. Eng. md.* **99B** 250–6
[66] Webster M N and Sayles R C 1986 A numerical model for the elastic frictionless contact of real rough surfaces *Trans. ASME, J. Tribol.* **108** 315
[67] Sayles R S and Bailey D M 1987 The modelling of asperity contacts *Tribology in Paniculate Technology* ed B J Briscoe and M J Adams (Bristol: Hilger)
[68] Schofield R D and Thornley R H 1972 Calculating the elastic and plastic components of deflection of plane joints formed from machined surfaces *Proc. 12th MTDR Conf.* (London: MacMillan)
[69] Thomley R H, Connolly R and Koenigsberger F 1967/68 The effect of flatness of joint faces on the static stiffness of machined tool joints *Proc. IMechE* **182** 18
[70] Sherif H A 1991 Parameters affecting the contact stiffness of nominally flat surfaces *Wear* **145** 113–21
[71] Nagaraj H S 1984 Elasto plastic contact of bodies under normal and tangential loading *ASME, J. Tribology* **106** 519
[72] Tsukizoe T and Hisakado T 1965 On the mechanism of contact between metal surfaces. The penetrating depth and the average clearance *J. Basic Eng. Trans. ASME* **87 D** 666
[73] Mitchell L A and Rowe M D 1967/68 Assessment of face seal performance based on the parameters of a statistical representation of surface roughness *Proc. IMechE* **182** 101
[74] Mitchell L A and Rowe M D 1969 Influence of asperity deformation made on gas leakage between contacting surfaces *J. Mech. Eng. Sci.* **11** 5
[75] George A 1976 *CEGS Rep.* November
[76] Johnson K L, Kendal K and Roberts A D 1971 Surface energy and the contact of elastic solids *Proc. R. Soc.* **A 324** 301–13
[77] Fuller K N G and Tabor D 1975 The effect of surface roughness on the adhesion of elastic solids *Proc. R. Soc.* **A 345** 327–2
[78] Whitworth K 1840 British Association, Glasgow
[79] Tyndall J 1875 *Proc. R. Soc. Inst.* **7** 525
[80] Derjaguin B V, Muller V M and Toporov Y P 1975 Effect of contact deformation on the adhesion of particles *J. Colloid Interface Sci.* **53** 31–26
[81] Muller V M, Yushchenko VS and Deriaguin B V 1980 On the influence of molecular forces on the deformation of an elastic sphere and its sticking to a rigid plate *J. Colloid Interface Sci.* **77** 91–101
[82] Pashley M D and Pethica J B 1985 The role of surface forces in metal–metal contacts *J. Vac. Sci. Technol.* **A 3** 757–61
[83] Maugis D and Pollock H M 1984 Surface forces deformation and adherence at metal microcontacts *Acta. Metall.* **32** 1323–34
[84] Pashley M D, Pethica J B and Tabor D 1984 Adhesion and nucromechamcal properties of metal surfaces *Wear* **100** 7–31

[85] Chang W R, Etsion I and Bogy D B 1988 Adhesion model for metallic rough surfaces *J. Tribol. Trans ASME* **110** 50~2

[86] Chang W R, Etsian I and Bogy D B 1987 Elastic plastic model for the contact of rough surfaces *ASME, J. Tribol.* **109** 257–63

[87] Muller V M, Derjaguin B V and Toporov Y P 1983 On two methods of calculating the force of sticking of an elastic sphere to a rigid plane *Colloids Surf.* **7** 251–9

[88] Greenwood J A 1997 Adhesion of elastic spheres *Proc. R Soc. LoM. A* **453** p 277

[89] Bradley R S 1932 The cohesive forces between solid surfaces and the surface energy of solids *Phil. Mag.* **13** 853

[90] Bowden F B and Tabor D 1964 *Friction and Lubrication of Solids, Parts I, II and 111* (Oxford; Clarendon)

[91] Sayles R S and Thomas T R 1976 Thermal conduction of a rough elastic solid *Appl. Energy* 1249–67

[92] Tsukizoe T and Hisakado T 1968 On the mechanism of contact between metal surfaces Part 2 *Trans ASME* **90F** 81

[93] Thomas T R and Probert S D 1965 Thermal resistance of pressed contacts *UKAEA Rep* TRG 1013 (RIX)

[94] Kraghelsky I V and Demkin N B 1960 Contact area of rough surfaces *Wear* **3** 170

[95] Yovanovich M M 1982 Thermal contact correlations, spacecraft radiative transfers and temperature control (ed T E Horton) **83** 83–95 (New York: AIAA)

[96] Yovanovich M M and Nho K 1989 Experimental investigation of heat flow rate and direction on control resistance of ground lapped stainless steel interfaces *24th AIAA Thermophysics Conf, (Buffalo)* paper 89–1657

[97] Majumdar A and Tien C L 1991 Fractal network model for contact conductance *Trans ASME; Journal of Heat Transfer* **113** 516

[98] McWaid T H and Marschall E 1992 Application of the modified Greenwood and Williamson contact model for the prediction of thermal contact resistance *Wear* **152** 263

[99] Lambert M A and Fletcher L S 1995 Thermal contact conductance ofnon flat rough metals in vacuum *ASME/ISME Thermal Engineering join. conf. ASME New York* **1** 31

[100] Torii K and Nishino K 1995 Thermal contact resistance of wavy surfaces *Revista Brasileira de ciencias mecanicas* **17** 56

[101] O'Callaghan P W, Babus' Haq R F and Probert S D 1989 Prediction of contact parameters for thermally distorted pressed joints *24th AIAA Thermophysics Conf, (Buffalo)* paper 89–1659

[102] Song S and Yovanovich M M 1988 Relative contact pressure dependence on surface roughness and Vickers microhardness *J. Thermophys.* **12** 4347

[103] O'Callaghan P W and Probert S D 1987 *Proc. IMech E* **201** 45–55

[104] Cetinkale T N and Rishenden M 1951 Thermal conduction of metal surfaces in contact *Proc. mt. Conf on Heat Transfer (London)* pp 271–5

[105] Jones M H, Howells R I L and Probert S D 1968 Solids in static contact *Wear* **1** 225–40

[106] Stubsted W 1964 Thermal contact resistance between thin plates in vacuo *Eng Rep* (Cedar Rapids, Iowa: Colins Radio Cp)

[107] Aron W and Colombo G 1963 Controlling factors of thermal conductance across bolted joints in vacuum *ASME* paper 63–WA 196

[108] Mikesell E P and Scott R B 1956 Heat conduction through insulating supports in very low temperature equipment *J. Res. NBS* **57** 371–78

[109] Thomas T R 1968 *PhD Thesis* University of Wales.

[110] Holm R 1958 *Electrical Contacts Handbook* 3rd edn (Berlin Springer) (4th edn 1967)

[111] Thomas T R and Probert S D 1971 Correlations for thermal contact conductance in vacuo *HTAA, ASME, J. Heat Transfer* Jan, 1–5

[112] Howells R T L and Probert S D 1968 Contact between solids *TRG Rep.* 1701 (R/X)

[113] Clausing A M 1966 *Int. J. Heat Mass Transfer* **9** 791

[114] Brodie D E 1961 *PhD Thesis* McMaster University, Hamilton, Ontario

[115] Turner M J B 1965 Sliding electrical contacts *IEE, SQJ* pp 67–76

[116] Dowson D 1977 Early concepts of the role of surface topography in tribology *Proc. 4th Int Leeds/Lyon Symp. on Tribology (IMechE)* pp 3–10

[117] Archard J F 1975 Fact and friction *Proc. R. Inst. GB* **48** 183–206

[118] Dowson D 1973 The early history of tribology *First European Tribology Congr. (IMechE)* paper C253/73

[119] Euler 1750 *Histoire de L'academic Roy ale des Sciences et Belle–Lettres* Annee 1748, Berlin 1950, p 122

[120] Coulomb C A 1785 Theory des machines simples *Memoire de Mathematique et de Physique de I'Academie Royale p* 161

[121] Hardy W B 1936 *Collected Works* (Cambridge: Cambridge University Press)

[122] Bowden F B and Leben L 1940 The nature of sliding and the analyses of friction *Phil. Trans. R. Soc. A* 2391

[123] Temoshenko S 1934 *Theory of Elasticity* (New York: McGraw–Hill)

[124] Mansfield D 1999 Profile reconstruction *Taylor Hobson Tech. Note*

[125] Wilson (in [98] Part II, p41)

[126] Quin T F J 1967 *Trans. ASLE* **10** 158

[127] Quin T F J 1968 *Proc. IMechE* **182** 201

[128] Quin T F J 1969 *Proc. IMechE* **183** 124

[129] Childs T H C 1980 *Tribal. Int.* **12** 285

[130] Archard J F 1985 Friction between metal surfaces *Seminar, Friction & Contact Noise, (Delft, June 1985)*

[131] Shahinpoor M and Mohamed MAS 1986 On the effect of asperity pair elasticity joints on friction resistance *Wear* **112** 89–101

[132] Ganghoffer J F and Schultz J 1995 A deductive theory of friction *Wear* 188 88

[133] Moslehy F A 1997 Modelling friction and wear phenomena *Wear* **206** 136

[134] Proctor T D 1993 Slipping accidents in Great Britain an update *Safety Science* **16** 367

[135] Taneeranon P and Yandell W O 1981 Micro texture roughness/effect on predicted road tyre friction in wet conditions *Wear* **69** 321

[136] Wakuri, Hamatake T, Soejima M and Kitari T 1995 Study on the mixed lubrication of piston rings in internal combustion engine *Nippon Kikai Gakkai Ronbu Shu* **61** 1123

[137] Blencoe K A and Williams J A Friction of sliding surfaces carrying boundary fields

[138] Hum B, Colquhoun H W and Lenard J G 1996 Measurements of friction during hot rolling of aluminium strips *J. Mat. Proc. Tech.* **60** 331

[139] Jonasson M, Pulkkenen T, Gunnerson L and Schedh E 1997 Comparative study of shot blasted and electrical discharge textured rolls with regard to frictional behaviour of the rolled steel sheet surfaces *Wear* **207** 34

[140] Saha P K, Wilson W R D and Timsit R S 1996 Influence of surface topography on the frictional characteristics of 3104 aluminium alloy steel *Wear* **197** 123

[141] Skarpelos P and Morris J W 1997 The effect of surface morphology on friction during forming of electrogalvanized sheet steel *Wear* **212** 165

[142] Yevtkushenko A A 1993 non stationary heat formation on the slipping contact of rough surfaces under the condition of combined friction *Dopovidi AN Ukrainy* **11** 51

[143] Yevtkushenko A A and Ivanyk E G 1995 Stochastic contact model for rough frictional heating surfaces in mixed friction conditions *Wear* **188** 49

[144] Michael P C, Rabinowicz E and Iwasa Y Thermal activation in boundary lubricated friction

[145] Fri 15 Azarkhin A and Devenpeck M 1997 Enhanced model of a plowing asperity *Wear* **206** 147

[146] Fri 16 Torrance A A, Galligan J and Liraut G 1997 A model of the friction of a smooth hard surface sliding over a softer one *Wear* **212** 213

[147] Polycarpou A A and Soom A 1995 Application of a 2D model of continuous sliding friction to stick–slip *Wear* **181** 32

[148] Vatta F 1979 On the stick slip phenomenon *Mech. Res. Commun.* **6** 203

[149] van De Velde F and De Baets P 1996 Mathematical approach of the influencing factors on stick–slip induced decelerative motion *Wear* **201** 80

[150] Hirst W and Hollander A E 1974 Surface finish and damage in sliding *Proc. R. Soc.* **A 337** 379–94

[151] Suh N P 1993 Delamination theory of wear *Wear* **125** 111

[152] Poon C Y and Sayles R S 1992 The classification of rough surface contacts in relation to tribology *J. Phys. D: Appl. Phys.* **25** 249–56

[153] Burwell J T 1957/58 Survey of possible wear mechanism *Wear* **1** 119

[154] Kragelskii I V *Friction and Wear* (Engl. Transl. L Randon and J K Lancaster) (Washington, DC: Butterworths) 98

[155] Akamatsu Y, Tsuchima N, Goto T and Hibi K 1992 Influence of surface roughness skewness on rolling contact fatigue life *Trib. Trans.* **35** 745

[156] Dong W P and Stout K J 1995 An integrated approach to the characterization of surface wear 1: Qualitative characterization *Wear* **181–183** 700

[157] Torrance A A 1995 Using profilometry for the quantitative assessment of tribological function P C based software for friction and wear prediction *Wear* **181–183** 397

[158] Rosen B G, Ohlsson R and Thomas T R 1996 Wear of cylinder bar microtopography *Wear* **198** 271

[159] Chandrasekaran T 1993 On the roughness dependence of wear of steels a new approach *J. Mat. Sci. Lett.* **12** 952

[160] Meng H C and Ludema K C 1995 Wear models and predictive equations their form and contact *Wear* **183** 443

[161] Holm R 1946 Electrical contacts *Almqvist and Wiksells Boktrycheri Uppsala*

[162] J F 1953 Contact and rubbing of flat surfaces *J. Appl. Phys.* **24** l2

[163] Barwell F T 1957 Wear of metals *Wear* **1** 317

[164] Podra P and Anderson S 1997 Wear simulation with the Winkler surface model *Wear* **207** 79

[165] Rapoport L 1995 The competing wear mechanisms and wear maps for steel *Wear* **181** 280

[166] Reynolds O 1886 On the theory of lubrication and its application to Mr Beauchamp's towers experiments including an experimental determination of the viscosity of olive oil *Philos. Trans.* **177** 157–244

[167] Grubin A N 1949 *Central scientific research institute for technology and mechanical engineering,* trans Moscow Book no 30, Engl Transl. DSIR 115–6

[168] Elrod H G 1977 A review of theories for the fluid dynamic effects of roughness on laminar lubricating films *Proc. 4th Leeds/Lyon Symp. IMechE* pp 11–26

[169] Tender K 1977 The lubrication of surfaces having a cross striated roughness pattern *Proc. 4th Leeds/Lyon Symp. on Tribology (IMechE)* pp 807

[170] Tender K 1977 Lubrication of surfaces having area distribution isotropic roughness *ASME J. Lubr. Technol.* **99F** 32330

[171] Aouichi A, Biouet J and Vinh J 1981 On numerical signal processing of an evaluative friction force: characterization of successive stages of wear regime in dry sliding of steel on steel *Proc. 8th Leeds/Lyon Symp. on Tribology* (Guildford: Butterworths) pp 77–83

[172] Cameron A 1981 *Basic Lubrication Theory* (Chichester: Wiley)

[173] Christensen H 1969/70 Stochastic models for hydroelliptance lubrication of rough surfaces *Proc. IMechE* **184** 1013–22

[174] Greenwood J A 1982 *Rough Surfaces* (ed T R Thomas) (London: Longmans) p 196

[175] Shukla I B and Kumar S 1977 *Proc. 4th Leeds/Lyon Symp. on Tribology (JMechE)*

[176] Tzeng S T and Saibel E 1967 Surface roughness effect on slider bearing lubrication *ASLE Trans* **10** 334–8

[177] Kuo–Kuang Chen and Dah–Chen Sun 1977 On the statistical treatment of rough surface hydrodynamic lubrication problems *Proc. 4th Leeds–Lyon Symp. on Tribology (IMechE)* p41

[178] Chen K K and Sun D C 1977 On the statistical treatment of rough surface hydrodynamics *Proc. 4th Leeds–Lyon Symp. on Tribology (IMechE)* pp41–5

[179] Rohde S M and Whicker D 1977 Some mathematical aspects of the hydrodynamic lubrication of rough surfaces *Proc. 4th Leeds–Lyon Symp. on Tribology (IMechE)* 32–40

[180] Sun D C and Chen K K First effects of Stokes roughness on hydrodynamic lubrication *Trans. ASME* 99 (1)

[181] Bayada G and Chambat M 1988 New models in the theory of hydrodynamic lubrication of rough surfaces *Trans ASME, J. Tribol.* **110** 402–7

[182] Phan–Thien N 1981 On the effects of the Reynolds and Stokes surface roughness in a two dimensional slide bearing *Proc. R Soc.* A **377** 349–62

[183] Phan-Thien N 1982 Hydrodynamic lubrication of rough surfaces *Proc. R Soc.* **A 383** 439–46

[184] Tonder K 1987 Effects of skew unidirectional straight roughness on hydrodynanuc roughness Part 2 Moving roughness *Trans. ASME, J. Tribol.* **109** 671–8

[185] Seabra J and Berthe D 1988 Elastohydrodynamic point contacts, pt I Formulation of numerical solution *Wear* **130** 301

[186] Dowson D and Higginson G R 1966 *Elastohydrodynamic Lubrication* (Oxford Pergamon), 1959 A numerical solution to the elastohydrodynanuc problem *J. Mech. Eng. Sci.* **1** 615

[187] Pinkus O 1987 The Reynolds Centennial—A brief history of the theory of hydrodynamic lubrication *Trans ASME* **109** 22–15

[188] Dowson D 1977 *History of Tribology* (London: Longmans)

[189] Dowson D 1962 A generalised Reynolds equation for fluid film lubrication *J. Mech. Eng. Sci.* **4** 159–70

[190] Dowson D and Whomes T L 1971 The effect of surface roughness upon the lubrication of rigid cylindrical rollers *Wear* **18** 129–40

[191] Berthe D and Godet H A 1973 More general form of Reynolds equations, application to rough surfaces *Wear* **27** 345–57

[192] Barragan Deling FdM, Evan H P and Sniddle R W 1989 Micro–elasto hydrodynamic lubrication of circumferentially finished rollers the influence of temperature and roughness *Trans. ASME* **111** 7306

[193] Johnson H L and Higginson J G 1988 A non Newtonian effect of sliding in micro–ehi *Wear* **128** 249–64

[194] Sadeghi F and Sui P C 1990 Thermal elastodynamic lubrication of rough surfaces *Trans. ASME J. Tribol.* **112** 3416

[195] Dyson A 1976 The failure of elastohydrodynamic lubrication of circumferentially ground discs *Proc. IMechE* **190** 699–711

[196] Cheng H S and Dyson A 1976 Elastohydrodynanuc lubrication of circumferentially ground rough discs *Trans ASLE* **21** 2540

[197] Zhu D and Cheng H 5 1988 Effect of surface roughness on the point contact ehl *Trans ASME. J. Tribol* 110327

[198] de Silva G M S, Leather J A and Sayles R S 1985 The influence of surface topography on the lubricant film thickness in end point contact *Proc. Leeds/Lyon Conf*

[199] Patir N and Cheng H S 1978 An average flow model for determining effects of 3D roughness on partial hydrodynamic lubrication *Trans. ASME, J. Tribol.* **100** 12–17

[200] Goglia P, Conry T F and Cusano C 1984 The effects of surface irregularities on the elasto hydrodynamic lubrication of sliding line contacts, pt I, single irregularities, pt II, wavy surfaces *Trans. ASME J. Tribol.* **106** 104–112, 113–19

[201] L 1 Turaga R, Sekhor A S and Majumdar B C 1996 Stochanstic FEM model of one dimensional hydrodynamic bearings with rough surfaces *Wear* **197** 221–227

[202] Christensen H and Tonder K 1971 The hydrodynamic lubrication of rough bearings of finite width *ASME J. Lub. Tech.* **93** 324

[203] Lubrecht A A, Ten Nopel W E and Bosnia R 1988 The influence of longitudinal and transverse roughness on the EHL of circular contacts *Trans. ASME J. Tribol.* **110** 421–26

[204] Evens C R and Johnson K L 1987 The influence of surface roughness on elastohydrodynamic traction *Proc. IMechE* **201** 145–50

[205] Bush A W, Skinner P H and Gibson R D 1984 The effect of surface roughness in electrohydrodynamic lubrication *Wear* **96** 177–202

[206] Tander K 1984 A numerical assessment of the effect of striated roughness on gas lubrication *Trans ASME, J. Tribol.* **106** 315–20

[207] White J W and Raad P E 1987 Effect of a rough translating surface on a gas film lubrication a numerical and analytical study *Trans ASME, J. Tribol.* **109** 271–5

[208] Crone R M, Jhon M S, Bhushan B and Karis T E 1911 Modelling the flying characteristics of a rough magnetic head over a rough rigid disc surface *Trans. ASME, J. Tribol* **113** 739–49
[209] Hughes T J R, Frenencz R M and Hallquist 1987 *Comput. Methods Appl. Mech. Eng.* **61** 215–48
[210] White J W, Raad P E, Tabrizi A H, Ketkar S P and Prabhu P P 1986 A numerical study of surface roughness effects on ultra thin gas films *Trans. ASME, J. Tribol.* **108** 171–7
[211] Greenwood J A and Morales–Espejel G E 1994 The behaviour of transverse roughness in EHL contacts. *Proc. Instn Mech. Engrs, Port journal of Engineering Tribology* **208** 121–132.
[212] Morales–Espejel G E, Greenwood J A and Melgar J L 1995 Kinematics of roughness in LHL. In *Proceedings of the 22nd Leeds–Lyon Symposium on Tribology* Lyon, pp 501–523
[213] Venner C H and Lubrecht A A 1992 Transient analysis of surface features in an EHL line contact in the ease of sliding. Trans. ASME *J. Tribology* **116** 186–193.
[214] Lubrecht A A and Venner C H 1992 Aspects of two–sided surface waviness in an EHL line contact. In *Proceedings of the 19th Leeds–Lyon Symposium on Tribology* Leeds, pp 205–224
[215] Venner C H, Couhier F, Lubrecht A A and Greenwood J A 1997 Amplitude reduction of waviness in transient EHL line contacts. In elastohydrodynamic 96, *Proceedings of the 23rd Leeds–Lyon Symposium on Tribology* Eds D Dowson *et al*, pp 103–112 (Amsterdam: Elsevier)
[216] Fohl J 1975 Dissertation Universitat Stuttgart
[217] Uetz H, Khosrawi M A and Fohl J 1984 Mechanism of reaction layer formation in boundary lubrication *Wear* **100** 301–13
[218] Briscoe J B Private communications over a number of years
[219] Rihaczek A M 1969 *Principles of High Resolution Radar* (New York: McGraw–Hill)
[220] Bishop I F and Sniddle R W 1981 Some experimental aspects of running in and scuffing failure of steel discs operating under elastohydrodynamic conditions *Proc. 8th Leeds/Lyon Symp. on Tribology* (Guildford: Butterworths) pp 62–70
[221] Sreenath A V and Raman N 1976 Mechanics of smoothing of cylinder linear surface during running in *Tribol. Int.* **9** 55
[222] Heilman P and Rigney D S 1981 Running–in process affecting friction and wear *Proc. 8th Leeds/Lyon Symp. on Tribology* (Guildford: Butterworths) pp 2–32
[223] Montgomery R 1969 Run–in and glass for motion on grey cast iron surface *Wear* **14** 99
[224] Thomas T R 1978 The characterisation of changes in surface topography during running-in *Proc. 4th Leeds/Lyon Symp. on Tribology (IMechE)* p 99
[225] Whitehouse D J 1978 Surface topography and quality and its relevance on wear *MIT Conf.* p 17
[226] Whitehouse DJ 1976 Approximate methods of assessing surface roughness parameters *Ann. CIRP* **25** 461–5
[227] Schmidt W 1980 Private communication
[228] Westkamper E and Hoffineister H W 1996 Function orientated lapping and polishing of ceramic rolling elements through characterization of the workpiece surface *Annals CIRP* **45** 529
[229] Schneider Yu G 1984 Formation of surfaces with uniform micro patterns on precision and instrument parts *Precision Engineering* **6** 219
[230] Schneider Yu G 1982 Service characteristics of work pieces with regular micro relief *2nd Ed. Machine–building Leningrad*
[231] Tonshoff K K and Kappel H 1998 Surface modification of ceramics by laser machining *Annals CIRP* **47** 471
[232] Man H C 1995 Modification of ceramics surfaces by excimer laser for adhesive bonding *Proc. ICALEO* p 449
[233] Ludema K C 1984 A review of scuffing and running in of lubricated surfaces with asperities and oxides in perspective *Wear* **100** 315–31
[234] Salmon G and de Gee A W J 1981 The running–in of concentrated steel contacts, a system oriented approach *Proc. 8th Leeds/Lyon Symp. on Tribology* (Guildford: Butterworths) pp 15
[235] Archard J F 1959 The temperature of rubbing surfaces *Wear* **1** 438–55
[236] Welsh N C 1957 Frictional heating and its influence on the wear of steel *J. Appl. Phys.* **28** 960–1
[237] Archard J F and Rowntree R A 1981 Metallurgical phase changes in the wear of steel, a critical reassessment *Proc. 8th Leeds/Lyon Symp. on Tribology* (Guildford: Butterworths) pp 285–93
[238] Gans B F 1987 *Proc. Trans. ASME, J. Tribol.* **109** 427–31
[239] Arnell R D, Da Vies P B, Hailing J and Whomes T L 1991 *Tribology: Principles and Design Application* (London: Macmillan)
[240] Suh PM 1973 *Wear* **25** 111
[241] Dawson P H 1962 On the effect of metallic contact on the pitting of lubricated rolling surfaces *J. Mech. Eng. Sci.* **4** 16–21
[242] Onions R and Archard J F 1974 Pitting of gears and discs *Proc. IMechE* **158** 673–82
[243] Tallian T E, McCool J F and Sibley L B 1965/66 Partial elastohydrodynamie lubrication in rolling contact *Proc. IMechE* **180** 169
[244] Leaver R H, Sayles R S and Thomas T R 1974 Mixed lubrication and surface topography of rolling contacts *Proc. IMechE* **188** 461–9
[245] Higginson G R and Leaver R H 1967/70 Fluid lubrication of tapered roller bearings *Proc. IMechE* **184** 18
[246] Lundberg G and Palmgren A 1947 Dynamic capacity of rolling bearings *Acta Polytech. Mech. Eng. Ser. 1* 3; 1952 *Ada Polytech. Mech. Eng: Ser.* **2** 4
[247] Tallian T E and McCool J I 1978 A new model of fatigue life dispersion in rolling contact *Wear* **47** 147–54

[248] Tallian T E, Chiu Y, Huttenlocher D, Sibley L, Kamenshine J and Sindlinger N 1963 Lubricant films in rolling contact of rough surfaces *ASME/ASLE Lubr. Conf. (Rochester, NY, Oct. 1963)* paper SALE 63LC21

[249] Snidle R M and Archard J F 1969/70 Theory of hydrodynamics lubrication for a spinning sphere *Proc. IMechE* **184** 839–48

[250] Kolar D and Dohnal M 1986 Fuzzy description of ball bearing *Wear* **110** 3547

[251] Hirano F, Yamashita N and Kamitani T 1972 Effect of surface roughness on work hardening and rolling fatigue in relation to EHL *Proc. IMechE Conf. on Tribology* paper C8/71, pp 151–63

[252] Gentle C R and Posdari M 1983 Computer simulation of starvation in thrust loaded ball bearings *Wear* **92** 125–34

[253] Martin H M 1916 Lubrication of gear teeth *Engineering* **102** 199

[254] Nayak P R 1972 Surface roughness effects in rolling contact *Trans. ASME, J. Appl. Mech.* **39** 456–60

[255] Carter F W 1925 On the action of a locomotive driving wheel *Proc. R. Soc.* **112** 151

[256] Garnell P 1966/67 Further investigation on the mechanics of roller bearings *Proc. IMechE* **181** 1–10

[257] Comment by Jackobsen according to Sayles 1993 Private communication

[258] Briggs O A D and Briscoe B J 1976 Effect of surface roughness on rolling friction and adhesion between elastic solids *Nature* **260** 313–15

[259] Linfoot E H 1955 *Advances in Geometric Optics* (Oxford Clarendon)

[260] Berthe D, Michan B, Flamand L and Godet M 1977 Effect of roughness ratio and Hertzian pressure on micropits and spalls in concentrated contacts *Proc. 4th Leeds/Lyon Symp. on Tribology (IMechE)* pp 233–8

[261] Gray G and Johnson K L 1972 The dynamic response of elastic bodies in rolling contact to random roughness of their surfaces *J. Sound Vibr.* **22** 323–42

[262] Crandall S J 1963 *Random Vibrations* (Cambridge, MA: MIT Press)

[263] Ramanathan S, Radhakrishnan V M and Vasudevan R 1976 The influence of vibration on the rolling contact fatigue damage of case carburised steels *Wear* **40** 319–24

[264] Hagfors T J 1966 *J. Geogr. Res.* 71 379

[265] Glodez S, Winter H and Stuwe H P 1997 A fracture mechanics model for the wear of gear flanks by pitting *Wear* **208** 77

[266] Sayles R S and Macpherson P B 1982 Influence of wear debris on rolling contact fatigue *ASTM Tech. Publ.* **771** 255–74

[267] Webster N M, Lannide E and Sayles R S 1985 The effect of topographical defects on the contact stress and fatigue life in rolling element bearings *Proc. 12th Leeds/Lyon Symp. on Tribology* (London: MEP)

[268] Kimura Y and Sugimura J 1984 Microgeometry of sliding surfaces and wear particles in lubricated contact *Wear* **100** 3345

[269] Carson R M and Johnson K L 1971 Surface corrugations spontaneously generated in a rolling contact disc machine *Wear* **17** 59

[270] Prakash J 1984 On the lubrication of rough rollers *Trans. ASME J. Tribol.* **106** 211–17

[271] Tallian T E 1978 Elastohydrodynamic effects in rolling contact fatigue *Proc. 5th Leeds/Lyon Symp. on Tribology* (London: MEP) pp 253–81

[272] Rosenberg R C and Trachman E G 1974 Ball grease bearings *Mech. Eng.* Dec., pp 27–9

[273] Hoepnner D W and Goss G L 1972 Corrosion fatigue *NACE Publ.* p617

[274] Waterhouse R B and Allery M 1966 *Trans. Am. Soc. Lubr. Eng.* **9** 179

[275] Fellows L J, Nowell D and Hill D A 1997 On the initiation of fretting fatigue cracks *Wear* **205** 120–129

[276] *American Society of Metals Handbook* p 671, fig. 7, reference [4]

[277] Greenfield P and Allen D H 1987 The effect of surface finish on the high cycle fatigue strength of materials *GEC J. Res.* **5** 129–0

[278] *American Society for Metals Handbook* 9th edn, **11** 126

[279] Taylor D and Clancy O M 1991 The fatigue performance of machined surfaces *Fatigue Fract. Eng. Mater. Struct.* **14** 329–36

[280] Taylor D 1989 *Fatigue Thresholds* (London: Butterworths)

[281] Kitagawa H and Takahashi S 1976 Applicability of fracture mechanics to very small cracks or the cracks in the early stage *Proc. 2nd Int. Conf. on Med. Behav. Mater., (Boston)* pp 627–39

[282] Taylor D 1990 Crack like notches and notch like cracks *Proc. ECF–8, (Turin, Italy)* (EMAS, UK)

[283] Siebel E and Gaier M Z 1957 *Ver. Dtsch mg.* 98 *Eng. Dig.* (transl.) **18** 109

[284] Finnie 1995 Some reflections on the past and future of erosion *Wear* **186–187** 1–10

[285] Magnee A 1995 Generalized law of erosion application to various alloys *Wear* **181** 500–510

[286] Achter M R 1966 *ASTM Conf. on Crack Fatigue, (Atlantic City, NJ)*

[287] Evans U R 1960 *Corrosion and Oxidisation of Metals* (London: Edward Arnold)

[288] Tomashov N D 1966 *Theory of Corrosion and Protection of Metals* (New York: Macmillan)

[289] Ehrlich C and Tumbull D 1959 *Physical metallurgy of stress corrosion fracture* (New York: Interscience) p 47

[290] Ogilvy J A 1991 *The Theory of Wave Scattering from Random Rough Surfaces* (Bristol: Hilger)

[291] Beckmann P and Spizzichino A 1968 *The Scattering of Electromagnetic Radiation from Surfaces* (Oxford: Pergamon)

[292] Bass F O and Fuchs I M 1963 *Wave Scattering of Electromagnetic Waves from Rough Surfaces* (New York: Macmillan)

[293] Teague C, Vorberger T V and Maystre D 1981 Light scattering from manufactured surfaces *Ann. CIRP* **30** 563

[294] Maystre D 1988 Scattering by random rough surfaces in electromagnetic theory *Proc. SP1E* **1029** 123

[295] Maystre D and Saillard M 1988 Rigorous theory of metallic dielectric non–periodic rough surfaces and its numerical implementation *Proc. SPIE* **29** 94

[296] Crock P and Prod'hamme L 1980 Contribution de la technique d'immersion a l'analyse de la diffusion optique par des surfaces tres rugeuses *J. Opt. (Paris)* **II** 319–22

[297] Whitehouse D J 1970 *PhD Thesis* University of Leicester

[298] Thomas T R (ed) 1982 *Rough Surfaces* (London: Longman)

[299] Rakels J H 1998 Reconstruction of surface when only differentials are known *J. of Inst. Physics Nanotechnology* **9** 49–53

[300] Ament W S 1960 Reciprocity and scattering by certain rough surfaces *Trans. Ant. Prop.* AP–8 167–74

[301] Welton P J, Frey H G and Moore P 1972 Experimental measurements of the scattering of acoustic waves by rough surfaces *J. Acoust. Soc. Am.* **52** 1553–63

[302] Lynch P J and Wagner R J 1970 Rough surface scattering shadowing and multiple scatter, and energy considerations *J. Math. Phys.* **11** 3032–42

[303] Osanna P H 1982 Deviations in form and workpiece accuracy *Wear* **83** 255

[304] Wagner R J 1967 Shadowing of randomly rough surfaces *J. Acoust. Soc. Am.* **41** 138–47

[305] Hardin J C 1971 Theoretical analyses of rough surfaces shadowing from point source radiation *J. Acoust. Soc. Am.* **52** 227–33

[306] Thorsos E I 1988 The validity of the Kirchhoff approximation for rough surface scattering using a Gaussian roughness spectra *J. Acoust. Soc. Am.* **83** 78–92

[307] Crandall S H 1969 First crossing probabilities of the linear oscillator *Applications and Methods of Random Data Analysis (Inst. of Vibration Research, Southampton, July, 1969)* paper V

[308] Whitehouse D J, Jungles J and Goodwin E 1972 Modern methods of assessing surfaces *Proc. Jpn Inst. Prod. Eng.* p 503–24

[309] Dainty J C (ed) 1975 *Laser Speckle and Related Phenomena* (Berlin: Springer)

[310] Jakeman E and Jefferson J H 1986 Scattering of waves by refractive layers with power law spectra *Proc. Inst. Acoust.* **8** 12–19

[311] Church E L, Jenkinson H A and Zavada J M 1977 Measurement of the finish diamond turned metal by differential light scattering *Opt. Eng.* **16** 360

[312] Franklin S P and Schneider C R A 1988 Rough crack like defects in NDT morphology and ultrasonic cracking *CEG 13 Rep.* OED/STN/88/20003/R

[313] Underwood E E and Banerji K 1986 Fractals in fractography *Mater. Sci. Eng.* **80** 1–14

[314] Pande C S, Richards L E, Louat N, Dempsey B D and Schwoeble A J 1987 Fractal characterisation of fractured surfaces *Acta Metall.* **35** 633–7

[315] Berry M V and Lewis Z V 1980 On the Weierstrass–Mandelbrot fractal function *Proc. R. Soc.* A **370** 459–84

[316] Mitchell M W and Bonnell D A 1990 Quantitative topographic analysis of fractal surfaces by scanning tunnelling microscopy *J. Mater. Res.* **5** 22–54

[317] Sayles PhD p 920

[318] Berry M V 1979 Diffractals *J. Phys. A: Math. Gen.* **12** 781–97

[319] Mandelbrot B B 1977 *The fractal geometry of nature* (New York: Freeman)

[320] Wormington M and Bowen D K 1992 Computer simulation of X–ray reflectivity data *American Crystallography Association, (Denver, Aug 1992)*

[321] Ogilvy J A 1985 Ultrasonic wave scattering from rough defects the elastic Kirchhoff approximation *UKAEA AERE Ri 1866* (HMSO)

[322] Wormington M, Bowen D K and Tanner B K 1992 Principles and performance of a PC based program for simulation of grazing incidence X–ray reflectivity profiles MRS Symp. Proc. **236** 119–24

[323] Parratt L G 1954 *Phys. Rev.* **95** 359

[324] Ber A and Yarnitzky Y 1963 Functional relationship between tolerances and surface finish *Microtechnic* **XXII** 449–51

[325] Osanna P H 1962 Surface roughness and size tolerance *Qualitat* (Vienna: Technical University of Vienna)

[326] Bjorke O 1978 *Computer–aided Tolerancing* (Oslo: Tapir)

[327] Osanna P H and Totewa P 1992 Workpiece accuracy: the critical path to economic production *Int. J. Mack. Tool Manuf.* **32** 45

[328] ONORM M .1116 1988 Surface roughness and ISO tolerance quality

[329] ISO 286–1 1988 ISO system of limits and fits. Part 1 Bases of tolerances deviations and fits; Part 2 Tables of standard tolerances, grades and limit deviations for holes and shafts

[330] Martin A M 1916 The lubrication of gear teeth *Engineering* **102** 119–121

[331] Berthe D and Vergne P H 1987 An elastic approach to rough contact with asperity interactions *Wear* **117** 211 [35] Tabor D 1975 A simplified account of surface topography and the contact between solids *Wear* **32** 268–71

[332] Hisakado T 1976 The influence of surface roughness and work hardening layers on the contact between a rough and flat surface *Wear* **37** 41–51

[334] Hisakado T 1975 Effect of surface roughness on contact between a rough and flat surface *Wear* **35** 53–61 [38] Hisakado T 1974 Effect of surface roughness on contact between solid surfaces *Wear* **28** 217

[335] Schofield R E and Thomley K H 1973 Calculating the elastic and plastic components of deflection of joints fanned from mechanical surface with flatness errors *Proc. 13th MTDR Conf.* (London: Macmillan) pp 67–74

[336] Hailing J and Nuri S 1985 The elastic contact of rough surfaces and its importance in the reduction of wear Proc. *IMechE* **199** 139

[337] Semenyuk N F 1986 Average height of roughness asperities and density of contact spots in the contacting of a rough surface with a smooth surface *Trenie 1, Iznos* **7** 85–90 *(Sov. J. Friction Wear)*

[338] Semenyuk N F 1986 Average values of total and mean summit curvatures and heights of asperities of an isotropie rough surface *Trenie I, Iznos* **7** 830–40

[339] Thomas T R 1973 Influence of roughness on the deformation of metallic surfaces in static contact *6th Int. Conf. on Fluid Sealing* paper B3–33

[340] Back N, Burdekin M and Cowley A 1973 Review of research on fixed and sliding joints *Proc. 13th MTDR Conf.* (London: Macmillan) pp 87–98

[341] Schofield R E and Thomley R H 1967/68 Mathematical expressions of surface finish characteristics *Proc. IMechE* **182** 446

[342] Babus' Haq R F, George H E, O'Callaghan P W and Probert S D 1990 Thermal template for the prediction of the thermal resistance of a pressed contact *Advanced Computational Methods in Heat Transfer, Heat Conduction, Convection and Radiation, Proc. 1st Conf, (Portsmouth)*

[343] Marks W D 1972 Characterisation of statistical transients and transmission media. The measurements of power moments spectra *J. Sound Vib.* **22** 149–296

[344] Whitehouse D J and Zheng K G 1992 The use of dual space frequency functions in machine tool monitoring *Meas. Sci. Technol.* **3** 9

[345] Cans B F 1987 *Proc. Trans. ASME, J. Tribol.* **109** 427–31

[346] Nilsson L R K 1978 The influence of bearing flexibility on the dynamic performance of radial oil film bearings *Proc. 5lh Leeds/Lyon Symp. on Tribology* (London: MEP) pp 311–19

[347] Sharratt F and Whitehouse D J 1992 Surface texture and fatigue *Rep.* no 8, Microengineering Centre, University of Warwick

[348] Boit M A 1957 On the reflection of acoustic waves on a rough surface *J. Accoust. Soc. Am.* **30** 479–80

[349] Grasso, F, Muumeci F, Scordmo A and Ronchi L 1983 Second and higher order scattering of electromagnetic waves by a rough metal surface *Opt. Acta* **11** 1–22

[350] Rakels J H 1986 The use of Bessel functions to extend the range of optical diffraction techniques for in–process surface finish measurement of high precision turned parts *J. Phys. E: Sci. Instrum.* **19** 76

[351] Rakels J H and Hingle H T 1986 The use of optical diffraction techniques to obtain information about surface finish tool shape and machine tool condition *Wear* **109** 259

[352] Mulvaney D J, Newland D E and Gill K F 1986 A characterization of surface texture profiles *Proc. Inst. Mech. Congrs* **200** 167–78

[353] Kobayashi T 1991 Reconstruction of crack history from conjugate fracture surfaces *Fatigue Crack Measurements. Techniques and Applications* ed K J Marsh, R A Smith and R O Richie (UK: Eng. Materials Advising Services) p389

[354] Kobayashi T and Shockey D A 1989 Computational reconstruction of environmentally accelerated cyclic crack growth in reactor steels *Proc. NACE Corrosion 89 Sympos.* (Houston: NACE) paper 563

[355] Love A E H 1944 *A Treatise on the Mathematical Theory of Elasticity* (1892) (Cambridge: Cambridge University Press)

[356] Prescott J 1961 *Applied Elasticity* (1924) (London: Longmans Green)

[357] Puttock M J and Thwaite E G 1969 *National Standards Lab. Tech. Paper No* 25, (Melbourne: CSIRO)

[358] Kellog O D 1929 *Foundations of Potential Theory* (New York: Murray)

[359] Williamson J B P 1967/68 Microtopography of solid surfaces *Proc. IMechE* **180** 21–30

[360] Archard J F 1957 Elastic deformation and the laws of friction *Proc. R. Soc. A* **243** 190–205

[361] Dyson J and Hirst W 1954 The true contact area between solids *Proc. Phys. Soc. B* **67** 309–12

[362] Greenwood J A 1967 The area of contact between rough surfaces and flats *Trans. ASME, J. Lubr. Technol.* **89F** 81–9

[363] Greenwood J A and Trip J H The elastic contact of rough surfaces *Trans. ASME, J. Appl. Mech.* **34E** 153–9

[364] Greenwood J A *Tribology Rough Surfaces* ed T R Thomas (London: Longman) p 191

[365] Greenwood J A Johnson K L and Matsubara E 1984 S surface roughness parameter in Hertz contact *Wear* 47–57

[366] Mikic B B and Roca R T 1974 *Int. J. Heat Mass Transfer* **17** 205

[367] Rossmanith H P 1976 Stochastic analysis of non top–contacting rough surfaces *Wear* **37** 201–8

[368] Whitehouse D J 2001 Some theoretical aspects of structure functions, fractal parameters and related subjects *Proc. Inst. Mech. Eng.* **215** 207–10

# Chapter 8
# Nanometrology

## 8.1  Introduction

This chapter is an exposition of the meaning of nanometrology and how it fits in with surface metrology. There is a description of concepts rather than the reporting of lists of specifications and data. The situation is that the subject is changing so rapidly that factual information quickly becomes out of date; the underpinning concepts have so far had too little time to be formulated convincingly. A primary objective therefore is to consolidate these underlying concepts. However, factual information has not been neglected. It has been distributed appropriately throughout the book rather than being lumped, out of context, in this chapter.

Nanotechnology started as a name in 1974. It was introduced by Taniguchi [1] to describe manufacture to finishes and tolerances in the nanometre region. He extrapolated the specifications from existing and past machine tools, such as lathes and grinders, to the new generation of machine tools. He concluded quite correctly that in the late 80s and 90s accuracies of between 0.1 $\mu$m and 1nm would be needed to cater for industries' needs (figure 8.1). It soon emerged that the only way to achieve such results was to incorporate very sophisticated instrumentation and ultra precision dimensional metrology into the design [2].

### 8.1.1  Scope of nanotechnology

This move towards instrumentation and metrology was already developing in the electronics industry where the drive was towards miniaturization for higher packing densities and faster switching. As a result, highly controllable and stable processes such as lithography were introduced. This meant a need arose for very accurate positioning of the specimen. In turn this resulted in an interest in miniature actuators, motors and accurate slideways for which new technologies have had to be developed. In particular, new materials and thin film research were pre-eminent.

As well as electronics and manufacture, new developments on the nanoscale are taking place in biology and chemistry [3]. In terms of disciplines therefore nanotechnology encompasses more than engineering. It is the meeting point at the atomic scale incorporating biology, physics, engineering and chemistry. This overall position is shown in figure 8.2. Not only has the science base increased but also the disciplines within each science.

### 8.1.2  Nanotechnology and engineering

Because nanotechnology, in name at least, started in engineering, it is probably most informative to follow and investigate the growth of the subject in this discipline.

Developing from the need to machine more accurately as demands grew, for example in the fields of compact discs, gyroscopes, etc, came new methods of fabrication with different materials. Together with these applications came the need to make smaller sensors and transducers to enable non-intrusive control. In engineering applications, 90% of transducers are concerned with the measurement of displacement,

position or their derivatives such as strain, pressure and acceleration. This has resulted in a mechanical microworld which has emerged from the technology developed for integrated circuits. Already many small mechanisms are being made, some including miniature motors of mm dimensions. Highly reliable accelerometers are already in use.



**Figure 8.1** Development of achievable machining accuracy.

These devices are fabricated on silicon substrates using extensions of such integrated circuit manufacturing processes as photolithography, thin-film deposition and the like. Microstructures are more ambitious but are now being developed in many parts of the world. Using innovative processes involving bulk micro-machining—sculpting silicon with chemical etchants and surface micromachining—etching layers of thin films deposited on the substrates or bombarding with ion beams, has changed the concept of machining well beyond conventional methods.



**Figure 8.2** Nanotechnology: the meeting point.

Also, the new use of parts so machined has created new names in the technical literature. Micro-electromechanical systems (MEMS), microdynamics, micromechanics, etc, are words which have emerged. They are perhaps best put into perspective by means of typology mapping as shown in figure 8.3.

Typology of disciplines based on dimensional characteristics



Key

o—o—o Nanotechnology
x—x—x Micromechanic
– – – Microdynamics
•—•—• Precision engineering
△—△—△ Conventional engineering

**Figure 8.3** Typology of disciplines based on dimensional characteristics.

In this figure conventional engineering includes turning etc; precision engineering comprises diamond turning and grinding. In micromechanics lapping and polishing are employed, and in nanotechnology the new methods mentioned above. The important point is that nanotechnology also utilizes advances in the conventional processes such as ductile (not brittle) grinding of materials, such as silicon, glass and ceramics. Because of the ever-increasing tendency to form the essential shape of a component from composite material, bulk material removal methods are becoming used less often. Finishing processes, old and modern, are coming into their own. Hence it can be seen that nanotechnology, as seen in figure 8.3, crosses the boundaries of the other disciplines. For this reason it can be regarded as an 'enabling technology'.

One consequence of the growth of nanotechnology is that the interests of engineers and scientists are beginning to overlap. The engineer is striving for more and more refined manufacture and measurement whilst the physicists and chemists are trying to synthesize. The engineers are attempting to move top-down and the physicists bottom-up. Nanotechnology is where they meet. This causes some problems in communication.

Nanotechnology is not a science of the future—it exists today in many ordinary domestic appliances such as tap washers and CD players, where the rotation has to be true to within a few nanometres. In fact in certain fields such as ruling and grating machines it has existed for almost a century.

### 8.1.3 Nanometrology

Nanometrology is usually reserved for measuring the geometrical features of size, shape and roughness. Also, although not a direct feature of the workpiece, its position relative to a co-ordinate system is often included. These are features usually required for engineering applications. However, the with the new generation of microscopes,(e.g. SPM) at least one of the features is not geometrical but physical (e.g. force). Aspects of nanometrology are measured selectively with respect to discipline. Figure 8.4 shows a simple breakdown in which each axis is in length. The y axis is here taken as unity. The point to note is that the engineer,

physicist, biologist and chemist have different relative proportions of the geometry to deal with. For example, in engineering applications the roughness is much smaller than the size or shape.

As mentioned above the nature of many disciplines begins to change character as the scale of size gets smaller. In particular as the scale approaches atomic and molecular dimensions the extra physical constraint of quantum mechanics has to be take into account. One embodiment of this is the tunnelling of electrons through a potential barrier.



**Figure 8.4** Relative relationship between geometric components (*a*) conventional engineering (*b*) micromechanics, (*c*) quantum engineering.

## 8.2 Effect of scale on manufacture, functionality and instrument design

### 8.2.1 *Machineability*

#### *Nanomachining*

Traditional precision finishing of hard and brittle materials such as single crystal wafer and magnetic heads consists of grinding, lapping and polishing. These depend too much on skilled labour. Current thinking suggests that these traditional processes could be replaced by a 'nanogrinding' producing damage–free surfaces.

The conventional production of form is based on pressure copying where the main parameters are (i) grain size, (ii) uniformity of grains and (iii) machining pressure. Material removal $M$ is based on Preston's formula $M = pvk$, where $p$ pressure, $v$ is velocity and $k$ is a constant. This refers to loose abrasive grains such as are found in lapping. There are three basic modes of material removed in this regime: (i) brittle mode machining, (ii) microcrack machining (iii) ductile mode machining. It is interesting to note that the mechanisms of microgrinding are paralleled by wheel wear mechanisms.



**Figure 8.5** Material removal, wheel wear modes.

What the figure describes is the amount of material removed. What it does not give is the amount simply moved. The regime where plastic flow occurs around the grain is where ductile machining fits.

The critical question is how can ductile machining, whether grinding or turning, be guaranteed? The answer to this in part is that it has already been achieved in diamond turning. Optical parts have been made using diamond turning; the problem now is to carry the techniques to the grinding process.

The removal of brittle materials has conventionally been considered to be caused by indentation and scratching by an abrasive powder as in lapping and polishing. Hence studies of material removal in the case of brittle parts have focused on fracture mechanics mechanisms. Indentation and scratch tests have been extensively used to help in this study. In the case of lapping and polishing, the pressure distribution between the workpiece and tool dominates. This had led to the concept of 'pressure copying' which has been the traditional method of surface generation to a high finish. The value of pressure on each grain (or average grain; $p_c$) has been found.

There is another copying method. This is called the 'motion-copying' technique in which the tool is forced to follow a given path determined elsewhere in the machine tool. The figure on the workpiece becomes an exact replica of this motion. Diamond-turning methods on soft materials using grinding are not simple. It depends on scale of size. This has meant the determination of a critical distance $d_c$ between the workpiece and tool which has to be maintained in order to guarantee ductile conditions. In an important paper Bifano *et al* [4] found experimentally the value of $d_c$ for germanium. Subsequently, the $d_c$ values were found for a variety of materials. It was found that $d_c$ could used to tighten the specification of grinding (or cutting) machines in terms of accuracy, feed, resolution and stiffness. From this requirement a new generation of grinding machines has been developed mainly based on the work of Miyashita and co-workers [5].

If the extraneous displacements, flexures, strains, or whatever, can be kept below $d_c$ for a given material, then ductile cutting can be achieved; if it is not, then the brittle regime takes over and the roughness and also the 'figure' to a lesser extent deteriorate. The figure shows how the nanogrinding domain relates to conventional grinding in terms of grain size. Notice the difference on the height distribution of active grains.

Some recent important work has been reported by Miyashita. He realized that the surface generated by the various material removal modes shown in figure 8.5 should be related to parameters of the process itself. Usually grinding is selected because it can be controlled confidently and is also most often used. Miyashita is a realist when he states that mechanisms of wheel wear and surface production are not yet based on 'good physics'.

A principal task is to try to predict surface roughness from the characteristics of the grinding process and include such factors as 'run-out' of the wheel. He has arrived at the formula given below.

If $dg$ is depth of cut of grains on the wheel

$$dg = 2a\frac{v_w}{v_s}\sqrt{\frac{d_w + \Delta d_w}{D_e}} + \Delta h \tag{8.1}$$

where $a$ = distance between sequential cutting edge
$v_w$ = peripheral speed of workpiece
$v_s$ = peripheral speed of grinding wheel
$d_w$ = wheel depth of cut
$\Delta d_w$ = wheel run-out or waviness of wheel periphery
$D_e$ = equivalent wheel diameter
$\Delta h$ = range of height distribution of grain tops
$\lambda$ = critical wavelength between waviness and random components = $2\pi a$



**Figure 8.6** Process parameters in terms of spectral density.

The question is whether the $\Delta h$ is related to any of the common surface parameters such as $R_a$. In another form: are the surface statistics of the wheel related to the surface statistics of the workpiece?

The criterion of ductile mode grinding,

$$dg \sim \Delta h \leq d_c \tag{8.2}$$

where $d_c$ is the ductile to brittle transition in the depth of cut, summarizes nanomicrogrinding in terms of material removal (a), surface finish and size tolerance (b) and related processes (c).

It can be seen that progress is being made at understanding ductile machining and in particular how the surface finish can predicted given the process parameters. One question which has to be asked is how $\Delta h$, the range of cutting edge heights, is measured. If a stylus method is used, the range of heights measured $\Delta h$ is always smaller than the true value, because the stylus rides on the grain flanks rather than over the summits. It could be argued that the same argument applies when measuring the finish on the surface. This does not follow because the grain impression in the surface is smoother than predicted from the geometry of the grains because of the ductile flow produced!

$Z'$ [mm$^3$/mm s], $v_s$:const. (30m/s)

| $10^2$ | $10^1$ | $10^0$ | $10^{-1}$ | $10^{-2}$ | $10^{-3}$ | $10^{-4}$ | $10^{-5}$ |

Grinding

Microgrinding

Size tolerance [μm]

Surface roughness [nm]

$R$a

10

10

$R$a

Size

1

0.1

0.1

0.01

0.01

0.001

Work materials: hardened bearing steel

Material removal rate—MRR [mm$^3$/mm sec]

| $10^2$ | $10^1$ | $10^0$ | $10^{-1}$ | $10^{-2}$ | $10^{-3}$ | $10^{-4}$ | $10^{-5}$ |

Performance of grinding process

Frontier ⟨⟨⟨ ⟩⟩⟩ Frontier

High MRR grinding       Convertional       Fine MRR grinding

Lapping, honing
smoothing, super finishing

Polishing

**Figure 8.7** Nano/microgrinding and its relationship to similar processes.

The real problem with 'ductile' machining is that it is a function of 'scale of size' and is not fully understood. The very small dimensions involved (~$10^2$nm) suggest that it may be linked to the density of dislocations in the material. At the instant of cutting, the energy between kinetic energy of the material and the surface free energy required to produce discrete chips becomes unbalanced. This means that the surface free energy is too high for individual chips to be produced, with the result that a continuous chip is produced indefinitely. Molecular dynamic simulations might give the answer.

### 8.2.2    Dynamic considerations—how force balance depends on scale

Figure 8.8 shows how the balance of forces in structures changes with scale of size.

As the scale of size is reduced there is a profound change in the relative importance of the components of the force equations. At the mm level and above the inertial component dominates. However, as the scale is reduced to the nanometre level, the inertial term falls off faster in magnitude than either the damping or elastic terms. The most important component becomes the elastic force, shown in figure 8.8*b*. However, elastic

**Figure 8.8** How the force balance changes with scale of size.

properties are controllable because material properties are known. The inertial term is negligible leaving the damping term which is areal and is consequently dependent on the surface finish. Adhesion and friction depend a great deal on the relatively uncontrolled surface. The change in the relative importance of the forces is the reason why insects have such different movement capabilities than mammals. Insects can change direction instantly indicating little inertia. They can walk upside down on ceilings indicating that adhesion forces are much greater than forces on mass, (e.g. gravity).

It is essential that this miniature force balance is understood if MEMS and related objects are to function properly. Furthermore, the intricate measurements in 3D required to characterize MEMS is far more complicated than for planar measurement based on semiconductor technology.

### 8.2.3 Structural considerations—how stiffness and resonance change with scale of size

The scale of size is important in two aspects of structures. The dynamics of the system changes balance with scale as indicated above, which has important implications in microdynamic structures such as pumps and valves. However, there is another scale effect on structures which relates to static behaviour. This is important in the manufacture of nano-structures as well as their measurement (i.e. the nanometrology). The force loop in machining and the metrology loop in instruments depend on scale (figure 8.9).

Taking the loop in figure 8.9 in its simplest form as a cantilever fixed at one end at earth and subject to forces at the other, in this case the cutting force, the deflection $\delta$ of the tool due to a force $P$ is

$$\frac{Pl^3}{3EI} \qquad (8.3)$$

at the point of application.

A similar formula for deflection in the middle of a loop uniformly loaded $W$ per unit length is

$$\frac{Wl^4}{EI} \cdot \frac{5}{24} \qquad (8.4)$$

$E$ is Young's modulus and $I$ is the area of inertia.

**Figure 8.9** Mechanical loop of simplified machine tool.

The stiffness of the loop is $\dfrac{P}{\delta} \sim \dfrac{K}{l^3}$ (8.5)

From equation (8.3) if the scale is reduced by a factor of two, the stiffness increases eight times for the same load $P$ at (i.e. $P$ at $\dfrac{l}{2}$ producing $\delta'$ (figure 8.10)).



**Figure 8.10** Cantilever deflection.

So if accurate dimensions and shapes are required to nanometre accuracy the machine tool should be small. Also deformation due to thermal gradients across the loop is less likely. This means that any metrology loop maintains its shape and so preserves surface metrology fidelity but not dimensional fidelity.

The size dependence of deformation of a cantilever is not only dependent on the scale as shown above, but the elastic properties themselves have been shown [6] to depend on scale. Furthermore the effect is surface dependent which is not obvious at mm and $\mu$m scales. The important length scale is identified to be determined by the ratio of the surface modulus $S$ to that of the bulk $E$, $\left(\text{i.e. } \left(\dfrac{S}{E}\right)\right)$. In addition a further factor is the geometrical shape of the arm represented by $\alpha$ and the length represented by $l$. Taken as whole there is a non-dimensional parameter which influences the effective stiffness by a significant amount.

$$\text{i.e. } \left(\dfrac{S}{E}\right) \cdot \left(\dfrac{\alpha}{l}\right)$$ (8.6)

Such a relationship can be predicted by continuum mechanics but breaks down at quantum levels. It is interesting to note that the expression involves both the surface and bulk moduli. Very often these are assumed to be the same. It is at the nanometre level where the differences become apparent.

### 8.2.4 Functional dependence on scale of size

Figure 8.11 shows where the behaviour patterns of static and dynamic forces fit the function map concept explained in chapter 7 and, in particular, figure 7.165. In many practical functions these two components are complementary. In terms of magnitude it is the dynamic component which dominates behaviour. However, consider now figure 8.12. This is the same as figure 7.166 except that the line representing the components of a second order force equation can now be reassigned to represent the change in emphasis of the force within a system as the scale of size reduces. Volume domination (inertia terms) give way to areal and finally to linear dominance. Miniaturization to the nanoscale pushes the emphasis towards the origin of the function map. This constitutes a profound change in the importance of the component forces and their resultant movements and deflections. In these regions nanometrology 'enables' semiconductor and often planar applications by means of position, size and roughness control. Size, shape and roughness albeit at somewhat larger dimensions are effective in 3D applications such as MEMS. In the former applications the typical scale of size is nanometres or less, whereas in the latter it can extend from tens of nanometres up to hundreds of micrometres.



**Figure 8.11** Force boundary domain split up into static and dynamic components.



**Figure 8.12** Force regime map showing force balance trend.

### 8.3 Metrology at the nanoscale

#### 8.3.1 *Nanometre implications of geometric size*

In precision engineering dimension, shape and roughness are limited by the manufacture and the process mechanisms. For example, in general manufacturing, the time of machining, depth of cut and tool path are intended to produce the desired shape, size and roughness. Lack of control of any of these process parameters results in scrap. The situation is different at molecular and microcluster scale. Roughness as such is hardly definable. It is shape and size which dominate as will be seen.

A regime of scale is given in Table 8.1.

**Table 8.1** Groupings of atoms.

|  | Bulk | Particles | Microcluster | Molecular |
|---|---|---|---|---|
| Size Dimension | $\mu$m - $\mu$m | 100nm | 10nm | 1nm |
| Size Number | > 106 | 104 to 106 | 102 to 104 | 1 to 102 |

Table 8.1 refers more to physics and chemistry than it does to engineering. Breaking the categories down gives the following.

(a) 'Bulk' in the table includes MEMS (micro electro mechanical systems), micrometers as well as large structures.
(b) Particles include powders and composites, for example aluminum.
(c) Microclusters include colloids, catalysts, nanotubes, fullerines.
(d) Molecular scale includes thin films such as Langmuir-Blodgett, also self-assembly proteins.

The question is whether the traditional geometric breakdown into roughness, shape and size still holds in the regimes outlined above. This issue will be considered next.

#### 8.3.2 *Geometric features and the scale of size*

Let the total geometry $G$ be represented by a scalar sum of the components taken as estimates only.

$$G = S_i \cap S_{sh} \cap R \tag{8.7}$$

In macroengineering $S_i$ (size) is much greater than $S_{sh}$ (shape) on $R$ (roughness). Also as shown in equation (8.7) they are nominally independent of each other. i.e. the values of the shape and roughness do not substantially affect the value of the size or the size affect the roughness.

Equation (8.7) can be represented by a graph (figure 8.13).

Although only shown simply, figure 8.13 demonstrates the change in relative positions of the geometrical components as the scale of size is reduced. The bunching together in the practical case of figure 8.13(*c*) makes it difficult to assess size shape and roughness independently. The reason for this is largely due to the difficulty in scaling down roughness.

**Table 8.2** Variability of components.

| Feature | Size | Shape | Roughness |
|---|---|---|---|
| No. of principal variables | One variable | Two variables | Multivariables |

Table 8.2 is a guide illustrating the reason for the bunching; whereas reduction in size is usually straightforward in the sense that there is usually one major variable, in roughness there are many. For example, in

**Figure 8.13** Relative components of geometry vs scale (*a*) macro engineering (*b*) ideal scaled down macro engineering (*c*) actual micro engineering.

grinding there is not a great range of grain sizes available so reduction of the scale of size does not necessarily have a suitable size of grain. Similarly, shape is determined by the path of the tool and has more than one input not rigorously balanced.

Under these conditions the geometric factors cannot be considered to be independent. For example, the size of an element is usually taken between surface mean lines as shown in figure 8.14.



**Figure 8.14** Definition of size (*a*) with small surface roughness (*b*) with large relative roughness.

### 8.3.3 How roughness depends on scale

It is clear from figure 8.14 that determination of the mean lines is difficult because of the large roughness and small sample. The spacing between the mean lines has a large error tolerance. If instead of the mean lines, the peak envelopes are used the error can be bigger. At the same time the estimation of the roughness is prone to error again because of the limited sample length $W$. Some examples of micromechanical devices showing the disproportionate size of the roughness are given in figures 8.15 and 8.16.

Roughness has to be considered at the nanoscale in three cases,

(1) on large objects such as mirrors
(2) on small objects such as in microdynamics
(3) at the molecular and atomic scale.

**Figure 8.15** (*a*) Micro gear, (*b*) micro stud.



**Figure 8.16** LIGA machined micro gear.

In many cases in nanometrology the role of roughness and waviness is conventional. It is quite possible to have roughnesses on silicon wafers and mirrors that are well below nanometre level and can even reach angstrom values. As the object is large compared with the roughness, traditional machining methods can be used such as grinding. The roughness is generated by the production process and consequently can be used to control the process in the normal way. The same is true for waviness only with respect to the machine tool as shown in chapter 6.

There is a growing number of applications, however, where the roughness is on small parts, sometimes less than a fraction of a millimetre.

It is useful at this stage to review how the various components of the geometry of a workpiece behave as the scale is reduced. This is analogous to the balance of forces discussed in section 8.2.2.

Roughness of large particles and associated powders has not been the subject of much investigation. Roughness is regarded as irrelevant. Probably the only obvious effect of roughness is to increase the surface area and hence the surface energy. So there is a limit to the extent that small angular particles can be fragmented.

At the cluster stage it becomes questionable whether roughness has any meaning. The cluster by definition has a smaller molecule/atom density in the bulk rather than on or near the surfaces so the surface is very important in imparting properties. However, this is one situation where roughness and shape take on different roles. The shape is determined by the configuration of molecules. It is automatically generated. Roughness or what could be called roughness has to be ascribed to problems in shape generation (i.e. the presence of defects).

Roughness at the molecular scale usually means the roughness at the substrate interface and effects resulting from it which relatively speaking is roughness on a large object. In one instance the interface roughness causes stiction between the polysilicon layer which can inhibit the release of one from the other. On the

other hand [7] a certain amount of substrate roughness can be used via soft x-ray scatter to measure the film thickness and the position of assembling proteins in biological switches. Other than these rare applications roughness changes its definition at the nanotechnology scale to mean the presence of defects. It could also mean the disorder in the layers of whiskers of carbon nanotubes and fibres [8] when generated. In these situations 'roughness' can be avoided by careful and controlled processing. This is exactly the opposite of what happens in macro-engineering where the roughness is inevitable, being a relic of the process.

It seems from above that roughness as such is difficult to identify in subnanometre applications except for substrate roughness and optical scatter.

Although the patterning of surfaces by etching or lithographic methods could be construed as roughness, in this chapter it is not. The laser patterning of rollers for sheet steel reduction is a recognized roughening technique so here again the definition of roughness changes at the nanoscale; transverse patterns with no height variation may well be roughness in the future. There is some confusion because the traditional importance of roughness has been in contact and flow situations and some areal patterns on the surface influence the tribology, however similar types of patterns can be used in lithography for determining structure.

In the same way that conventional machining is subtractive (i.e. in stock removal) the bulk micromachining by etching of silicon wafers is also subtractive. However, when micromachining the surface the process is often additive, e.g. by sputtering.

There are other factors about the surface which are important and yet do not fall into the category of surface roughness. One such feature is the orientation of the electronic structure at the surface. The growth of molecularly thin layers of oxides, sulphides and the like depends on symmetry of electron structure at the surface. This is not a process as such but it is a surface effect of reactivities which can influence the shape of local self assembled islands [9].

Semiconductors often have a symmetry at the surface which is different to the bulk. This produces a new 'unit cell' at the surface which is called 'surface reconstruction' and is a consequence of the structure seeking a low potential energy. The reconstruction influences deposition of material on the surface [10].

The selectivity and efficiency of a photoelectron chemical reactance, e.g. semiconductor colloids, can be improved by modifying the surface of the semiconductor particles being manipulated with a noble metal like platinum or ruthenium [11].

The surface, its orientation, azimuthal orientation and patterns, particularly of proteins, is becoming one of the most important aspects of nanoscale biological assemblies. As before this is not roughness in the accepted sense but it is certainly surface influence and behaviour [12].

An interesting point is that the coverage of the surface with proteins or enzymes has to be high to get high activity. However, it has been found that there also has to be a certain freedom of movement to allow proper protein fixing and settling as, if the coverage gets too high, it can inhibit activity.

From the above it is clear that surface effects can be vital in biological and chemical activity at the nanometre level. Roughness is not negligible. What happens is that certain different aspects of roughness are more important in non-engineering applications than they are in conventional engineering.


### 8.3.4 Surface and bulk properties dependence on scale

This scenario is about the most difficult to understand and it puts extreme constraints on surface boundaries and bulk volume which have to be interpreted in terms of analytical or operational definitions. Both the surface and the bulk have properties which are scale of size dependent.

The roughness of nanopowder and clusters is not, as already mentioned, significant. It is like asking for the height of the tallest hill on an asteroid; it is not relevant. The real problem is that of controlling the production of such particles such that 'reactance' and stability can be maintained. These are determined to a large extent by the surface properties, which are more chemical than geometrical. The actual reactions are a function of the first atomic layer in the surface and its interactions with the outside world. As in tribology, the geometry determines the amount of 'exposure'. The chemistry determines what happens.

It has been shown that 'clean' surfaces can be obtained on single crystals by cleavage. This property is used for step height calibration. The environment is not critical. However, in the reactance case the environment is critical; cleavage should only take place in a vacuum to limit oxidation or other extraneous chemical reactions.

When exposed to the atmosphere various chemicals get adsorbed into the surface and can sometimes diffuse to the bulk material. Adsorption can drastically change electrical and magnetic properties. Notice that it is the surface which is of paramount importance not the bulk.



**Figure 8.17** Surface atoms as ratio of bulk atoms.

The ratio of surface to bulk atoms is very dependent on size as can be seen in the figure. In general $R$ is very low $\sim 10^{-5}$ for cluster-like particles of spherical shape but for small diameters $\sim$1-2nm $R$ can be 0.5. The actual ratio depends on atom or molecule size.

As well as size, the density of defects on the surface is important. If this gets high the defect effect can swamp the characteristics. For this reason cleanliness is very important in preparing surfaces.

Defects take over from roughness as the important surface characteristic in the nanometre and sub-nanometre domains.

Defects can be characterized in terms of physical dimension. For example zero dimensional (point) defects have negligible spatial size, one dimensional defects have significant size in one direction only and so on. The maximum dimension of any defect depends on that of the parent sample which obviously fixes the upper limit. At the atomic scale every defect is three dimensional but at a very small size. This is not taken to be significant on a macroscale; another case where the definition changes with scale! At what size this change of definition occurs is open to question. It is probably acceptable to take the nanoscale as the breakpoint.

The other sort of point defect is the vacancy or omission of atoms or molecules, which causes some distortion of the surrounding lattice. It is sometimes necessary to describe the distortion in the lattice of a crystalline material due to a defect. The description is scalar, vectorial and tensorial. The orientation of a dislocation for example is the Burgers vector, which specifies the amplitude and the orientation of the distortion.

### 8.3.5 Shape and (scale of size)

#### 8.3.5.1 Engineering shape

There are a number of shapes in engineering which are especially relevant. Of these roundness is probably the most important. Practically every dynamic application involves rotation. The measurement and control of roundness is well understood at the macro level as demonstrated in chapter 4. However, as the size reduces it becomes more difficult, at the millimetre level and below, to include MEMS. The measurement is tricky. This

is mainly due to centring errors. Also, the signal as seen by the instrument is not the same as that obtained at the macro level. Referring back to the roundness signal $p(\theta)$

$$\rho(\theta) = e\cos(\theta - \varphi) + \sqrt{R^2 - e^2\sin^2(\theta - \varphi)} \tag{a}$$

which for $R \gg e$ (where $R$ is the radius of the part and $e$ is its eccentricity relative to the centre of rotation of the instrument) becomes a limaçon

$$\rho(\theta) = R + e\cos(\theta - \varphi). \tag{b}$$

However, in the miniature case $R$ and $e$ can be the same magnitude because $R$ is small and $e$ is often larger than at the macro level because of filtering difficulties, so

$$\rho(\theta) = R + e\cos(\theta - \varphi) - \frac{e^2}{2R}\sin^2(\theta - \varphi) \tag{c}$$

which has a second harmonic term, so that harmonic analysis to get rid of the eccentric term is no longer a simple operation.

Offset against this is the fact that the circular shape is quite easy to generate at the miniature level. An example is the use of a shower ion beam which uses a collimation beam of ions rather than a focused one.



**Figure 8.18**

Figure 8.18 shows a rough cylinder, e.g. a micro-wave stub rotated randomly within the shower. Invoking the central limit theory to describe the machining mechanism allows a circular part to be produced naturally as illustrated in figure 8.18(*b*).

For machined objects shape is usually determined by the movement of the tool. This involves movement in two or three dimensions. The curve produced in space is usually simple, conic forms such as parabolas and spheres are often required. In chapter 7 one application, aspherics, explains the complicated curve sometimes needed to correct the optical wavefront before the object plane. The aspheric form can conveniently be included within the understanding of surface shape but there could be a problem in the nanometrology regime because the definition of ideal shape can be completely different from the conventional Euclidian one.

### 8.5.3.2  *Molecular shape and scale of size*

Take for example the concept of shape in the cluster and molecular regimes. An example of the estimated shape of particle shapes or molecular groups has been given the symbol $R_s$ [13].

$$R_s = \left(\frac{L^2}{48N} - 1\right) \tag{8.8}$$

where $L$ is the length of a periphery of the cluster and $N$ is the number of atoms or molecules enclosed within it.

In equation 8.8 the number 48 is included to reflect the concept that the regular hexagon is taken to be the perfect shape giving $R_s = 0$. Presumably, the hexagonal shape of the benzene ring is taken as the reference. For an engineer the perfect shape is more likely to be a circle or straight line.

Another aspect of shape at the nanoscale is that it can have an intrinsic meaning i.e. it is predetermined by energy or other considerations. Whereas in engineering the shape usually determines whether the workpiece can move (i.e. roundness can determine how well a shaft revolves in a bearing). At the nanoscale, the shape of an object can be a result of the internal energy balance between, say, chemical and elastic energies. In clusters or even liquid droplets the energy balance is between Coulomb attraction and surface energy. From these comments it is clear that in nanotechnology shape is often an output from energy considerations rather than an input to an engineering function. Shape can therefore have completely different implications in macro- as compared with nano-regimes.

In the molecular domain shape is determined by the chemical composition and the structures have a fixed shape depending on the particular molecular configuration. In other words, the shape of a structure is preordained. The question is whether or not it is inhibited by defects. Shapes are discrete, not continuous as in macro-engineering. The measurement of shape therefore is usually not the issue: failure to achieve the expected shape is.

In quality control terms in the nano- and molecular regimes shape is often an attribute and not a variable, e.g. is it the correct shape rather than the 'value' of the shape?

The shape of interfaces can also influence the quantum behaviour of materials. The solution of Schrödinger's electron wave equation can be influenced by the boundary conditions; an ideal parabolic shape for the potential barrier in a quantum well depends on the shape of the boundary, which could be an interface or a surface.

### 8.3.5.3 *Molecular and biological shape considerations*

Lithographic methods have been used effectively in fabricating computer elements but they are approaching a limit: it is difficult so see how much smaller silicon based binary switches can be made. Also detection and control of defects during fabrication gets more difficult.

The molecular engineers [14] suggest using supramolecules (where two or more sub-units come together) or more generally molecular machines (i.e. to synthesize from simple chemical blocks more intricate molecules which respond rapidly to small controlled changes in their environment).

Ideally the switching mechanism should not be based on covalent bond formations because they are slow and the mechanism is complicated. Noncovalent interactions can be much faster. Molecular self assemblies have been suggested. These are mainly based on hydrogen bonded organic assemblies. However, more interest is being shown now in transition metal based mediation. The result is that they are now being systematically employed in the construction of discrete nanoscale sized structures of well defined shapes and sizes.

An example is the molecular square, which derives from a two dimensional polygon. The actual shape of the sub-unit is determined by the turning angle at the chain link. For example 90° for molecular square, 60° for molecular triangle, 108° for pentagon and 120° for hexagon (figure 8.19).

To see how this works a molecular triangle requires three angular and three linear components assembled in a cyclic manner.

The ability to readily interchange the building blocks is an advantage of self assembled systems.



**Figure 8.19** Molecular geometries.

Probably, apart from the very small sizes, it is the fact that such assemblies are defect-free otherwise there is no connection: the shapes are precise. If they exist they are perfect. Molecular machines such as supramolecules have only been used in very simple structures such as switches or shutters.

3D forms based on these sub-units are possible, for example, a tetrahedron or a helix. In fact most biomolecules exhibit some recognizable 3D spatial orientation.

Proteins are used for self assembly sub-units. Their sensitivity to small changes is the environment makes them suitable for use in microcircuits and sensors. A discussion of the various proteins and their self-assembly mechanism is outside the scope of this book. However, reference [14] contains information.

One of the difficulties with making self-assembling nanostructures is the condition of the substrate. The most obvious factors are the roughness of the substrate and the choice of immobilizing agent. Other points are the surface energy and the chemical nature of the surface.

### 8.3.5.4  Carbon molecular shapes

It is due to their extremely precise shapes and sizes that carbon nanotubes and similar structures such as fullerenes have such remarkable electronic and mechanical properties. In principle they can be used in gears, rack and pinions (figure 8.20). It seems incredible that these structures relate to the parent planar graphite.



**Figure 8.20** Simulated shapes for rack and pinion.

Some enhanced mechanical properties include better stability, strength, stiffness and elastic modulus. The helicity, which is the arrangement of carbon hexagons on their surface layer honeycomb lattice, is determined by the symmetry of the structure, the helix angle and the diameter of the tube (i.e. the shape and size).

It should be noted that it is not only carbon which has the potential for forming tubes. Molybdenum disulphide ($MoS_2$) the well-known lubricant also has a similar structure.

### 8.3.5.5  Simulation and scaling up of molecular machine elements

Another aspect of scaling relates to the bottom-up regime rather than the top-down regime, in which the effect of scaling on dynamic systems has already been discussed. This aspect is simulation. Ultimately, atomic behaviour has to be related to nanobehaviour which in turn has to be related to macrobehaviour. Because of the difficulty of dealing theoretically with many variables and boundaries, it rests with computer simulation such as molecular dynamics to do the extrapolation from atomic to behaviour. Unfortunately the computation depends heavily on the number of units $N$ (say atoms) and levels $n$ (basis functions i.e. energy levels).

The Hartree-Fock procedure results in a total of $(n^2 + n)/2$ kinetic and nuclear attractions together with $(n^4 + 2n^3 + 3n^2 + 2n)/8$ repulsion integrals so that the complexity is $(10N)!$ if $n = 10$. In practice, the over all complexity is proportional to $n^4$ and $N!$ both of which are enormous quantities for very modest values of $n$ and $N$. Difficulty can be experienced with as few as 20 atoms. Different strategies have to be adapted

depending on *N* and *n* in order to make the simulation tractable [13, 14]. Computer time has to be balanced against accuracy.

As *N* and *n* increase the possibility of errors creeps up. One of the biggest problems is that of checking the programme. Satisfactory behaviour at $N = 20$ does not guarantee it at $N = 200$. It is well to remember that many critical phenomena occur at the surface and therefore *N* need not take on bulk values. So surface effects can more readily be simulated than bulk effects. It should be possible to extrapolate bulk properties from surface properties; the bulk can be considered to be a 'stack' of surfaces.

At the heart of the problem of predicting, macro-behaviour from atomic properties is the mixing of quantum mechanics with classical mechanics. One way of tackling this type of problem is to separate it into two components: the unit event and the distribution of events. This is the procedure for contact as shown in chapter 7. Instead of a 'contact' being the unit event the interaction of electron fields of two atoms could be the atomic unit event and their distribution the macro-picture. The unit event here is evaluated rigorously by quantum/wave mechanics using Schrödinger's equation and the latter conventionally worked out with molecular dynamic simulations. In other words, the breakdown of a problem in tribology could well be the model for the nano-surface behaviour.

### 8.3.5.6   *General comments on the scale of size*

The metrological factors in equation 8.7 all change in some way as the size reduces to below the nanometre level. Remember that the balance of forces also changes with size, with each component changing disproportionately. However, the metrological components above change form rather than value. Roughness, as usually defined if it exists, is not a feature of great importance— or so the lack of technical papers on it suggest—substrate roughness seems to be the exception. For the purely mechanical applications in microdynamics, roughness, shape (form) and size become scrambled. It is at this scale of size and with mechanical applications it is becoming more sensible to lump the components together as in *G* above. Forces are small enough not to have to consider sub-surface stress and not to worry about waviness. The top-down approach for miniaturization keeps to its traditional roles but shape and size become discrete, in examples of molecular synthesis in nanotube technology, the shape and size are fixed in advance.

Figure 8.4 illustrates the diversity of the geometric relationships. What it does not show is the difficulty of establishing a common framework of standards to back the different types of measurements. Whilst it is easy to understand the need for standardizing and establishing traceability for size (length, position) it is not so easy to do the same for roughness and shape. Length calibrations are included in shape and roughness but in addition, in the latter two, there is a frequency element that has to be included. It could be argued that roughness and shape (or form) have been extensively standardized in the past, as set down in ISO Handbook [16] so what is the problem? It is that at the very small dimensions involved in nanosurface metrology the availability of length and shape artefacts to act as standards is limited, as well as the extra environmental problems of noise, cleanliness etc. But, is the measurand, e.g. magnetic flux, traceable?

Also, as the size decreases the nature of the signal changes. What is a straightforward deterministic measurement in engineering gets replaced by spatial statistical averaging and, finally, temporal statistics brought about because of quantum mechanics. This is explained by figure 8.21 in section 8.4.

## 8.4   Stability of signal from metrology instruments as function of scale

### 8.4.1   *Length*

Having decided that the meaning of the geometrical components can change with size, is there anything else which changes? Unfortunately the answer to this question is yes. It is the means whereby a stable signal is obtained from the test piece by the instrument [16].

(a) Deterministic

Engineering metrology

$S$

$$\dfrac{S}{N} \gg 1$$

(b) Spatial integration

Nanometrology

$$\dfrac{\int S\mathrm{d}x}{N} > 1$$

$m_1$ $S$ $m_2$

$\uparrow x$

(c) Temporal integration

$p_1$ $p_2$

$S$

Quantum mechanics metrology

S is signal

N is noise

$$\dfrac{\int S\mathrm{d}t}{N} > 1$$

**Figure 8.21** Signal form.

Figure 8.21 illustrates the point. In (a) the system is conventional engineering and is stable relative to the signal $S$. Uncertainty is small so $N'$ is small. What is seen as the signal $S$ is acceptable. In (b) the actual signal at any point changes with position so, to get a meaningful signal $S$, the geometry has to be smoothed by integration giving $m_1$ and $m_2$. The signal is now $S = m_1 - m_2$ where $m_1$ and $m_2$ are mean planes. The distance $S$ can be accurate to much less than a nanometre because of the spatial integration.

In (c) at molecular and atomic levels the available signals passing between the two points is limited by the distance between points $p_1$ and $p_2$ which are molecular or atomic in size. In this regime tunnelling across the gap or through a barrier is subject to quantum laws. The only way to increase the probability of an electron or photon moving between $p_1$ and $p_2$ and so closing the measurement loop is to increase the observation time (i.e. to integrate temporally). This automatically makes the instrument sensitive to the environment. This temporal problem is not the same as the bandwidth of electrical noise, which will be discussed later in section 8.8.

Figure 8.22 shows how the signal processing of figure 8.4 fits into the scales of size as shown as the ordinate axis. The lateral plane has an axis in terms of areal information on the surface whereas the other axis is of heights of the surface. These particular lateral arguments are included to illustrate how the various instruments used for surface and nanometrology can be grouped in terms of capability. Figure 8.23 is virtually the same except that the scale of size has been replaced by time. Although the instrument groups are converging on the same specification, there has been a tendency to concentrate on spatial rather than height information; improvements in height information have been straightforward but this is not the case for lateral information. Stylus methods have had to improve considerably from the simple area scans. Just having an areal picture is not enough. Technology based on semiconductor requirements, such as lithography, demands

**Figure 8.22** Measurement for different types of instrument as function of scale of size.



**Figure 8.23** Measurement trends of surface-nano instruments

precise position and length monitoring. Now that high resolution is possible with SPM the pressure to position exactly has increased dramatically.

The increase in resolution and accuracy of modern instruments forced by application pressure has brought their specification into the nano regime—in effect almost the same scale of size spatially as that of the surface roughness. This is why the subjects of surface metrology and nanometrology are blending and why position and spacing are being lumped into the new discipline of surface nanometrology.

### 8.4.2 Nano position sensing

This is defined as the technology of moving and measuring with subnanometre precision [17] and has values which are comparable to very fine surface finish. It is also obviously related to length. Optical methods can be used to measure position using white light fringes as in chapter 4. However, a very common method is the use of piezoelectric crystals suitably energized by a known voltage to produce the movement, and some gauge, usually capacitive to control the movement. The control is necessary because, although piezoelectric

crystals are very stiff, they suffer greatly from non-linearity and hysteresis. It needs a very precise and accurate gauge to find out where the end of the crystal is in space!

As there are three components involved in the geometry of a workpiece each having its idiosyncrasy it is wise to review the calibration procedure in order to see if the three can be brought together in a traceable way at the nanoscale.

## 8.5 Calibration

### 8.5.1 General

'Metrology is only as good as the calibration procedure.' This statement is just as true for nanometrology as it is for conventional engineering and physics. Quite often calibration procedures have evolved rather than having been thought out.

The real problem lies in the way engineering surface metrology and microscopy have evolved. Traditionally engineering applications, such as those which occur in tribology, have had more weight placed on height variations than on spacing or lateral dimension. Microscopists on the other hand have looked mainly for areal structure and position in the plane of the object. This dichotomy has resulted in two different traditions of calibration which are only now being integrated [18] (figure 8.23). This shows that although starting from different backgrounds all methods are converging into instruments which have height and spacing capability to about the same degree.

### 8.5.2 Length and position calibration and measurement at the nanoscale

Following Teague [19] there are three types of length calibration artefacts. The simplest is *line spacing* (figure 8.24). In this type of calibration the effect of the shape of the probe and its interaction with the surface are small because any aberration introduced due to the interaction is common to both lines if measured from the same direction. Another artefact is line width and a third is size or extension. These three types are shown in figure 8.24.

Take first the case of the contacting probe. The finite slope of the probe causes the sharp edges of the line spacing (a) and the line width (b) artefacts to become indistinct (bd). The edges can still be seen but the signal



**Figure 8.24** Length standards showing the effect of the interaction between the probe and the artefact.

needs some processing to ensure that the edge can be identified. This is usually achieved by differentiation. In the case of the line spacing (gd) the tip dimension is not critical but it is in the case of the line width standard (ae). In this case the only way that the line width can be estimated is by taking the tip size (if known) into account by crude deconvolution. Both these artefacts suffer from a degree of integration (i.e. smoothing). In the case of the extension (af) the distance is measured between two contacting probes. The problem here is that the finite contact force at both ends tends to make the extension ($l$) smaller than it really is. The actual deformation can be worked out using Hertz equations given in chapter 4 and is shown graphically in figure (bf) giving the modified extension $l'$.

In the case of non-contacting probes utilizing electromagnetic the problems are of a different nature. Diffraction effects cause sharp edges to become enhanced. This edge enhancement can be regarded as advantageous if spacings are being measured but deleterious if step height is needed. Contact probes are the exact opposite!

As in the use of contact probes, measuring extension as seen in figure 8.24 (cf) suffers from distortions. These exist because of the penetration of the waves into the surface to an extent determined by Maxwell's equations. It is not a question of which method is best. Both contact and non-contact are correct according to the rules with which the value of spacing (or height) is obtained. From what has been said contact methods tend to integrate whereas optical methods differentiate. These discrepancies are usually small enough to be ignored when the standards are in general use. However, in the realm of nanometrology these errors can be important so it is evident that more appropriate ways need to be found. At the very least the composite geometry and physical properties of probe and standard have to be known. In the case of contact probe radius, slope and composite elastic modulus are named and, in the case of the optical method, the numerical aperture of the optic, the wavelength of light and the electrical conductivity of the material need to be known so that compensation can be made.

SEM techniques have similar problems to optical methods but also extra ones because asymmetric images can occur due to the position of the detector relative to the incident beam (figure 8.25).

The actual image produced can easily differ from the true surface shape, as shown in figure 8.25.

In figure 8.25(a) secondary electrons ejected due to the primary electrons hitting the surface (from the gun) have to drift over the hump in order to be received at the detector (figure 8.25(a)). Some of these electrons inevitably get lost. On the other hand, when the gun electrons are on the side of the hump near to the detector, all get captured (figure 8.25(b)). This means that the apparent signal (c) at the detector is asymmetrical relative to the hump (figure 8.25(c)): a characteristic of the position of the detector rather than the shape of the surface!



**Figure 8.25** Distortion of SEM image due to detector position.

Scanning tunnelling microscopes using a probe can also produce false results, as can be seen in figure 8.26.

This shows tip hopping which can occur with an imperfect probe. Although this effect looks obvious from the figure it is very difficult to detect in practice.

In the new scanning microscopes the probe size and shape is critical because the tip itself is of the same order of magnitude as the features being examined (i.e. nanometres).

Because of the importance of the tip, various methods have been employed either to measure the tip or to deduce the real surface profile from the measured one—assuming a stylus shape and size. Methods of doing this are described in chapter 5.



**Figure 8.26** Distortion of STM picture due to poor probe shape (peak hopping).

Recent work [20] uses the concept of surface plasmons. These are charge oscillations at a metal surface which can couple to the incoming on-plane wavefront of light, thereby producing anomalies in the angular reflectance dependency of the object. By using a metal coated diffraction grating it is argued that, not only can the surface profile be obtained, but also the probe curvature—by a process of de-convolution. The problem is that these plasmons are only produced at metal/dielectric interfaces so that the technique is hardly practical for day to day purposes. Also the specimen has to be smooth enough and have spacings wide enough to ensure that the probe always bottoms. This is the exact opposite of the groove specimen mentioned earlier, in which the probe cannot be allowed to bottom. In practice a bit of both situations occur and as a result of this the 'inversion problem' cannot be adequately solved.

### 8.5.3 The probe

The probe is so important in the new generation of scanning microscopes that a lot of effort has gone into their manufacture. In particular the controlled manufacture of the geometry of the tip has been an objective. To this end 'electrochemical' etching techniques (ECE) have been developed [21] to fabricate tungsten and platinum/iridium tips for use in STM. The tips have specifically been made with a high aspect ratio (i.e. subtend a very small angle at the surface). In addition much effort has been expended to keep the radius of curvature small.

The technique of ECE to make STM is a natural progression from field ion microscopy. The use of platinum rather than tungsten for the tips is due to the fact that platinum, although being a softer material, is inert to oxidation. Iridium is infused to make the tip stiffer.

Other unconventional STM tip preparations have been explored. These include ion sputter thinning to get better resolution and reliability. Tips capable of achieving atomic resolution spatially have been fabricated out of pencil lead by coating tungsten tips with colloidal graphite. Also shearing Tt/Ir wire gets a cut tip which is often good enough to use in measuring roughness in micro-areas of silicon for example.

Possibly the most interesting probes for STM are carbon nanotubes. These, intrinsically, seem suitable for nanometrology because they have a small diameter, large aspect ratio $\left(\dfrac{\text{Length}}{\text{Diameter}}\right)$ and have high stiffness to lateral movement. The only difficulty is that of attaching them to the body of the pick-up [22]. This is a similar problem to that of attaching carbon nanotubes to pickup shanks.

Figure 8.27 shows just how sharp STM probes can now be. What is needed is a method of avoiding damage to the tips and even knowledge of the damage when it occurs!

Figure 8.28 illustrates the relative sizes of probes.

(*a*)                                    (*b*)



**Figure 8.27**  STM tips of less than 100 nm radius.



**Figure 8.28**  (Where LS is line spacing and LW line width).

There are optical 'super tips' [23], which can be smaller than 10 nm, possibly single molecules that can absorb light energy and transcribe to different optical properties.

Special silicon tips capable of being used in the fabrication of semiconductor devices, e.g. by atomic force microscopy, have such small size that movements atom by atom are claimed, resulting in possible locations of $10^{14}$ GB/m2 on silicon wafers. Angles of 10° at the very tip are possible.

**Table 8.3** Probe limitations in nanometrology. (Uncertainties (nm) for Measurements)

|  | Line Spacing | Line Width | Extension |
|---|---|---|---|
| **STM** |  |  |  |
| Diameter, d, of Probe-Specimen Interaction = 0.2 nm |  |  |  |
| Bending, B, of probe shaft = 5 x $10^{-2}$ nm | $\sqrt{2} \times 0.1$ d = $1.4 \times 10^{-2}$ | $2 \times 0.1$ d + 2B + = $0.15 - 0.2$ | 2B + 2U + tip dia. Uncer. = 1.2 |
| Uncertainty, U, in location of surface = 0.1 nm |  |  |  |
| **SEM** |  |  |  |
| 1 kV onto features on silicon, 30 nm electron range, ER | $\sqrt{2} \times 0.1$ ER = 4.2 | $2 \times 0.1$ ER to $2 \times$ ER = $6 - 60$ | ~~~~~ |

As many nanometre devices use a contact probe much effort has gone into making probes suitable for a wide range of applications. NPL [24] has developed a 3D low force probe incorporating miniature capacitance transducers and a light weight flexure structure made from tungsten carbide and beryllium-copper strips. The probe assembly weighs about 0.33 g and has a 0.1 mN probing force which allows contact to be made without significant damage. This probe is attached to a hybrid CMM (coordinate measuring machine).

### 8.5.4   *Height measurement calibration at the nanometre scale of size*

It was soon realized that, for semi-conductors and in general nanotechnology, traditional height measurement standards were too coarse. For this reason height artefacts based on molecular or atomic constants were sought (see chapter 5). As early as 1973 a technique using lattice spacings was tried. First mica was cleaved mechanically, the idea being that, when cleaved, some jumps of lattice spacing would be present on an otherwise molecularly smooth surface (figure 8.29). This technique had to be rejected because the mica was springy and often introduced long wavelengths into the trace due to the lack of stiffness under the load of a typical stylus probe. Another crystal, that of topaz was used which was very hard and seemed satisfactory but was expensive. In this method it is not necessary for unit lattice faults to be found. Multiple lattice spacings and the method of exact fractions can be used. It seems highly possible that some of the new ceramics might be used instead of topaz. Even rock salt is a possibility.

One big problem associated with height calibration at the nanometre level is that there is a tendency to try to relate the calibration procedures devised for contact methods with those of optical methods. In particular, optical methods have been used to measure artefacts developed for stylus contact methods. This attempted comparison has been carried out for roughness standards. If the surface is rough the optical and mechanical probe gives the surface geometry fairly adequately. The difference between optical and mechanical techniques is not usually significant. But if the surface is smooth then even small differences can cause a serious divergence between the two methods.

Although it seems sensible to use optical interferometry to measure position, displacement [25] and step height, typical results with interferometry where d is measured height gave an uncertainty of measurement of $10^{-4}$d nm + 2.5 nm whereas stylus methods gives 1.5 x $10^{-3}$d nm + 0.7 nm. These are typical results. Errors in

**Figure 8.29** Cleaved topaz.

the stylus method (Nanostep) can be attributed to small non-linearities in the transducer, whereas the errors in the optical methods such as Wyko, Zygo, Zeiss and UBM are usually attributed to difficulties in fringe position measurement, surface contamination and roughness of the surfaces of the steps. Attempts to improve the evaluation of the step heights by means of software have been used extensively. One [26] uses a step height as part of the least squares best fit algorithm and gets remarkably consistent results.

Many articles in the literature compare different instruments for step height standards. Some use the average of many readings as a basis for the comparison. This approach should have disappeared years ago. The true way to compare instruments properly is to measure along *exactly* the same track by using kinematic methods together with track identification techniques, for example, the double cross method (figure 8.30).



**Figure 8.30** Relocation method for comparing different instruments.

Using this relocation method there is no ambiguity of trace position. The yaw angle is adjusted until just the two marks appear on the trace. One reason why there is sometimes noise present, which buries the signal is because there is a mismatch between what is being measured and the metrology unit. One example of this mismatch is the use of light having a nominal wavelength of $0.5\mu$m to measure nano detail; the mismatch in this case being about 500:1. This gives a very low signal to noise ratio which is never advantageous. Also, there is a limit to what can be extracted from the noise despite clever algorithms. Earlier it was suggested that the crystal lattice spacing is an appropriate unit being comparable with nanoscale features. Another idea is to build an X-ray interferometer [27, 28]. The actual device is shown in figure 8.31(a) and (b).

In this arrangement the interferometer is a monolithic block of silicon. This is machined to project three blades from one face (figure 8.31(a)). The third blade is part of a ligament hinge. X-rays enter from the left in the figure 8.31(b) and are focused via the second blade B onto C, the movable blade. Movement of C varies the absorption of X-rays by the blade. This is detected by a scintillation counter.

If C is moved linearly upward the voltage output varies sinusoidally. This signal is the master which, when related back to the lattice spacing can be used to compare with the test transducer signal whose contact point or focal point is positioned on the top of the moving ligament hinge configuration of the crystal. The actual value of the X-ray wavelength is not critical but the parallel faces of the blades are. Interpolation of the sinusoidal signal gives much better than sub-nanometre accuracy and resolution.

(a)

80 mm

25 mm

60 mm

Drive

**Figure 8.31(a)**

(b)



A          B          C

Incident
x−ray beam

R          RR

T                    H

O

Diffracting planes          d=0.3135 nm

**Figure 8.31(b)** X-ray interferometer.

This X-ray interferometer is very small. The crystal itself is a few cubic centimetres so that it can be carried from place to place. Silicon is an almost perfect material for such an interferometer because it does not suffer from hysteresis and is perfectly elastic. This means that if it works at all it has to be perfect. Overstrain simply snaps the hinge [29].

The use of X-ray interferometers in the calibration of nanometre instruments has recently been extended at the National Physical Laboratory [24]. Their basic premise is that interpolating fringes from an optical interferometer presupposes that the fringe has a perfect form so it is beneficial to linearly interpolate the optical fringe using an X-ray interferometer [30].

Such a system makes use of the attributes of both types of interferometer. It has a long range of 1 mm using optical fringes of 158 nm together with the high resolution of the X-ray interferometer of 0.192 nm.

## 8.6  Noise

### 8.6.1  Signal limitations

Apart from the problems of movement and positioning the signal itself is subject to noise. This noise can be due to electronic, thermal or spatial effects.

### Electronic

This is due to the electronic fluctuations inherent in the input stage of the amplifier or in the probe itself.

The basic electronic noise is Johnson noise and shot noise [31].

The equation of RMS voltage for Johnson noise is

$$E\left(V_j^2\right) = \frac{KT\Delta B}{R}$$

(8.9)

For a typical system this noise is equivalent to subnanometres.

Shot noise is the variation in the current due to the random movement of electrons. It has a bandwidth $\Delta B$ as above and gives about the same level of equivalent noise as Johnson noise.

### 8.6.2 Thermal effects

Brownian motion can cause deflection of the end of an AFM cantilever [32]. Calculations show that for a simple detection system based on viewing the end of the cantilever, the expected variation in $z$, (i.e. $z$ is given by equation 8.10, where $E$ is expectation.

$$E\left(\Delta z^2\right) = \sqrt{\frac{kt}{K}} = \frac{.065nm}{\sqrt{K}} \tag{8.10}$$

$K$ is $0.25\ E\ wh^2/L^3$ where $h$ is thickness, $w$ is width and $L$ is the length of the cantilever.

In cases where the deflection is measured by a laser beam reflected from the cantilever end, e.g. an optical lever, then $dz(L)/dx$ is the operative measurand.

$$\text{thus} \quad Z(L) = \frac{2L}{3}\frac{dz(L)}{dx} \tag{8.11}$$

By applying thermal vibration theory to cover the two cases when (a) the one end of the cantilever is free and (b) when it is making contact with the surface, it is possible to deduce the different modes of vibration of the beam using standard formulae.

When all vibrational modes are taken into account the root mean square deflection $\sqrt{\overline{z}^2}$ is, given that the optical lever is used for measurement,

where
$$\left.\begin{array}{ll} \sqrt{\dfrac{4kt}{rK}} & \text{for free end} \\[3mm] \sqrt{\dfrac{kt}{3K}} & \text{for contact end} \end{array}\right\} \tag{8.12}$$

$k$ is Boltzman's constant, $t$ is temperature in degrees Kelvin.

This shows that the RMS value of variation for the free end is twice that of the contacting cantilever.

For a temperature of 22°C and a spring constant such as given by $K$ in Nm the values of the resultant thermal noise are just less than $10^{-10}$ m. This means that the mechanical noise and the electronic noise are about the same and are uncomfortably close to the desired resolution. For small structures, e.g. small cantilevers for AFM, the internal friction has been broken down into surface effects and bulk effects, both of which dissipate energy independently showing yet again the importance of the surface and its roughness [33].

## 8.7 Calibration artefacts

It has always been difficult to provide specimens with which to calibrate precision instruments. In engineering applications the obvious example of the use of artefacts is the use of gauge blocks, which are used as a practical intermediary between the test piece and length standards whose values can be traced back to the standard metre. Unfortunately, while these gauge blocks were suitable for use in most engineering practices of the early part of the century their use is restricted today. This has meant that below the range of micrometres and extending down to nanometres a real gap has emerged. Despite many efforts the need has not yet been adequately met.

One of the problems is the fact that there is a need for an artefact of regular pattern. This type is needed for nano-structure spacing as well as for lateral information in scanning microscopes and for filter characteristics in micro-roughness.

Traditionally, artefacts for roughness or structure have been made by using ruling techniques common in diffraction grating manufacture. Nowadays diamond turning has been used to generate sine waves in copper. These are usually chrome plated to give the artefact hard-wearing properties. Unfortunately this type of standard cannot be used equally effectively for stylus methods and optical methods. This is because the stylus method tends to ignore remnant turning marks in the surface (figure 8.32(a) and (b) whereas the light in the optical techniques gets diffracted (figure 8.32(c)by such marks).



**Figure 8.32** Attempt to calibrate optical instruments with tactile standard.

This illustrates the difficulty of getting artefacts that are suitable for both stylus and optical calibration. One interesting technique for making artefacts within the sub-micron to micron range has been described [32]. This range is extremely close to the nanometer range so sometimes these artefacts are tried out with bizarre results such as shown above.

This method is quite different from that using micro-spheres made of latex [26]. In this new technique use is made of the fact that phospho-silicate glass (PSG) can thermally re-flow at comparatively low temperatures ~850°C-1100°C. This property can be used for making artefacts but with a limited pitch capability.

Dimensional accuracy can be achieved using conventional techniques and the micro-sphere (the latter melted) suffers from various problems [32] which include the lack of a straight edge. Also temperature effects limit the use of photopolymers which can only satisfactorily be used in replica form.

For these and other reasons artefacts made by conventional methods do not take spherical shapes or assume accurate periodic surfaces.

With PSG stable curved surfaces and surfaces in a variety of shapes down to the nanometre scale (Figs. 8.33(a), (b), and (c)) can be achieved.

The method involves depositing PSG onto silicon wafers by the chemical vapour deposition process. Electron beam methods allow a periodic structure to be formed on the surface. This smoothes out to a sinusoidal shape when the wafer is heated to 900°. Periodicities down to 0.5 $\mu$m are easily achievable, which is getting close to a useful nanometre standard.

Hemispherical and oblate hemispherical shapes can be produced with very smooth surfaces by using a variant of the above method. This involves using PSG with a thin layer of poly-silicon.

Other techniques for complex artefacts include dry etching of $SiF_6$, which enables cusp-like periodicities to be formed with spacings of 0.5 $\mu$m or better.

Wet etching of silicon in potassium hydroxide can produce a very realistic vee-shaped periodic surface. The actual shape depends on the mask design (figure 8.33(d)).

It is interesting to note that the very first roughness standard had a triangular shape. General Motors used this standard, which was ruled in copper and had a chrome layer deposited on it (Calyblock).

This standard was found to be especially useful even in 1936 because it could be used as a height standard, a roughness standard and also for stylus wear. It did the latter by means of a finely ruled standard of 2 $\mu$m p-v and $R_a$ (arithmetic average) values of less than 0.5 $\mu$m (the true value) obtained from the standard were a measure of diamond wear.

Ruled or diamond-turned periodic surfaces are also used to test the high pass filters used in surface instruments. A range of about 3:1 in spacing can be obtained simply by tracking the stylus at an angle across the ruled standard (figure 8.34). In this method the instrument has to have high cross axial stiffness and the stylus should ideally not be rectangular at the tip.



Track 1

Track 2

**Figure 8.33**



(*a*) Sinusoidal grating structure



(*b*) Vertical periodic



(*c*) Reflowed PSG hemisphere



(*d*) Sawtooth-like grating

**Figure 8.34** Various attempts at generating curved and periodic surfaces at the nanometer level.

Various attempts have been made to calibrate the height and distance accurately. One designed to calibrate sinusoidal standards has recently been developed at NIST. The probe height is calibrated relative to slip gauges and the position by means of a laser interferometer [33]. By this means a calibration of both the height parameters, e.g $R_a$ of 1 $\mu$m to ± 1.2% and $S_m$ of 800$\mu$m to ± 0.06%.

## 8.8 Dynamics of calibration at nanometre level

Calibration of length rarely involves consideration of dynamics but surface metrology does, the probe having to collect many successive data points in any one measurement of, say, roughness or roundness. It is very difficult to achieve the necessary control of position accurately. Sometimes however, it is possible to go half-way towards a calibration under dynamic conditions. In a previous section reference was made to sinusoidal specimens for use in roughness calibration. These standards are still difficult to make with a small spread so alternative methods have been advocated for example by putting the probe of the roughness instrument on a calibrated vibrating table whose movement is perpendicular to the traverse [34]. This movement is highly controllable by means of an interferometric system [35]. Unfortunately the movement of the probe along the surface, with its attendant motor and gearbox noise is not taken into account. Actually the problem of roughness calibration should be spatial and not temporal. The temporal nature of the instrument (e.g. constant speed tracking) sometimes introduces problems rather than solving them.

One of the few attempts to take into account the translation dynamics was made by Futami [36], who designed a complex system involving non-contacting translating units with a linear bearing and using a ball bearing guide. His paper takes account of the different equations of motion from start-up to free movement and demonstrates extremely competent design. Taking the dynamics into account truly allows nano-positioning without backlash on fine control.

One of the usual modes of operation of the AFM is vibrating the cantilever. This has the effect that surface forces exerted on it give an apparent change in the frequency of oscillation that enables the force to be estimated. This has been carried further in a development called impact oscillation to stop stiction. There is meant to be an impact at the tip in order to promote the correct mode [37, 38], which becomes in effect a tapping operation and, in principle, improves the sensitivity.

Various attempts have been made to make SPM techniques more accurate. The basic problem is the hysteresis and non-linearities of the piezoelectric movements, so for positional accuracy there has to be a recourse to independent position gauges which are usually capacitive. Also one trick is to put the signal voltage onto two identical PZT movements; one in the instrument and the other nominally remote. The latter is then used as the reference (i.e. its measured movement is taken as correct). The assumption is made that the equipment is doing the same—a sort of open loop bridge!

The results illustrate exactly what happens in practical instrumentation. Lateral structure is much better defined in relative terms than height parameters.

One of the problems of calibrating height and spacing is the tendency to use two or more techniques. One is the stylus and optical [39, 40]. Another is using AFM and optical interferometers [41], as well as using STM and X-ray interferometers [42]. An optical method is usually included because it gives a clear path to the standard optical wavelength, and hence traceability.

Some recent results have achieved tremendous accuracy in the length measurement in an absolute sense. For example, by using a dye laser, an accuracy of ± 1.3 nm and resolution of 0.2 nm has been achieved over a range of 100 mm. This constitutes a range to resolution ratio of $5 \times 10^8$ which is exceptional [43] and is the same as the wavelength accuracy! The seemingly impossible high ratio has been obtained by measuring the beat frequency between the frequency of the dye laser and that of a stabilized He-Ne laser, which counts as a primary standard of length.

## 8.9 Software correction

One standard way to achieve high accuracy in an instrument or a machine tool is to remove systematic errors usually by software. The idea is to find the systematic errors, store them and then offset subsequent measurements by the value of the systematic error at that point: it being assumed that the position of the systematic error is known. On this basis the absolute accuracy becomes, in effect, the precision of movement. This is usually perhaps five times better than the uncorrected measurements.

In terms of metrology stages this technique has been implemented in two and three dimensions [44, 45]. A metrology stage is self calibrated by an artefact plate, the positions of which are not precisely known. By assuming rigidity of the artefact plate an algorithm extracts the stage error map from comparison of three different measurement views of the plate. When there is no random error the algorithm exactly calibrates the stage error. In the presence of random noise, the algorithm introduces an error of about the same size as the noise itself.

The artefact plate in this analysis is positioned basically (a) at an arbitrary origin (b) at a given, rotation and (c) at a given translation. A considerable calculation has to be used to compute the *actual* values of the stage but this is really not an issue: the calculations can be rapid and accurate. This application of software demonstrates very clearly that a great increase in accuracy can be obtained even if the integrity of the measured data is not good. Software error compensation is used in electron beam lithography, and extensively in coordinate measuring machines.

The messy compensation computations required in this type of exercise show how difficult it is to go from a one dimensional situation to a two dimensional application. Three dimensional equivalents are even more difficult but are routinely being implemented.

Some attempts at 3D scanning have been reported [46]. In this work, a lead zirconium titanate (PZT) is used as a tube scanner measuring the three dimensions' positions interferometrically. Using this technique and using the tube in eight segments, angular deviations which normally limit the accurate use of PZT can be catered for. Sub-nanometre linearity of movement of the tube is claimed.

In practice it is always possible to get a calibration of both the test piece and the artefact by making measurements when the two have various configurations relative to each other. The inversion method, used in linear and rotational instruments is the simplest. This is shown schematically in figure 8.35. If the errors in the test piece are $e[x]$ and those in the master $m[x]$ then the output is $0_1[x]$ where

$$\left. \begin{array}{ll} \text{In case (a)} & 0_1(x) = m(x) + e(x) \\ \text{In case (b)} & 0_2(x) = m(x) - e(x) \end{array} \right\}. \tag{8.13}$$

The method does require that $m(x)$ are stable and that the $(x)$ values are the same in both equations.

This type of accuracy improvements implies that the instrument or master is stable. This in itself implies a certain minimum quality of instrument. Inversion or any other compensatory device can only be carried out with good instruments. This obvious fact is sometimes forgotten.



**Figure 8.35** Error reversal method of calibrating straightness.

The other use of software is not so much in correcting instruments with systematic error as reported above but in the interpretation of the data obtained. This aspect of the use of scanning microscopes is becoming very important especially in biological applications where software interpretation involves the use of high performance computer graphic work stations.

## 8.10   Nanometre metrology systems

### 8.10.1   *Scanning force microscope—nanometrology version*

The term metrology in this and similar devices refers to the possibility of pinpointing the position of the probe or tool in all three axes, and making this position traceable to the international length standard.

Many STM and AFM have produced realistic pictures of surfaces with very fine resolutions. However, as industrial uses increase a more quantitative output is required: the required output is for high accuracy rather than just high resolution. The idea is that if quantitative information from all axes is possible then the microscope becomes a measuring instrument at the atomic level.

Most present day atomic force microscopes use piezoelectric actuators (PZT) for both *x, y,* and *z* profiles. Commonly, translators are in the form of either a tube or tripod for small scans (see earlier).

A major problem with PZT elements is their inherent non-linearity caused by creep, tilt etc. The result is that the voltage applied to the PZT does not translate into position readily. In fact the picture can be highly distorted, even approaching 50% distortion. This has to be reduced. Sometimes it can be achieved by means of master grids and/or refined software.

In some cases lattice spacings of fractions of nanometres can be used. Use of FFT routines on the raw data has been tried [47] but only for relatively short ranges of 0-10 nm.

One way of linearizing is via a compensating non-linear drive voltage, which has to have some knowledge of the distortion, in order that software can be matched to the problem. Also a charge drive can be used [48]. These methods tend to be open loop in the sense that each axis is separately linearized. One of the best ways is to use an integrated metrology loop that uses the co-ordinates of the probe directly. Trying to fabricate nanostructures without having the necessary metrology, to define the position of the probe/tool is doomed to failure.

There are many instruments having position capability. As a rule they rely at least in one direction on a laser interferometer. It is somewhat disappointing to note that the improvement of position measurement, either within the metrology, or outside it, is almost always achieved by refinement of existing laser technology rather than new technology.

In one method [49] the positions *x, y, z*, are obtained using capacitance gauging. In this metrology (AFM) clever use is made of kinematic principles and the use of elastic compliance in various parts of the instrument. PZT tubes are deliberately left out of the design. It is interesting to note that thermal effects are a major factor in determining the performance limit—even if the instrument is kept constant to $\pm 0.2^0$ C.

One very ambitious project due to Teague [50] is perhaps an ultimate, metrology-based (CMM), called the molecular measuring or $M^3$. The specification for this is to obtain a point-to-point resolution of 0.1 nm or the distance between any two points within 50 mm x 50 mm x $100\mu$m volume with a net uncertainty of 1 nm. This is more of a planar measuring device aimed at masks, thin films and wafer characteristics than it is a mini CMM.

The real criteria which have to be met in order to get atomic capability are:

(i) The ability of the material to be stable in the presence of its thermal conductivity and coefficient of expansion.
(ii) The ability of the material to absorb thermal shock. (This is determined by the diffusivity—the ratio of the thermal conductivity and the density multiplied by the specific heat.)

(iii) Satisfactory response of the structure to dynamic and static loading. (This effect is minimized by the material having large specific stiffness.)

(iv) Response to transient loading—the device has to have controlled damping.

(v) Deformation of the material should be elastic.

Long range motion of the probe relative to the specimen is obtained by mounting the probe assembly and the specimen on separate carriages, which move on crossed linear hinges. These are monolithically part of the core assembly.

There are many clever features incorporated into the design to reduce mechanical and acoustic vibration and thermal effects—for instance, the body of the instrument is a hollow sphere of copper 350 mm in diameter.

The metrology reference system is the core of the instrument figure 8.36; (the metrology box).

The datum for displacement measurements is shown in figure 8.36; also shown is the integration with the interferometers for the two major axes. The inside differential interferometer configuration combined with the metrology box enables measurement of the probe position relative to the specimen mounted on the box independently of spurious drift or strain. The external differential interferometer design is shown in figure 8.37.

The outstanding feature of $M^3$ instrument is the totally integrated approach of the design. Because of this integration the desired specification looks achievable.



**Figure 8.36** MMM metrology reference system.

A less ambitious but more practical scheme has been developed at the NPL in Teddington [24]. This comprises a commercial AFM head but, instead of relying on piezoelectric actuators for the stage movement, three separate interferometers have been used to monitor the position of the stage in space. The reason for using laser interferometers (in this case plane mirror Jamin type [51]) is to enable direct traceability of the lasers' frequency against the primary laser standards at the NPL, thereby having a very short calibration chain. Mirror blocks are attached to the $z$ tube carrying the probe and the sample stage. In effect, the three lasers provide an independent metrology reference frame for the instrument. It has been called 'a metrological AFM' because of the independent axes, which do not rely on calibration using images of calibration transfer standards or any other indirect method. Its performance relative to $M^3$ is not known.

### 8.10.2 Traceability

There is no great problem in making the positions along the $x$ and $y$ axes of the metrology instrument traceable to international standards because both axes are of linear position. Instruments have been and

**Figure 8.37** MMM external differential interferometer configuration.

are being built to take advantage of this capability. It is, however, not as straightforward to make the *z* axis traceable.

If the *z* axis is for measuring topography of a surface then traceability is readily possible. The *z* movement of the probe will be measuring height and there should be little difficulty in relating surface height at *x, y* with the vertical position of the probe at *xy*. It is true that there may be some small discrepancy if the probe touches the surface with a force sufficient to cause elastic deformation but this will be ignored in what follows.

In principle it should always be possible to make the vertical position of the probe traceable by using interferometers or other calibrated transducers. This vertical (as well as *x* and *y*) traceability confers the title metrology to the instrument as for example the $M^3$. The question is, to what extent is this *z* position a true measure of the variable? Has the *z* position of the probe been calibrated absolutely with respect to the variable, be it charge density in STM or a specific force in the AFM? Should it be calibrated absolutely or is a closed loop mapping of constant current (or force or whatever) sufficient? In open loop mode SPM the variations in current may be related to the *z* gap, e.g. exponentially, but the assumption about this relationship has to be made; charge variations (or whatever) are inferred and indirect unless a direct calibration has been made. It may be that this aspect of the traceability has been addressed and perhaps even discounted but information is worryingly sparse.

## 8.11 Methods of measuring length and surfaces to nanoscale results with interferometers and other devices

### 8.11.1 *Conventional*

The use of the Michelson, Linnik and Tolansky interferometers for measuring length and surface geometry are well known and described adequately elsewhere [52]. The heterodyne methods are evaluated in chapter 4.

### 8.11.2 *Heterodyne*

As the name implies, there are two independent factors that influence the optical path of such systems. These factors can be different polarizations or different frequencies. Sommargren [53] uses a method similar to Nomarsky but more sophisticated. The basic unit comprise a Wollaston prism, a microscope objective and orthogonally polarized beams with a frequency difference of 2MHz produced by the Zeeman effect. This splitting method is a classical example of the heterodyne principle, in which the combination of two slightly different signals produces a beat of frequency much lower than the individual signals.

As shown in figure 8.38 the signals refract slightly on either side of the normal. The objective lens produces two focus points about 100 $\mu$m apart and 2 $\mu$m in size.

If $\Delta z$ is the optical path difference $\Delta \varphi = \dfrac{2\pi}{\lambda} \Delta z$ where $\Delta \varphi$ is phase change and $\Delta z = \dfrac{4 \Delta h}{(NA^2 + 4)^{\frac{1}{2}}}$ where $\Delta h$ is the height difference on the surface of the two points.

From these

$$\Delta h = \frac{1}{8\pi}(NA^2 + 4)^{\frac{1}{2}} \Delta \varphi \cdot \qquad (8.14)$$

Various spurious signals have to be removed, such as tilt, but the output is a circular profile track of the one spot around the other.



**Figure 8.38** Sommargren heterodyne probe.



**Figure 8.39**

## 8.11.3 *Frequency tracking—Fabry-Perot etalons (FPE)*

This technique has been used to determine the long term stability of materials [54] but is now being used for distance measurement [55].

In figure 8.39, which is a schematic outline of frequency tracking Fabry-Perot etalon, in which light from a tunable laser is transmitted through the etalon. The frequency transmission peaks given by $m\lambda_o = 2qL$

cos$\theta$ where $\lambda_o$ is the wavelength in vacuum. $L$ is the effective length separation between the test and reference surface and $q$ is the refractive index of the media in the optical path of the interferometer. The locked-on frequency is compared with a reference laser. The frequency difference is counted. This is related directly with cavity length.

The change in frequency of a given transmission peak due to a length change $\delta L$ is given by [56].

$$\frac{\delta f}{FSR} = \frac{\delta L}{\lambda_0 / 2\mu} \tag{8.15}$$

where FSR is the free spectral range of the Fabry-Perot etalon and is the frequency difference between adjacent modes of the cavity = $c/(2, \mu L)$ for parallel plane mirror interferometer where $c$ is the velocity of light.

Although attractive, this method can usually only give a range to resolution ratio in the thousands due to practical limitations such as angular sensitivity [57].

### 8.11.4 Capacitive methods

Early on in surface metrology capacitance methods were used to measure roughness. The big advantage over stylus methods was the extreme sensitivity of the method coupled with the fact that an area assessment was possible. The method faltered due to practical considerations usually associated with the difference in tilt and/or shape of the reference electrode relative to the test surface.

Capacitance methods have been resurrected not for surface measurement but for distance measurement. See for example reference [58].

Capacitative micrometers are essentially made of two conducting electrodes separated by a millimetre or less. A typical capacitance is 2pF with a 10 mm diameter. Any change in the overlap area or the separation changes the capacitance. The use of three terminal capacitors and shaped electrodes gives improved performance [59, 60]. However, shaped capacitors can be a disadvantage if the capacitative micrometer is intended for general use because the shape may not correspond to the object.

These sensors are most useful for measuring small gaps to high accuracy, e.g. 100 $\mu$m. This gap can be increased or decreased as the application warrants, as can the bandwidth, usually from 50 Hz to 5 KHz.

Linearity for capacitative sensors is 0.01%. However, by utilizing a small part of the range of a large gap sensor a linearity of 0.005% is possible.

Probably the most comprehensive treatment of capacitative transducers is due to Queensgate instruments [17]. Working from the basic formula for capacitance $C$

$$C = \frac{\varepsilon_r \varepsilon_o A}{d} \tag{8.16}$$

between two electrodes area $A$ and separation $d$ where $\varepsilon_r$ is the relative permittivity of the medium between the plates, $\varepsilon_o$ is the permittivity in a vacuum. Obviously, changing either $A$ or $d$ changes the value of $C$.

However, changing $d$ is much more sensitive than $A$

$$\text{i.e.} \quad \left. \begin{aligned} \delta C &= \frac{\varepsilon_r \varepsilon_o}{d} \varepsilon A \\ \delta C &= -\frac{\varepsilon_r \varepsilon_o}{d^2} \varepsilon A \end{aligned} \right\}. \tag{8.17}$$

It makes sense therefore to use $\delta d$ as the mechanism for small range high sensitivity and to use $\delta A$ for long range low sensitivity.

In figure 8.40 the basic configurations are shown. Figures (a) and (c) are for linear micrometer and (b) and (d) for rotational micrometers. Using standard values for $\varepsilon_r \varepsilon_o$ and making the actual capacitance the same for all configurations, Hicks and Atherton [17] conclude that the sensitivity of the short range variant is about

(*a*) Linear micrometer          (*b*) Rotational micrometer

(*c*) Linear micrometer          (*d*) Rotational micrometer

**Figure 8.40** Configurations for capacitance micrometer.

$100 \times$ that of the long range version, adding the proviso that the range is proportionately shorter. They give the range of the long range rotation as 1.5 rad with sensitivity 6.9 fF mrad[-1] and the short range rotation for comparison as 17 mrad range and 580 fF $-\mu m^{-1}$.

The fact that the values of $C$ are very small (i.e. 10 pF) means that very careful design is needed to avoid edge effects and stray pick-up. Amplifier noise has to be minimized.

What seems incredible is that it is possible to get a worthwhile signal to noise ratio with the value of capacitance so low. The probable reason is that nowadays potential users of nanopositioning devices discuss their needs with the designer early on in the development stage so that it is possible to match the device to the required specification. The reason why capacitative gauges are so popular is that they are relatively simple even when factors such as non-linearity have been corrected.

Stability of capacitative sensors is good because of their simple construction. The sensors and controllers are stable to 50 nm per month or 10 nm for days.

Modern sensors have some problems with tilt of test surface to reference. For example a tilt of one milliradian can cause a change of 0.5% in the scaling factor. Long range sensors, e.g. 500 $\mu$m, are much less sensitive to tilt than the short range. A typical table of performance is given in table 8.4.

**Table 8.4**

| Range ($\mu$m) | Linearity % | Scale Factor ($\mu$m per Volt) | Noise (nm rms $\sqrt{Hz}$) | Thermal Drift (nm °C) |
|---|---|---|---|---|
| 20 | < 0.08 | 2 | < 0.001 | 1 |
| 100 | < 0.08 | 10 | < 0.005 | 4.4 |
| 250 | < 0.06 | 25 | < 0.013 | 11 |
| 500 | < 0.05 | 50 | < 0.075 | 22 |
| 1250 | 0.06 | 125 | < 0.188 | 55 |

Errors due to pick-up noise from external sources, either vibrational or electromagnetic can usually be minimized for different applications.

### 8.11.5 *Stylus instruments*

There are a number of instruments working at the sub-nanometer level of roughness and step heights which follow on from the pioneering work of Dr. Reason of Taylor Hobson with the Talystep in 1966. One deriva-

tion is the Nanosurf I and II being developed at the NPL. These are conspicuous by their long, accurate traverse units of 50 mm. The essential development feature is the kinematic design of slideway in which the inertial forces and frictional resistance forces are balanced through the centre of mass of the unit [61]. In the Nanostep [62] electronic and mechanical noise are reduced to 0.03nm RMS. As in the $M^3$ machine the vertical range is limited to about $20\mu$m.

Some instruments such as the Veritekt 3 by Zeiss [63] can operate as hybrid instruments, in which a sharp stylus of nominally $0.2 \rightarrow 2$ $\mu$m can be used in a conventional mechanical tactile probe or in AFM mode according to requirements. Also such instruments operate in either continuous mode or stepping mode, in which the stylus force is changed from $10^{-8}$ N to $10^{-7}$ for stepping.

In all modem instruments there has been a utilization of low expansion ceramics and glasses such as Zerodur, which considerably reduce thermal problems. In the earlier instruments the mechanical loop had to allow expansion due to heat but, by a clever use of materials such as ebonite, the loop was made stable but dynamic rather than static.

### 8.11.6   *Electron and x-ray diffraction*

X-ray diffraction used as the basis for height or length calibration has been successful. However, the specimen, e.g. the silicon crystal and the diameter of the x-ray beam, are millimetres in size which severely restricts the applications and usefulness of the technique. Using intense synchotron sources brings down the size to about 1 $\mu$m but diffraction patterns of the object still require long durations, e.g. tens of minutes. The only way to investigate small size objects is by using electron beams instead of X-rays. Electrons are used in transmission electron microscopes and scanning electron microscopes but they suffer sometimes from the strong scatter produced in the bulk. Surface scatter gets lost at normal incidence. Also surface resolutions can be poor.

Probably the ultimate use of electron diffraction methods as far as nanosurface metrology is concerned is lateral crystal structure of periodic nature such as the surface of silicon. So, to extend the effective range and reducing the penetration into the bulk, the reflection mode is used. In this the angle of incidence is less than the Bragg angle. This has to be for the dominant lattice spacing on the surface (or suspected to be on the surface).

This configuration produces diffraction spots perpendicular to the surface and at right angles to the incident beam.

Sometimes also obtained are the so-called Kikuch lines, which are the result of incoherent electron scattering. To get an image reconstruction from the diffraction pattern requires all the diffraction spots. What is more usual is to form an image produced by selecting one or more diffraction spots and retransforming this diffraction spot. This does give an image of what on the surface has the periodicity to generate the angle of the diffraction spot. Care must be used in interpreting these images because any randomness of the true surface is blocked; the surface looks convincingly regular.

## 8.12   Summary

Most of the technology outlined above has been concerned with simple measurement such as displacement and position. These and other requirements involve scans of many $\mu$m and often of millimetres. In the vertical direction there is usually a probe with very high sensitivity and a much smaller range than $x$ and $y$. As the requirement for more flexible measurement systems grows the pressure on using repeatable and accurate slideways is very high. This presents problems because slideways fall into the realm more of micro-mechanics than nanotechnology. This meeting of scales of size is always difficult because one scale may be in the domain of physicists, while a larger scale could well be in the domain of engineering: their respective backgrounds are entirely different. There is a similar interface between nanotechnology and atomic scale phenomena. Here the mode of treatment is different because quantum

mechanics is involved in the atomic scale. In other words, nanotechnology is concerned more with the statistics of events (i.e. average distance between planes) while quantum mechanics is concerned more with probabilities. For the generation of measurements based on VLSI technology instruments having wide scans and limited vertical range are absolutely acceptable. The M³ instrument developed by NIST is a very good example of this. However, it is conceivable that the generation of measurements required for micro-motors and MEMS in general, as well as chemical effects such as catalysis will require what can be called a micro-or nano co-ordinate measuring machine in which all axes have nominally the same range and resolution. Any measurement of complex units such as biological cells, rather than the measurement of surface structure may well require such a new instrument. From the notes in this chapter the uncertainties build up tremendously from ID to 2D and also if the 4 requirements of Teague [57] are taken into account.

It could in certain cases involve pathological situations occurring. For example, if an atomically sharp probe were to pass close to another similar probe they could completely ignore each other. Because of quantum mechanics no energy would have had time to be exchanged. In other words, for the finest of measurement temporal bandwidth has to be taken into account — it takes time to build up a signal!

In the case above, the temporal bandwidth would have to be dramatically reduced in order to get a signal. This therefore increases the possibility of noise getting into the system from the environment.

Measurement in the nanotechnology regime is important but whether the new problems posed can be overcome remains to be seen. Even in the areal measurements where expertise already exists there are big problems. One such problem is the need to measure the 400 mm wafers. These are so large yet so thin that roughness, flatness and thickness measurement has yet to be achieved.

The very important issue of calibration and traceability in nanometrology has been raised. That movements in the $x, y$ and $z$ directions of STM and other members of the SPM family can be calibrated with traceable certainty to the international length standard is not in question. What is less certain is whether true fidelity between the actual surface variable of interest and the $z$ movement (which it purports to stimulate) has actually been achieved.

One of the most important points to emerge in nanosurface metrology is that the surface and interface technology is taking on a big role for the future. This is not necessarily in quantum effects such as tunnelling through shaped barrier walls. The main advance is in the biological and chemical development of nano-machines and computers where the surface dominates activity and reactance and where the influence of external stimuli act; bulk properties are somewhat less important. Workers with colloids, catalysis, thin films and boundary lubrication already know this.

This chapter has considered the way in which the traditional geometric features of size, shape and roughness found on typical engineering workpieces are affected by the scale of size. Under the general heading of shape, form and waviness (which embodies machine tool problems) are not likely to be important at small scales where the machining or the process is more likely to be additive, rather than subtractive as in general machining. Additionally, high stresses causing sub-surface problems are unlikely.

Figure 8.41 is an attempt, purely figuratively, to bring together some of the most important threads. One important point is that at scales down to micromechanics, which include microdynamics and MEMS, traditional meanings can be allotted to shape and roughness. The objective is to remove material to get down to size and shape, incidentally but inevitably incurring roughness (process marks). However, at smaller scales the size and shape are inevitable because they are dictated by the atomic/molecular/cluster structure. So size and shape tend to be automatic once the chemical structure has been decided on. What is in question is any deviation or defect. At atomic-size even shape loses its meaning because quantum effects are concerned more with probabilities rather than geometries. Also, at these very small sizes, curves tend to look straight; shapes are effectively linearized at the nanoscale.

One fundamental problem is that at the molecular/biological level it is more likely to be surface reaction which is more important, e.g. in catalysis than surface geometry. There are many other surface possibilities

such as charge densities, measured by the scanning tunnelling microscope, Maxwell strain and local magnetic or capacitative effects.

The various types of scanning probe microscopes which are now in use have basically the same design, which is included in chapter 4 on instrumentation. Here issues which are concerned with how engineering factors can be transcribed down into the nanometric region are discussed.

In figure 8.41 the reason why size, shape and roughness are thin lines at the top is to reflect the point that in normal engineering A1 and A2, the features can be considered independent yet the ratios of size between them does not change dramatically. On the other hand the ratios can become close to unity for MEMS where there are miniature parts being made conventionally.

Similar remarks hold for semiconductor-based surfaces. So on the graph there is a broadening of the size and shape lines. For different applications the lines depicting the geometric features can be moved horizontally.

From the graph it is seen that, over the basic engineering discipline A1 and A2, the geometric feature which falls mainly within the nanotechnology range of 0.1 $\mu$m to 0.1 nm is the roughness! At the smaller scales of size the shape and size move progressively into the nanoband displacing the less noticeable roughness.

An important point is that where traditional engineering mechanics is encroaching on semi conductor technology (i.e. where A3 and B in the figure meet) there is a definite sudden reduction in the average values of size, shape and roughness. This reduction is not uniform because it is a bigger reduction in size than it is in shape and roughness. As a result of this disproportionate reduction a 'bunching' of the geometric features occurs at A3—B domains. This means that size, shape and roughness should be measured with one instrument (i.e. surface and dimensional metrology) converge. However, it is desirable to retain their distinction within the measurement range which is difficult considering their closeness; stiffnesses of conventional instruments would probably be too low to guarantee the necessary resolution.

There is a solution because conventional instruments when reduced in size increase their stiffness even with the same design. Using the Rayleigh approach for stiffness shows that compliance reduces as the cube of the force loop length. An instrument of half size would have a nominal stiffness increase of eight to one—more than adequate for resolving the lumped geometrical parameters: working with smaller scale size has its advantages!

At the next interface—between B and C on figure 8.41 chemical and biological effects take precedence over engineering. Hence, the interpretation of the figure changes below the B regime.

At this interface size and shape are determined by chemical structure and not by machining! This results in the following trends:

a) Size and shape become discrete attributes not continuous variables below B. In effect, at this point the shape and the size become 'quantized', an example of what could be called quantum engineering.

b) Roughness has no meaning at this scale; molecules and atoms and clusters have no roughness in the traditional sense. What they do have are defects and unequal distributions and positions of the components such as nanotubes.

c) Shape can take on another meaning. For example perfect shape would not necessarily be a circle or straight line as in engineering. Perfect shape could be a hexagon thereby recognizing the importance of the benzene ring in organic chemistry and biology.

d) At smaller scales the combinations of different size, shape and roughness becomes great which results in a broadening of the geometric lines in region B.

Summarizing, where the top-down disciplines (engineering) meet the bottom-up disciplines (such as physics), there is a metamorphosis; engineering nanotechnology (Taniguchi) cannot be extrapolated below the semiconductor (B) line in a straightforward manner.

**Figure 8.41** Relation of features to scale and discipline.

# References

[1]   Taniguchi N 1983 *Ann. CIRP* **32** 573

[2]   Bryan J B 1979 *Precision Engineering*, **1** 129–132

[3]   Franks A 1989 Nanotechnology opportunities *J. Phys. E Sci. Inst.* **20** 237

[4]   Bifano T, Blake P, Dow T and Scattergood R O 1987 *Precision Machining of Ceramic Materials* Proc. Ceramic Soc. ASME (Pittsburgh PA: AES) **99**

[5]   Miyashita M and Yoshioka J N 1982 *Soc. Prec. Eng.* **16** 43

[6]   Millar R E and Sheney V B 2000 Size dependence of elastic properties of nanosized structural elements. *Nanotechnology* **11** No. 3 139–147

[7]   Ishikawa T, Lida A and Matsuchito T 1986 *NUCL, Instrum. Methods Phys. Res.* **A256** 348

[8]   Che G, Lakshumi B B, Fisher E R and Martin C R 1998 *Nature* **393** 346

[9]   Zangwill A 1988 *Physics at Surfaces* (New York: Cambridge University Press)

[10]  Duke C B 1996 *Chem. Rev.* **96** 1237

[11]  Uchihara T, Matsumara M, Yamamoto A and Tsubomora H 1989 *J. Phys. Chem.* **93** 5870

[12]  Bayburt T, Carlson J, Godfrey B, Shank–Retxlaff M and Sliger S G 2000 Nanostructure, behaviour and manipulation of nanoscale biological assemblies *Handbook of Nanostructured Materials and Nanotechnology* Ed. H S Nalwa (San Diego: Academic Press) **5** Ch. 12 537

[13]  Sugano S 1991 *Microcluster physics Springer Series in Material Science* No. 20 (Berlin: Springer Verlag) **3**

[14]  Nalwa H S Ed. 2000 *Handbook of Nanostructured Materials and Nanotechnology* (San Diego: Academic Press)

[15]  ISO Handbook No. 33 *Limits Fits and Surface Properties* ISO Geneva 1996

[16]  WhitehouseD J 1994 *Handbook of Surface Metrology* 1st Edn (Bristol: Institute of Physics Publishing)

[17]  Hicks T R and Atherton P D 2000 *Nano positioning* [London: Peuton Press]

[18]  Whitehouse D J 1991 Nanotechnology instrumentation *J. Meas. Control* **24** 37

[19a]  Teague E C 1989 *J. Vac. Sci. Technol*. B7 1898–1902

[19b]  Whitehouse D J 1993 Horns of metrology *Quality Today* November 22–23

[20]   WattsR A, Sambles J R,Hutley M C , Preist T W and Lawrence C R 1997 *Nanotechnology* **18** 35–39
[21]   Musselman I H, Peterson P A and Russell P E 1990 *Precision Engineering* **12** no. 1, 3–6
[22]   Stevens R M D, Frederick N A,  Smith B, Arse D E, Stucty G D and Hansma P K 2000 *Carbon Nanotubes as Probes for Atomic Force Microscopy* (see Ref. 13)
[23]   Tan W and Kepelman R *Nanoscopic Optical Sensors and Probes* (see Ref. 15) **4**
[24]   Leach R, Haycocks J, Jackson K, Lewis A, Oldfield S and Yacoot A 2001 Advances in traceable nanometrology at the national physical laboratory *Nanotechnology* **12** 1–6
[25]   Brand U and Hillman W 1995 *Precision Engineering* **17** no. 1, 22–33
[26]   Edwards H 1997 *Nanotechnology* **8** 6–9
[27]   Bonse U and Hart M 1965 *Appl. Phys. Lett.* **6** 155–156
[28]   Bowen D K, Chetwynd D G, Schwargenberger D R and Smith S T 1990 *Precision Engineering* **12** no. 3, 165–171
[29]   Chetwynd D G,Krylova  N O and Smith S T 1996 *Nanotechnology* **7** 1–12
[30]   Basile G 2000 *Combined Optical and X–ray Interferometer for High Precision Dimensional Metrology* Proc. Roy. Soc. Lond. A **456** 701–729
[31]   Warmington J M and Bowen D K 1992 *American Cryst. Assoc. Conf.* Denver
[32]   Ensell G, Evans A, Farooqui M, Haycocks J A and Stedman M 1992 *Micromach. Microeng.* **12**, 179–180
[33]   Rast S, Waltinger C, Gysia U and Meyer E 2000 Noise of cantilevers *Nanotechnology* **11** no. 3 169–172
[33]   Vorburger T V, Song J F, Giauque G H W, Reneger T B, Whitenton E P and Croarkin M C 1996 *Precision Engineering* **19** No. 2/3, 157–163
[34]   Thwaite E G 1973 *Messtechnik* **10** 317
[35]   Hartman A W and Fang Sheng Jing 1985 *Precision Engineering* **8** No. 4, 203–211
[36]   Futanii S, Furutami A and Yoshida S 1990 *Nanotechnology* **1** 31–37
[37]   Peterka F and J Vaclik 1992 *Sound and Vib.* **154** 95–115
[38]   Haitjema H 1996 *Precision Engineering* **19** No. 2/3, 98–104
[38]   vander Water and Mopenaar J 2000 *Dynamics of Vibrating* AFM Nanotech. **11** No. 3 192–199
[39]   Watt R A, Sambles J R, Hutley H E, Priest T W and Lawrence C R 1997 *Nanotech.* **8** 6–19
[40]   Fuji T, Imabari K, Kowakatsu H, Watanabe S and Bleuler H 1999 *AFM FOR DIRECT Comparison Measurement of Step Height and Crystalline Lattice Spacing. Nanotech.* **10** no. 4 412–419
[41]   Wang Lin, Kuetgens U, Becker P, Koenders L, Li D and Cao M 1999 Calibration of standards for precision pitch measurement in the nanometre region by combining stm and xray interferometer *Nanotechnology* **10** no. 4 412–411
[42]   Brand U 1995 Comparison of interferometer and stylus tip step height measurement on rough surfaces *Nanotechnology* **6** no. 3 81–86
[43]   OgitaE, Lkezawa K, Isozaki K and Mikuriya K 1995 *Nanotechnology* **6** 148–151
[44]   Ye J, Takac M, Berglund C N, Gowen and Pesse R F, 1997 *Precision Engineering* **20** no. 1, 16–32
[45]   Bulsomo A, 1977 *CMM application and calibration Ann. CIRP 4,* **1** 463–466
[46]   Fujii T, Suzuki M, Yamaguchi M, Kawaguchi R, Yamada H, and Nakayama K, 1995 *Nanotechnology*, **6** 121–126
[47]   Jorenson J 1993 *Private Communication*,
[48]   Sommargren E 1981 *Precision Engineering* **131**
[49]   Xu Y, Smith S T and Atherton P D 1996 *Precision Industry* **19** no. 1
[50]   Teague E C 1989 Molecular machine project, american vacuum society *J. Vac. Sci. Technol.* B **7**, no. 6, 1895–2092
[51]   Downs M J, Birch K P, Cox M G, and Nun J W Verification of a polarization—insensitive optical interferometer system with subnanometre capability *Prec. Eng.* **17** 1–6
[52]   Jenkins F A and White H E 1951 *Fundamentals of Optics* (McGraw–Hill)
[53]   Sommargren G E 1981 *Appl. Optics* **20** 610
[54]   Patterson R 1988 *Laurence Livermore Report* UCRL 53787
[55]   Beker P *et al* 1981 *Phys Rec. Lett.* **6** 155–156
[56]   Yariv A 1971 *Introduction to Optical Electronics,* Holt, (New York: Rhinehart and Winston)
[57]   Teague E C 1991 Nanometrology AIP Conf. Proc. 241 *Scanned Probe Microscopy* (Santa Barbara)
[58]   Queensgate Instruments Data Sheet, 1997
[59]   Clothier W K Bairnsfather H J 1967 *Sci. Inst.* **44** 931
[60]   Thompson A M 1958 *IRE Trans Instrum* I–7 245–253
[61]   Davies S T and Chetwynd D G 1992 *Foundations of Ultra Precision Mechanical Design* (London: Gordon and Breach)
[62]   Taylor Hobson Nanostep Technical Specification
[63]   Carl Zeiss Jena. Technical Specification

# Chapter 9
# Summary and conclusions

## 9.1 General

Surface metrology is developing rapidly. In addition to being used as part of engineering metrology it is taking an important role in the emerging discipline of nanotechnology. However, the requirements are different. Figure 9.1 shows the general picture of surfaces. It has been shown that the surface can be used to control the manufacture as well as help in the function (performance) of a workpiece.



**Figure 9.1**

In engineering the diagram applies mainly to situations in which two surfaces are involved, as in tribology, although single body applications such as optical scatter are also important. The scale of size is usually in the micron or submicron region. Often it is height information of the surface that is used. Also the surfaces are often curved due to rotational requirements.

With miniaturization, however, the application is usually planar, single surface and concerned more with areal (i.e. lateral) 2 dimensional structure. Many nanotechnology applications concerned with computing are under this heading; the manufacturers in these cases are usually derived from the semiconductor industry. The scale of size is nanometre or less.

Yet another developing area is in miniature electromagnetic systems (MEMS) in which conventional parts such as gears, cams and other typical engineering parts are made and used on a very small scale. The surface characteristics and requirements come somewhere in between the other two regimes mentioned above.

Each of these classes of applications have their specific needs yet all rely on surface metrology. This has had to develop with the new requirements.

Every chapter in the book shows that progress has been made, despite the fact that the forces for change have been from different directions. In engineering, the so-called top down approach, the requirement is for faster, more precise measurements and functionally realistic characterization. One of the biggest issues is whether to use stylus methods or optics, whereas for the physicists, chemists and biologists who deal with atomic and molecular scales of size there is a requirement for integration. Often the phenomenon to be measured is not topography but a different phenomenon such as charge density whose variation across the surface reveals structure. In these cases calibration is a problem addressed in chapter 5.

The metrology is being required to cope with completely different aspects of the surface. Some of these will be highlighted below.

## 9.2 Characterization

There are two basic questions; one is whether or not enough information has been obtained; the other is whether the information available or part of it can be condensed suitably to be relevant for a specific application. It is unrealistic to expect one mode of characterization to be useful for every application. Unfortunately in the past there have been cries for 'one number' surfaces. Some of the issues are given below in question/answer format.

1. Should characterization be based on discrete or continuous data?

It could be argued that, as the surface is continuous, this form of data should be used for functional consideration. Usually the data is in discrete (digital) form anyway, for convenience, but it can introduce problems, especially for characterizing areal (3D) data as seen in chapters 2 and 3.

2. What are the current trends?

Random process analysis was important because it enabled functional parameters such as peak curvature to be estimated from the profile. Also it could extract tool wear and other manufacturing process parameters from the workpiece data as shown in chapters 6 and 7. It was of more limited use in contact theory. Attempts to enhance its usefulness by using space frequency functions such as the Wigner distribution ambiguity function and wavelet analysis have been partially successful. Complex machine tool vibration problems and flaw detection are two examples.

However, the problem of top down behaviour with two surfaces is elusive (See chapter 7 for new possibilities), because it involves parallel rather than series operations.

Another trend has been to characterize areal character (sometimes called 3D). The argument for this is valid — more detail than profiles is often needed. Some of the attempts at characterization are in chapter 2. However, extension of random process analysis to include areal statistics has proved to be difficult. Instead so-called 'functional parameters' involving, in most cases, some variants of the material ratio curve are being suggested as possibilities. These are mainly height parameters taken with the conventional profile parameters such $R_a$ (extended to areal dimensions to give $S_a$). The minimum number of parameters available for one surface is proposed to be 14. (See chapter 2).

3. What is the situation regarding spatial phase?

Spatial characteristics involving areal 'lay' for machined surfaces and general structure' for nano and semiconductor surfaces are also required. This aspect of characterization when compared with the former shows the difficulty of trying to use the same parameters for different applications. For example, random process analysis in the form of autocorrelation or power spectral analysis *destroys* phase information in the profile (or area) intentionally in order to reveal the underlying statistics in the data.

However, spatial phase provides positional information about the surface, which is precisely what is required for semiconductor and nanometrology application (e.g line spacing or thickness). Phase information

revealed by cross-convolution of the surface system is a new possibility. So what is good for one application is useless for another.

4. What are active parameters?

Because of the difficulty and expense of testing the functionality of surfaces there are moves to use parameters as operators on the data to imitate the function. One example given in chapter 7 is 'cross convolution' to characterize the movement of one surface relative to another. This constitutes an 'active' parameter. Yet another trend is to have 'robust' parameters. These can change their specification according to the data. An example is in plateau honing, in which the low pass filtering simply ignores parts of the surface that have very deep valleys and would distort the remnant signal if included in the assessment.

5. Is fractal analysis useful?

There are conflicting aspects to this subject. It depends on whether the surface characteristic being investigated is scale invariant or not. Much effort is being expended in examining surfaces for fractal properties because fractals get to the very nature of surface generation and, although fractals are a way of looking at the surface spectrum, they do it in a way that highlights the generative mechanism. If fractal properties do exist then it is possible for the usefulness of surface properties to 'abseil' right down the scales of size from engineering metrology to atomic scales, a bonus which eases instrumental problems.

There is no doubt that growth mechanisms and surfaces produced by deposition (Fig. 9.2) could well have fractal properties. However, surfaces produced by abrasion have Poissonian and, hence, Markov properties which, although similar, are not the same. Gaussian lateral properties, so often invoked by investigators, are suspect theoretically. On the other hand Gaussian height characteristics are, of course, highly respectable. The overall picture is shown in figure 9.3.



**Figure 9.2** Fractals in growth mechanisms.



**Figure 9.3** Relationship between processes.

Poisson statistics are a result of random cutting and are symptomatic of Markov processes which have the spectrum characteristic of $(w^2 + \alpha)^{-n}$ rather than—$k/w^n$. The break point given by $\alpha$ gives abrasive information. The index $n$ is integral in the Markov process but not necessarily for fractals. The problem revealed in chapter 2 is that for large $w$ the two mechanisms can look the same.

There is a growing tendency to try to use fractal analysis to describe all characteristics. However, the criterion for relevance is whether the mechanism being investigated is scale invariant or not. For example the wear mechanism or lubrication is scale dependent. Philosophically it is pointless trying to characterize them with fractal analysis which is by definition scale invariant — wear is transparent to fractals!

6. Should waviness be measured separately from roughness?

Waviness is still contentious. There is no doubt that in some applications, for example in light scatter, it is the whole surface geometry that needs to be taken into consideration. In other functions it is not clear whether this simplistic approach is justified. This is because certain geometric deviations (i.e. roughness) are produced by reasonably well-defined causes. These impart other characteristics to the surface as well as the purely geometric. Subsurface effects are produced by high pressures and temperatures and material movement. The same argument applies to waviness, only in this case the conditions are less abusive. To lump them together implies doing so physically as well as geometrically. It is not at all clear that this is allowable. The question arises as to whether there is a 'hidden' subsurface property, which can in fact be identified by the geometrical characteristics of the surface boundaries above. If so, roughness and waviness must be separated. In any case it does no harm; the geometry can always be integrated at a later stage.

This separation is now very straightforward since the introduction of digital filters to surface metrology in 1963 by Whitehouse, and even more so since his phase-corrected filters arrived in 1967. These have since been modified somewhat to allow for a Gaussian characteristic and a transmission at the cut-off of 50% rather than 75%. This latter move is sensible because it enables the assessment of roughness and waviness to be arithmetically complementary. Another method of separation, which takes more account of specific functional or process problems is the 'motif' method. This basically allows the drawing of lines on the profile passing through or at peaks and valleys. The algorithms for generating these lines can be targeted very easily using digital methods.

7. Should characteristics other than topography be included in surface metrology?

This question has already been answered in engineering applications. Greenwood introduced the physical properties of the hardness and the elastic modulus of the surface into the characterization with the 'plasticity index'.

If other mechanical properties are required in the characterization it is natural to try to measure them at the same time as the topography. It is under these conditions that the stylus method of measuring scores over optical methods.

The very basis of contact is challenged in chapter 7. Ways of redefining peaks (or summits) to take into account the bulk of the surface supporting them have been investigated. The premise is that these 'weighted peaks' are more likely to be important in the load carrying capacity of the surface than just peak (or summit) height. This approach leads to cumulative weighted peak distributions instead of material ratio curves, giving different load/deflection characteristics.

From the point of view of the scanning probe microscopes favoured in nanometrology, it is likely that the variable being measured is not topography but could be force or another phenomenon. This puts a different complexion on the definition of surface metrology. It could be that conventional metrology should not be restricted to geometry in the normal direction but should encompass the value of any surface characteristic as a function of the areal dimensions $x$ and $y$.

## 9.3   Data processing

Because a parameter has been identified as being useful by some theory, it does not make it automatically measurable. In practice one very important fact that has emerged is that parameters are not intrinsic to the surface — they are very dependent on the scale of size picked out by the measuring system. This is the reason for the hunt for fractal properties — they are intrinsic.

What has brought this realization about non-uniqueness to a head has been the use of digital methods. There is no doubt that the use of digital techniques has been advantageous on the whole in surface roughness, but there are disadvantages.

There is an almost unlimited capability to measure (or rather, to attempt to measure) many different parameters. The term 'parameter rash' has been used to describe the current situation. The versatile computer allows many parameters to be measured but does not provide the unnecessary physical justification to go with them. The result is a bewildering array of parameters, many of which are questionable.

To make things worse the digital variables themselves, sampling, quantization and numerical mode, greatly influence the numerical values obtained. Take peak measurement, for example. Figure 9.4 illustrates the point vividly. It shows that the probability that an ordinate is a peak changes as a function of the sampling interval (which determines the correlation between ordinates). It seems that very little of the curve relating parameter value to sample interval is free from the effects that can result from the mechanism of digital analysis. Simply adopting an arbitrary sampling interval results in any value for the parameter. It is no good increasing the sampling rate (reducing the sampling interval) to the maximum possible value.

This action increases the risk of errors due to quantization. In the case of peaks, for example, under certain circumstances the count can be considerably reduced. In figure 9.4 the only safe region (B in the figure) is that corresponding to a correlation of about $\rho=0.7$ or thereabouts, a fact which can be found from the correlation function. To illustrate this point, look at the way in which the sampling rate affects the numerical value of the peak curvature (figure. 9.5). It is immediately obvious that, unless the sampling can in some way be tied down, the answer can be anything! Even the criterion that the digital sampling should be as small as possible (region A in figure 9.4) falls down because the finite size of instrument stylus ultimately determines the short-wavelength value and hence the parameter. Unfortunately, letting the stylus decide is not good enough because it can vary from one instrument to the next, purely by chance. Even today styluses are usually quoted only on their bluntest tolerance, not the sharpest, so unless the stylus is known and guaranteed not to have been damaged, it is dangerous to proceed in this manner. Optical methods do not suffer from this latter point — the focus spot cannot be damaged. Also, since the introduction by Williamson of areal mapping methods, questions about areal sampling are beginning to be raised. Conventional sampling grids no longer have to be used. Trigonal sampling shown in figure 9.11 (c) for example, is an optimized sample scheme which covers the area more efficiently and has the bonus that any digital definition —for a summit, for example — the member data points are all of equal weight. It is also better for defining edges. This is even more true for hexagonal grid sampling (figure 9.6 (b)).



**Figure 9.4** Digital analysis and relationship to surface parameters.

**Figure 9.5** Non-intrinsic nature of surface parameters.



**Figure 9.6** Sampling patterns in areal assessment.

Another serious problem of processing is a direct result of measuring the surface, using a linear scan, when high slopes and curvatures are present, such as in aspherics. Removal of the shape of the form can easily give distorted pictures of roughness. The reason for this is that components are normal to the geometrical shape of the surface and not to the direction of traverse. Such problems will be worse when miniature parts are being measured, such as is the case when MEMS are being measured.

The real breakthrough due to digital methods in recent years has been the ability to store masses of data, with the result that the total geometry of 3D workpieces can now be stored. This information can be used to test for functionality. Examples of this are in cylindricity. It is becoming possible now to include the size as well as the shape deviations in the computer. In the past only a few of the surface deviations could be stored. The size of the workpiece was 'suppressed,' not because of choice but simply because the instrument did not have the capability to measure and store all the necessary data. This is no longer true. There is no reason why trial assemblies of workpieces cannot take place on a computer. It may be the best course of action in critical cases. However, there is one big problem. The data obtained from and the algorithms used in CMMs and surface roughness instruments have to be made compatible. An example of the problem is the identification of the centre of the best-fit circle when viewed from inside the circle with a surface-measuring instrument and from outside with a CMM. The results appear to be different! In fact they are the same but it is necessary to know the nature and reason for any apparent differences. The two approaches have to be matched!

## 9.4 Instrumentation

There have been advances in instrumentation in the past few years. In conventional instrumentation, many improvements have been made in software so that almost every instrument is user friendly and has constant checks on results for fidelity. Contingency values for inputs are now usual.

Many instruments are made to measure as many features as possible, including roundness and form, as well as waviness. This is necessary to reduce set up and calibration times. Problem surfaces having complex shapes are usually measured with stylus methods because of their insensitivity to surface slopes which can be high. Figures 9.7 and 9.8 show how the components of the workpiece usually associated with surface metrology change with miniaturization. The nature of the signal will be discussed below (figure 9.7).

There has been a change of emphasis in what to measure. Engineers using stylus instruments have been used to talking about roughness in terms of height parameters — thinking vaguely about contact and fit problems. On the other hand physicists, chemists and biologists usually have been more involved with structure revealed by the areal view in a microscope. Height has been secondary. It is not a coincidence that ordinary microscopes are calibrated accurately in the plane of the surface but not normal to it.

There is a realization in both camps that measurement should be in both planes and, as a result, there has been a convergence.

One aspect of this is that, at the molecular level, size, shape and texture all have to be measured. Their traditional independence is no longer valid. Individual features are genuinely three dimensional. The aspect ratio is shown in figure 9.7.



**Figure 9.7** Aspect ratio of instruments.

The result is that the conventional stylus instrument has had to take up lateral measurements and the other types of instrument height measurement. This is shown schematically in figure 9.8, which has been called the 'horns of metrology' because of the shape of the trend.

Optical methods have developed in the past few years with phase detection methods reported in chapter 4. Software improvement has meant that areal as well as profiles can be produced from relatively complicated shapes. Rough surfaces of tens of micrometres can be measured, as well as tenths of nanometres but their performance, apart from near field instruments, is restricted by the laws of physical optics. A major development has been the use of white light interferometry for absolute distance measurement. Stylus methods do not have the same restriction. Probe methods embodying the scanning probe microscopes have swept away many traditional criteria. Probe methods can shrink in all directions, the limit being mechanical rather than physical.

The new probe methods are more than small versions of traditional stylus instruments. One reason is that the probe properties determine the type of signal from the surface: the tip size and shape are not the only factor: non-geometric features such as conductivity or, say, capacitance play a part. This means that the probe is more important than in conventional stylus instruments. There are more variables to control at the most susceptible part of the instrument. Interpretation of SPM pictures is an art. The scan is not so much of a problem, but it does represent a change of emphasis with stylus type instruments, where the horizontal axis has not previously been regarded as vital.

Another consequence of the minute probe is that the signal between the probe and the surface is subject to quantum physics. With tip dimensions being measured in atomic terms there can only be probabilities of electron tunnelling or photon passage. In other words, the signal to noise is determined by temporal averaging (waiting for the signal to build up) rather than the spatial averaging relied upon in conventional instruments to achieve high positional accuracy. This change is shown shadowed in figure 9.8. A block diagram (figure 9.9) shows the transition from deterministic measurement in traditional engineering (where uncertainty is small compared with the size of part) to probabilistic importance in nanometrology below:



**Figure 9.8** Trend in instrumentation.



**Figure 9.9** Situation concerning tools for instrumentation.

The resurgence of the stylus technique in the STM and AFM, albeit in a different mode from the conventional roughness instrument, has also been leading to a reappraisal of their design.

This is necessary to increase speed and hence to reduce environmental and vibration problems.

Significantly, the way to do this is by realizing that optimization usually results from matching. In this case it is the matching of the system equations to a spatial co-ordinate system rather than the usual temporal

one which can improve performance. The technique is described in chapter 4. It shows again that a change in philosophy from the conventional is necessary in order to meet current and future requirements. By using a spatial co-ordinate philosophy to optimize damping, considerable improvements including speed are possible.

With the requirement to measure very small parts which are not planar, raster scans will not always be sufficient (e.g. MEMS). Planar parts are based on semiconductor applications and are suitable for raster scans. It is highly likely that the attitude of the probe relative to the object and the position of the reference will become much more versatile than at present. The probe should be normal to the surface during all the measurement scan. Another requirement is that $x$, $y$, and $z$ measurement should have equal or near-equal range and resolution. Current thinking suggests that this will be best achieved with PKM designs (parallel kinematic machines) rather than by using slideways.

On a more mundane level, conventional instruments have been tending to compartmentalize. This is shown in figure 9.10. The ordinate of the figure is speed of response and the abscissa is the range-to-resolution ratio. Optical methods in the form of light scatter devices have tended to be used for in-process measurement, whereas stylus methods have been used more for integrated measurement, utilizing the very wide range-to-resolution ratio available in modern instruments. This breakdown is now being questioned with the introduction of fast and accurate optical probes, which have the advantage of being non-contact. All techniques are aiming at the vectorial sum of the two trends, shown as a broken line in the figure. However, this is proving to be very difficult in practice. The basic reason has been the very large range of size encountered in surface metrology, ranging from millimetres down to atomic and angstrom levels. No one length unit,



**Figure 9.10** Trend in optical versus stylus methods.

such as the wavelength of light, can be used to cover such a range, although some very commendable attempts have been made to extend the optical range downwards by using heterodyne methods. Because of the dual requirements of total geometry and nano-resolution, a unit capable of doing both is required. The wavelength of light is still acceptable for the larger sizes and typical resolutions, but not for the atomic resolutions. This is resulting in the use of two units, the wavelength of light for the larger dimensions and possible atomic lattice spacing for the smaller. It is difficult to see how both can be fitted into a CMM capability. Some aspects of this are considered in the next section.

## 9.5 Calibration

The subject of calibration and standards is fundamental to the dissemination of quantitative information, yet it is often neglected, mainly because it does not directly add to the value of the workpiece. Long-term credibility, however, hinges on it. Ensuring good calibration is therefore essential. The general trend in calibration methods is shown in figure 9.11.



**Figure 9.11** Tools of calibration.

The large-range problem mentioned in section 9.4 has been to adjust the calibration techniques to the scale of size, so instead of trying to use the wavelength of light to measure to atomic levels — a jump of more than $1:10^4$ — atomic-scale calibration is used. X-rays are now being employed, thereby enabling the lattice spacings of elements such as silicon to be used as a reference unit.

Another trend in calibration, especially at the nanoscale extreme, is to use more statistical methods rather than the purely deterministic ones used in general engineering practice. This follows the general trend of instrumentation. It is also following the increasing use of quantum effects.

One of the difficulties of calibrating at the molecular and atomic level is that it is both height and spatial characteristics that need to be calibrated, preferably at the same time and under the same conditions. Ideally this would be done with a specimen such as a molecule of DNA, for example, in which all physical dimensions are known. However, this poses problems not just about being sure that the particular molecule is absolutely known, but also about whether it can be put down and anchored without distorting it in any way. So far these problems have not been adequately solved.

There is another issue which is not obvious. This is concerned with the calibration of the $z$ direction in SPMS. If the probe is measuring topography the multiple path interferometers can be used for $z$ as for the $x$ and $y$ positions. However, if the $z$ variable is not topography how is the movement/position of the probe in the $z$ direction related to the $z$ variable (e.g. charge or force). It is understood that using closed loop control enables a profile of constant $z$ variable to be plotted but what does this value mean? It is exactly the same problem as calibrating surface geometry without knowing the size of the part. Ensuring that each axis is traceable to the international metre is not entirely satisfactory because the $z$ position is only indirectly topographic! chapter 5 takes a look at this. Is the AFM calibrated for *any* source of force in the surface or can it be calibrated more specifically? The problem is partly visual. Convincing pictures are not necessarily accurate pictures, as SEM users will know.

The other basic problem is that the traditional unit of length, the wavelength of light, is very large when compared with some of the dimensions being measured (e.g. the size of atoms). There has therefore been a search for a unit of length which is more compatible with atomic distances.

One technique tried in the 1970s was to calibrate height using single or multiple lattice step produced by the cleavage of either mica or topaz. It may be that this technique will be harnessed with the advent of new materials of which many could be suitable. The use of molecular properties in metrology calibration has been adequately solved by using X-rays and the lattice spacing of silicon as reported in chapter 5.

## 9.6 Manufacture

There has been considerable development in manufacturing processes over the past few years. The implications for surface metrology include the following:

Increasing use of multifunction surface to reduce the running-in period in dynamic products and similar attempts to improve performance has led to the need for measuring profiles generated by more than one process. This is an example of the swing toward matching the manufacturing process to a direct function — designer processes.

Increasing use of energy beam methods for fine finishes, that is the use of atomic-type abrasive units to produce atomic finishes. Again this is a case of matching the size of the manufacturing process unit to the function unit.

The use of new materials, sometimes with porosities, for example, can cause measurement difficulties.

In the same way that the process has been matched to cater for a specific requirement of function, so too is the method of measurement. The material ratio curve has been used in plateau honing to identify the various parts of the profile having functional use. Whether or not this is a good example of a process matched to a functional need is not the point. The basic issue is that it is a realization that the two do need to be linked together to get optimum results.

Perhaps the most significant advancement has been the realization that many of the process and machine tool parameters are embedded in the surface topography and that, used properly, this 'fingerprint' can be used to tremendous effect. Using random processes allows the process and machine tool to be investigated. Figure 9.12 shows this for single-point cutting. For grinding the situation is similar where the unit machining mark left on the surface by the average grain is actually preserved within the correlation function.

Some changes are taking place in manufacture. This applies to familiar processes as well as extending to new ones. For example, there is a move towards dry turning in which all or most of the coolant is missing.



**Figure 9.12** Process 'fingerprint' of single-point cutting

Coatings on tools are now making this possible. This makes in-process monitoring of the worked surface much more possible. Milling and casting are being carried out to unprecedented levels in size accuracy and fine surface finishes.

Existing machining methods are being investigated at the subnanometre scale using the computing technique of molecular dynamics. Chip formation and subsurface effects are being resolved by techniques such as this. One outcome is ductile grinding, which has different characteristics to conventional grinding. This is one situation where nanotechnology is bringing engineering and physics together. At the same time theoretical discussions about the various mechanisms for producing surfaces are common. When and where fractal analysis can be used or whether Markov processes can be used to characterize processes are typical questions being asked.

Considerable advances have been made in unconventional machining. The newer methods such as chemical mechanical polishing (CPM) resort to mixtures of mechanisms in order to get particular surface characteristics. Other atomic/molecular processes such as molecular beam epitaxy are now being used.

At the other extreme of finish, 'structured surfaces' are being used for a variety of applications. These surfaces have periodic surface patterns which are picked for a specific use. Examples range from golf balls, tyres, road signs, etc, to fresnel lenses and mechanical seals. The surfaces are generated by a number of methods ranging from laser texturing to etching.

Applications are growing rapidly for such surfaces. One important point is that the patterns are invariably areal (i.e. on the surface): height variations, except perhaps for aspherics are relatively few. This trend is a worry to instrument makers because the patterns in the surface have high depth to width aspect ratios, which makes them difficult to measure using conventional techniques.

Changes have also been taking place in consideration of where to measure in the manufacturing environment. For many years the emphasis has been to make measurements on a batch basis in the inspection room. Some attempts have been made to measure features in process with feedback, to adjust adaptively the manufacturer's parameters to optimize the texture. Efforts to do this have not been entirely successful. The ideal, but probably misconceived, concept from the point of view of the manufacturer is not to measure the workpiece at all, but to make the process so good in the first place that measurement is unnecessary once the machine has been set up. This is called establishing the 'capability' of the machine tool. Obviously measurement is required in the first instance but not afterwards. To the manufacturer measurement is time consuming in most cases, and in any case involves cost. It has been said that measurement only produces scrap! In fact it only identifies it. It is no wonder that in the short term at least some manufacturers shy away from the commitment to quality. In the long term, of course they cannot do this. The 'capability' philosophy precludes process monitoring using the surface.

The message here is that measurement should not stand alone but it should be as close to the manufacture and function as possible. Only then can the tremendous amount of data on the surface be properly utilized for better manufacture to give optimum rather than just adequate performance of the workpiece.

## 9.7 Surfaces and function

The two problems with linking surface characteristics to function are: the difficulty of getting reliable information and the problem of producing an integrated picture whereby the relationship between surface texture and function is consistent.

Is it possible to find generic relationships between the texture and function rather than the haphazard relationships found in the literature?

Chapter 7 of the book attempts to help in both problems. Wherever possible theoretical as well as practical results have been collated so that a broad picture might emerge. Looking into the theory helps to avoid making false links due to instrument availability rather than true correlation.

Some of the evidence is given below. This is followed by the function map approach given in the main body of the text.

There have been a number of major inputs to functional behaviour over the last two decades or so. Two of these are random process analysis and digital techniques. The former has helped the understanding and the latter the implementation of theoretical models in simulations.

Probably the most significant single contribution is due to Longuet-Higgins, who opened up what was more or less a communication subject — random process analysis — into one which could be translated readily into surfaces and moving surfaces at that. From his pioneering work many of the subsequent advances were built.

Another significant innovation, in contact theory at least, was due to Greenwood and Williamson, who at a very early stage started the trend which is accelerating even today: that of integration. Their idea of combining geometrical and physical properties into one index to predict functional behaviour made qualitative what was intuitively obvious: that the behaviour under load depended on more than just geometry.

The other discovery, which, although helping in modelling and the understanding of measured results, caused a great deal of concern and which today is now accepted as inevitable, is the analysis by Whitehouse, who investigated for the first time the discrete properties of random surfaces. It showed that there were no intrinsic properties of surfaces — such things as the mean peak curvature, for example. All depended greatly on the measurement parameters (sampling, quantization, etc). These issues tied in exactly with the digital measurement of surfaces then being introduced. The realization that there was no intrinsic value followed from this and has thrown considerable doubt on many results in contact and tribological verification.

It could be said that these three steps were definitive in the theory of surfaces: the first in characterization, the second in contact and the third, discrete methods in the measurement theory of surfaces.

Many excellent results stemmed from the work above, notably Nayak, who integrated the work of Greenwood, Williamson and Longuet-Higgins, and McCool who integrated that of Greenwood, Williamson and Nayak.

Significant research has been associated with the introduction of roughness into lubrication models.

As far as boundary lubrication is concerned the early work of Hurst and Hollander still has to be surpassed. Their approach included some very thorough experimentation.

In normal lubrication Christensen first used actual stochastic profiles in the analysis, but unfortunately only for cases where few contacts take place, so he and subsequent investigators along this line had the disadvantage of considering a rather specialized case. Cheng flow factors for a global Reynolds equation from exact local solutions for known areas of roughness seem more plausible. It seems a pity that the two approaches are not closer! Time will tell which is the most fruitful. Dowson and Higginson's work in effectively linking contact phenomena under load and with lubrication in EHD and PHD is still definitive.

One of the worrying trends in investigating the functional behaviour of surface is the enormous increase in the amount of theory and simulation at the expense of experiment. This is inevitable in some ways because experimentation is expensive and time consuming and simulation is getting cheaper every year. The worry is that the work could get further from reality. A balance needs to be kept between theory, simulation and experimentation. The trend seems to have swung in favour of simulation.

Despite all the original promise there do not seem to be any new contact theories, only highly refined versions of old ones. The same can be said for characterization. In fact the situation has now been reached where real surfaces are being mapped and stored on a computer, where the experiment is simulated point by point using such tools as Hertzian fields and finite element analyses to predict flow and compliance without bothering about characterization. The next logical step to this is to carry out the experiment directly and to forget the simulation and theory altogether. This is a worrying thought because it indicates that something in the theory is missing and that bedrock — the experiment — is the only answer to a better understanding.

Other functions hardly fare better. For example, in optics there is still a lot of theory being produced which is marginally useful in practical cases (e.g. explaining shadowing effects), yet the more rigorous methods involving vector theory have still to be able to cope with ordinary surfaces (this role has been left to the advocates of scalar theory because it is realistic and, more to the point, tractable, although less accurate).

The function map concept is an attempt to classify function. This critical step has been rather overlooked because most industrial applications and research projects are concerned with relatively confined

performance or functional requirements. It is impossible to pinpoint the relative importance of surface geometry to function — function is open ended. A possible answer is to characterize function itself. A first attempt has been suggested by Whitehouse. This is the 'function map' approach.

The two axes (a) the gap between surfaces and (b) their relative lateral velocity represent the simplest, but at the same time most comprehensive classification.

Figures 9.13, 9.14 and 9.15 show how many functions can be represented on this map. Figure 9.13 identifies directions of energy flow relative to the surface positions. Figure 9.14 applies force boundaries so that the boundaries between different function mechanisms can be defined. Finally, figure 9.15 shows some critical functional conditions.

Although these are in a dimensionless mode, by normalizing the axes, they place the various functions relative to each other. Obviously, to be of any practical use the axes need to be quantified.

The function map as shown has two bodies. Other possibilities are shown in figure 9.16. By taking the gap to be large, the two surfaces can be considered to be independent and single surface properties, such as semi-conductor applications, can be explored. In chapter 7 many properties of single surfaces such as optical scatter are explored.

Trying to match surface parameters to functional parameters is very difficult. It is possible to give cases but they are generally inconsistent. Figure 9.17 shows one way to put surface parameters on the function



**Figure 9.13** Function map flow regimes.



**Figure 9.14** Force boundary domain.

**Figure 9.15** Function map.



**Figure 9.16** Tribological regime.



**Figure 9.17** General surface relevance.

**Figure 9.18** Surface parameter map.

map. Notice how vague it is! In fact only *types* of parameter are represented; this is about all that can be said with certainty. More precise surface parameters can be put down as in figure 9.18 but these can only be tentative. As there is often confusion about the importance of surface geometry, the influence can be good or bad. Table 9.1 shows a simple breakdown. When referred to figure 9.17 it can be seen that, as the distance from the origin *increases,* the influence of the surface finish becomes more degrading.

**Table 9.1** Dual roles of surface metrology and relationship to the statistics.



chapter 7 ends with an attempt to illustrate how the function map can be put on a theoretical basis. At present this is a general systems approach but it should be pursued with urgency.

## 9.8   Nanometrology

Most of the technology outlined above has been concerned with simple measurement such as displacement and position. These and other requirements involve scans of many $\mu$m and often of millimetres. In the vertical

direction there is usually a probe with very high sensitivity and a much smaller range than *x* and *y*. This is a result of the objects being planar. In the future this will have to change as components become more three-dimensional.

As the requirement for more flexible measurement systems grows, the pressure on using repeatable and accurate slideways is very high. This presents problems because slideways fall into the realm more of micro-mechanics than nanotechnology. This meeting of scales of size is always difficult because one scale may be in the domain of physicists, while a larger scale could well be in the domain of engineering: their respective backgrounds are entirely different. There is a similar interface between nanotechnology and atomic scale phenomena. Here, the mode of treatment is different because quantum mechanics is involved in the atomic scale. In other words, nanotechnology is concerned more with the statistics of events (i.e. average distance between planes) while quantum mechanics is concerned more with probabilities. For the generation of measurements based on VLSI technology instruments having wide scans and limited vertical range are absolutely acceptable. The $M^3$ instrument developed by NIST is a very good example of this. However, it is conceivable that the generation of measurements required for micro-motors and MEMS in general, as well as chemical effects such as catalysis will require what can be called a micro- or nanocoordinate measuring machine, in which all axes have nominally the same range and resolution. Any measurement of complete units such as biological cells rather than the measurement of structure may well require such a new instrument. It is possible that the only way to achieve such a working volume is by means of a Stewart platform type of axis generator (i.e. PKM).

One of the most important points to emerge in nanosurface metrology is that the surface and interface technology is taking on a big role. This is not necessarily in quantum effects such as tunnelling through shaped barrier walls. The main advance is in the biological and chemical development of nano-machines and computers where the surface dominates activity and reactance and where the influence of external stimuli act; bulk properties are somewhat less important. Workers with colloids, catalysis, thin films and boundary lubrication already know this.

Chapter 8 has considered the way in which the traditional geometric features of size, shape and roughness found on typical engineering workpieces are affected by the scale of size. Fractals may be more relevant here but Markov processes seem better able to bridge the gap between micro and nanotechnology.

Figure 9.19 is an attempt, purely figuratively, to bring together some of the most important threads. One important point is that at scales down to micro-mechanics, which include microdynamics and MEMS, traditional meanings can be allotted to shape and roughness. The objective is to remove material to get down to size and shape, incidentally but inevitably incurring roughness (process marks). However, at smaller scales the size and shape are more invariant because they are dictated by the atomic/molecular/cluster structure. So size and shape tend to be automatic once chemical structure had been decided on. What is in question is any deviation or defect. At atomic size, even shape loses its meaning because quantum effects are concerned more with probabilities rather than geometries. Also at these very small sizes curves tend to look straight; shapes are effectively linearized at the nanoscale.

It is important to realize that even the fundamental geometries change meaning as the scale is reduced. Figure 9.19 is an attempt to show such changes as the scale of size reduces. The scale is the horizontal row and the application the column on the left.

The reason why size, shape and roughness are thin lines at the top of the figure is to reflect the point that in normal engineering A1 and A2, the features can be considered independent yet the ratios of size between them do not change dramatically. On the other hand the ratios can become close to unity for MEMS where there are miniature parts being made by conventional processes.

Similar remarks hold for semiconductor — based surfaces. So on the graph there is a broadening of the size and shape lines. For different applications, the lines depicting the geometric features can be moved horizontally.

From the graph it is seen that, over the basic engineering discipline A1 and A2, the geometric feature which falls mainly within the nanotechnology range of 0.1 $\mu$m to 0.1 nm is the roughness! At the smaller

Scale → Metre | Milli | Micro | Nano | Approach | Discipline

| Feature | Size | Shape | Roughness |

Conventional engineering — A1
Precision engineering — A — A2
Micro mechanics MEMS — A3
Semiconductor (Lithography) — B — B
Molecular nano tubes clusters — C1
Biological self assembly — C — C2
Atom quantum statistics — C3

Approach: Top down ↓ (Markov processes) (Fractal processes) Bottom up ↑

Discipline: Engineering, Mechatronics, Materials, Chemistry, Biology, Physics

Nano technology range

**Figure 9.19** Features as functions of scale and applications.

scales of size the shape and size move progressively into the nanoband displacing the less noticeable roughness.

An important point is that where traditional engineering mechanics is encroaching on semiconductor technology (i.e. where A3 and B in the figure meet), there is a definite sudden reduction in size than in shape and roughness. As a result of this disproportionate reduction a 'bunching' of the geometric features occurs at the A3 — B domains. This means that size, shape and roughness should be measured with one instrument (i.e. surface and dimensional metrology converge). However it is desirable to retain their distinction within the measurement range if possible, which is difficult considering their closeness; stiffnesses of conventional instruments would probably be too low to guarantee the necessary resolution.

There is a solution because conventional instruments when reduced in size increase their stiffness even with the same design. Using the Rayleigh approach shows that compliance reduces as the cube of the force loop length. An instrument of half size would have a nominal stiffness increase of eight to one — more than adequate for resolving the lumped geometrical parameters: working with smaller scale size has its advantages!

At the next interface — between B and C on figure 9.19 chemical and biological effects take precedence over engineering. Hence, the interpretation of the figure changes below the B regime.

At this interface size and shape are determined by chemical structure and not by machining! This results in the following trends:

a) Size and shape become discrete attributes not continuous variables below B. In effect, at this point, the shape and the size become 'quantized,' and are an example of what could be called quantum engineering.

b) Roughness has no meaning at this scale; molecules and atoms and clusters have no roughness in the traditional sense. What they do have are defects and unequal distributions and positions of the components such as nanotubes. It may well be that roughness in the traditional sense will be overshadowed by surface patterns in which height is constant and transverse structure matched to a design requirement is more relevant.

c) Shape can take on another meaning. For example perfect shape would not necessarily be a circle or straight line as in engineering. Perfect shape could be the hexagon thereby recognizing the importance of the benzene ring in organic chemistry and biology.

d) At smaller scale, the combinations of different size, shape and roughness becomes great, which results in a broadening of the geometric lines in region B.

Summarizing, where the top-down disciplines (engineering) meet the bottom-up disciplines (such as physics), there is a metamorphosis: engineering nanotechnology (Taniguchi) cannot be extrapolated below the semiconductor (B) line in a straightforward manner.



**Figure 9.20** Manufacturing map $\rightarrow$ function map.



**Figure 9.21** Metrology map $\rightarrow$ function map.

## 9.9 Overview

The following points can be made:

1. Size, shape and texture need to be integrated at the nanometre scale. There is no break on the workpiece.

2. The integrated involvement of geometrical and physical properties will increase, as will characterization based on this trend. There has already been great progress through early work on an index combining the two. This will continue, but probably at a more fundamental level involving atomic phenomena. The ultimate goal will include chemical properties.

3.  There is a growing difference in the philosophy of metrology between precision engineering, nanotechnology and atomic theory; at each of these boundaries bridges need to be built. In any event, common links have to be formed to leapfrog the scale of size jumps without losing the basic physics. The Wigner distribution as a common mathematical tool may be one way of linking quantum and probabilistic theory at the atomic level to random process analysis in general engineering.

4.  The whole basis of metrology instrumentation may have to move towards incorporating statistical enhancement methods for nanotechnology and even probabilistic considerations for atomic-scale measurements where atomic-type resolution is required in all co-ordinate axes.

5.  There has to be an integration of disciplines as the scale of size gets smaller because the atomic scale, usually the preserve of chemists, biologists and physicists, is now being encroached upon by engineers — all using the same atoms!

6.  Metrology is keeping pace with the technology — in fact leading it in some ways (i.e. X-ray metrology). In future metrology and instrumentation will have to get closer to both manufacture and function in order to fulfil the goal of manufacturing for a targeted function.

7.  There has to be a matching (to better than three orders of magnitude of scale to size of the unit of interest) in manufacture, metrology and function. If the function is atomic or molecular then the metrology unit must be within three orders of magnitude from it. Also, the unit of manufacture should be three orders from the atomic (e.g. using ion beam machining).

8.  All dimensions $x, y$ and z will have to be dealt with equal weight, and not be unbalanced as at present. New miniature CMMS with more degrees of freedom are called for.

9.  Surface is getting more important because of the trend to miniaturize. The role of surface metrology therefore becomes even more important as the scale of size is reduced.

10.  Calibration of non-topographic features for SPMS (in the $z$ direction) will have to become more formal.

11.  Theoretical work is needed to fill the gap, in certain functional cases, at present in explaining the dynamic 'top-down' behaviour of surfaces in contact. In particular the surface system (i.e. two or more surfaces) should be included in the surface characterization and not just the one surface. It seems pointless to specify 14 or more parameters for one surface and to ignore the other surface usually present. A possible way of doing this can be found in chapter 7. Also this book has attempted to unify the three operational blocks shown in figure 1.1: manufacture, measurement and function.

The vehicle for this has been the concept of the function map. Figures 9.20, 9.21 and 9.13 show that it is possible to show a framework. Figure 9.20 illustrates that contact effects can reasonably be expected to be due to the process marks whereas the path of the tool is more significant in determining lateral flow. Figure 9.21 shows that areal measurements of texture are needed to determine lateral flow patterns. To a first order, the manufacture, measurement and functional use are related. What is needed, apart from a considerable refinement of the arguments given above, is the inclusion of information science in either this or a similar mapping concept. It could be that the ordinate would become tunnelling gap and the abscissa the switching speed between areal elements.

Another main thrust of the book has been to identify and, wherever appropriate, to include nanotechnology and nanometrology in the general subject of surface metrology. This has meant investigation scale of size effects in metrology instrumentation as well as manufacture and function. Into this has crept arguments on fractal analysis; (i.e. when not to use it rather than when to use it).

Chapter 8 has been devoted to the metrological aspects of nanotechnology but it has been necessary to include nanotechnology in every other chapter of the book.

This handbook is possibly the first attempt to marry the engineering and the physical approach to surface and nanometrology. It has attempted to conjoin both disciplines with a common philosophy based on a systems approach. Throughout the book new and hopefully useful contributions have been put forward rather than presenting many disjointed facts. This stimulus may advance the boundaries of surface and nanometrology, as well as clarifying them.

# Glossary

| | |
|---|---|
| **2 *CR* network** | Analogue high pass filter used in early days to block waviness and form signals and thereby allow roughness to be measured. |
| **A/D converter** | Electronic device to convert analogue signal to a digital one. |
| **AA** | Arithmetic average roughness. Used as roughness measure in USA — equivalent to $R_a$ and CLA. |
| **Abbé error** | Misalignment between the sensitive direction of the metrology instrument and the dimension of the workpiece being measured. |
| **Abbott — Firestone curve** | Material ratio curve. Sometimes also called bearing ratio or bearing area curve. |
| **Abrasive machining** | Machining which uses multiple random grains as the cutting tools. |
| **Abusive machining** | Machining of such severity it causes subsurface stresses, usually tensile, and subsurface damage. |
| ACF | Autocorrelation function. |
| **Adhesion** | Force of attraction at subnanometre distances between bodies caused by atomic and molecular forces. |
| AFM | Atomic force microscope. Measures atomic forces. |
| **Aliasing** | Ambiguity introduced into digital signal due to infrequent sampling which causes a folding over of the spectrum in the frequency domain and gives false low frequencies. |
| **Ambiguity function** | A space frequency function based on the Fourier kernel. Achieves, in effect, a two-dimensional correlation. |
| **Amplitude discrimination** | A way of rejecting small peaks by imposing a minimum height restriction. |
| **Analytic signal** | A signal which only has positive frequencies; achieved by use of the Hilbert transform. |
| **Angular distortion** | Distortion of angles in roundness by measuring through apparent centre of part. |
| **Angular motion** | Angular component of error motion. |
| **Anisotropy** | Degree to which a surface has lay. Deviation from isotropy. |
| APDF | Amplitude probability density function. |
| **Asperity** | Peak. |
| **Areal** | Two-dimensional measurement of surface property over an area. Sometimes called three-dimensional measurement. |
| ARMA, MA | Autoregressive moving average and moving average. Recursive relationships used in time series analysis. |
| **Aspect ratio** | Ratio of vertical to horizontal magnifications used in instrumentation. |
| **Asperity persistence** | The reluctance of asperities to crush under load. Probably due to interaction between asperities. |
| **Assessment length (evaluation length)** | Length of surface over which parameter is assessed. Usually this distance is five sampling lengths. |
| **Auto correlation function** | The expected value of the product of signal with itself displaced. The function is the plot of the expected value for a number of displacements. |
| **Auto covariance function** | As above but with the mean values removed. |
| **Axial magnification** | Line of sight magnification. Axial movements in the object plane are squared in the image plane. |
| **Axial motion** | Axial component of motion error of machine tool. |
| **Axis of rotation** | Actual axis of rotation of spindle of machine tool or measuring device. |
| **Bandwidth** | Effective width of frequency spectrum. Usually measured to half power point. |

| | |
|---|---|
| **Barrelling** | Shape of type of cylindrical error. |
| **Basis function** | Unit machining or measuring function. |
| **Bearing area** | Material ratio value. The term area is a misnomer. |
| **Bearing ratio** | Material ratio. |
| **Beat frequency** | Frequency of difference between two signals. |
| **Best fit** | Usually refers to the least squares criterion or maximum likelihood. |
| **Beta function** | Function with two arguments used to classify probability density curves. |
| **BIFORE transformation** | Binary Fourier transformation. Related to Walsh function. |
| **Bistatic scattering** | Double scattering of ray of light from surfaces. |
| **Born approximation** | Used to determine reflectance coefficients in thin films. |
| **Bosses** | Surface protrusions taken as basis for scatter of light. |
| **Boundary lubrication** | Lubrication regime in which the bodies are separated by a molecularly thin film which reduces friction but does not support load. |
| **Box function** | A function representing a running average of data. |
| **Bragg equations** | Equations relating to the scattering of x-rays from a crystal. Elastic scattering. |
| **Bragg scattering** | X-ray scattering from crystal lattice. |
| **BDRF** | Bidirectional reflectance function. Optical scatter taken over many angles. |
| **Brillouin scattering** | Non-elastic scattering — acoustic mode of vibration. |
| **Brownian movement** | Random movement similar to random movement of molecules. |
| **BUE** | Built-up edge of material left on tool during cutting. |
| **Burnish** | Plastic movement of material by cutting tool, which usually smooths topography. Often produced in diamond turning. |
| **CAD** | Computer aided design. |
| **Calibration chain** | The chain of calibration from international to workshop standards. |
| **Calibration specimens** | Specimens used to calibrate roughness measuring instruments, usually for $R_q$. |
| **Calliper** | The basic system for a metrology instrument. |
| **Capability** | The guaranteed performance achievable by a machine tool or instrument. |
| **Cartesian coordinates** | Coordinate plotted on orthogonal axes. |
| **Caustics** | Reflections off surfaces caused by local curvatures. |
| **CBN** | Cubic boron nitride — cutting tool material. |
| **Cepstrum** | Inverse Fourier transform of the logarithm of the power spectrum. |
| **Characteristic depth** | Depth of profile from motif. |
| **Characteristic function** | The Fourier transform of a probability density function. |
| **Chatter** | Vibration imparted to a workpiece by elements of the machine tool being too flexible. |
| **Chebychev polynomial** | Type of polynomial which fully uses the separation between two boundaries. |
| **Chirp signal** | A signal whose frequency changes as the square of time. |
| **Chordal** | Measurement of roundness taken off chords of workpiece. |
| **Circularity** | Roundness. |
| **CLA** | Centre line average. Old roughness parameter equivalent to $R_a$. |
| **CMM** | Coordinate measuring machine. |
| **Coherence discrimination** | System which allows interference to be detected. |
| **Coherence modulation** | This is the form given to the envelope of a fringe pattern due to the finite bandwidth of the source. |
| **Coherence spatial** | Degree to which two rays from different spatial positions of the source are in phase. |
| **Coherence temporal** | Degree to which a signal is monochromatic. |
| **Compliance** | Elastic deformation of object when under load. |
| **Compression ratio** | Difference in magnification between horizontal and vertical axes. |
| **Compton scattering** | Non-elastic scattering. |
| **Concentricity** | Two $\times$ eccentricity. Locus of the centre of selected figure rotating around selected datum, e.g. axis. |
| **Condition function** | Boundary conditions in exchange algorithms. |

| | |
|---|---|
| **Confidence interval** | Boundaries on probability of a value to within a given significance. |
| **Conic** | Any geometric figure governed by second degree equations. |
| **Conicity** | The departure of an object from a true cone. |
| **Conversion efficiency** | The efficiency of a transducer in converting energy from one form into another. |
| **Convolution** | A type of integration of two functions in which the variable of integration gets folded, sometimes called 'Faltung'. |
| **Correlation** | The expected value of the product of two variables normalized with respect to their standard deviations. |
| **Correlation length** | The lag value over which the autocorrelation function falls to a small value, usually 10% or 1/e. |
| **Creep** | Difference in velocity between different points within contact region of rolling ball. |
| **Crest** | Taken to mean peak. |
| **Critical angle** | The angle below which internal reflection takes place. |
| **Critical distance in ductile grinding** | The depth of cut which allows the mode of grinding to be plastic and not fractural. |
| **Curvature** | Reciprocal of the radius of curvature. Sometimes approximated as the second differential. |
| **Cut-off length** | The wavelength along the surface which corresponds to the sampling length. |
| **Cylindricity** | Departure of workpiece from a truly cylindrical shape. |
| *D* **ratio** | Ratio of surface roughness to film thickness in pitting. |
| **D sight** | Whole sight techniques incorporating reflection from a diffuser, used to show up flaws. |
| **DAF** | Discrete ambiguity function. |
| **Damping factor** | Term representing energy loss in a second-order differential equation. |
| **Defect** | A departure from the expected statistics of a surface. |
| **Degree of freedom** | An independent movement. |
| **Designer surface** | A surface made specifically for a given function. |
| **DFT** | Discrete Fourier transform. |
| **DFTC** | Departure from true circle. |
| **Diametral** | Variations across diameters in roundness measurement. |
| **Difference operators** | Difference between discrete measurements. Can be central, forward, or backward. Used in numerical analysis. |
| **Differential logarithm errors** | The logarithms of the errors in a propagation formula for errors. |
| **Differential spatial damping coefficient** | The damping coefficient expressed in spatial terms. |
| **Diffuse reflection** | Light scattered from a surface at angles other than the specular angle. |
| **Digitization** | Taking discrete values of an analogue waveform usually at equal intervals. |
| **Digonal sampling and analysis** | Three point digital analysis. |
| **Dimensional metrology** | The measurement of the linear dimensions of workpieces. Usually including angles. |
| **Dirac comb** | A train of impulses. Can be used to represent sampling discretely. |
| **Directionality** | A measure of the way in which crests are pointing away from the normal. |
| **Discrete parameter** | A parameter of a surface obtained from the digital values. |
| **Discrete properties** | Properties of the digital form of surface. |
| **DMT** | Model based on Derjaguin to explain adhesion between objects, usually one rough and one smooth. |
| **Dual** | Alternative to primal method in exchange mechanisms. |
| **Ductile grinding** | Grinding in which the mode of material removal is plastic. |
| **Dynamic interaction** | Contact of two bodies, usually under lateral movement. |
| **E system** | System of measurement based on rolling ball across surface and measuring from the locus of the lowest point. |

| | |
|---|---|
| **Eccentricity** | The distance between the centre of workpiece and centre of rotation of measuring instrument. |
| ECM | Electro chemical machining. |
| EDM | Electro discharge machining. |
| EHD | Elastohydro dynamic lubrication. |
| EHL | Elastohydro dynamic lubrication. |
| **Eigen vector** | Mode of solution of a matrix. |
| **Eigen vector elastic scattering** | Scattering of radiation from an object which does not result in a change of wavelength. |
| **Ellipticity** | The maximum ratio of two orthogonal axes describing a near circular part. |
| **Energy gradient in instruments** | The rate at which energy is transferred as a function of displacement. |
| **Engineering metrology** | The overall subject of measurement as viewed from an engineering point of view. |
| **Envelope** | Curves generated according to a set of rules tangentially connecting peaks and/or valleys. Used as reference from which to measure roughness or waviness. |
| **Envelope methods** | Methods based on envelopes; usually the E system. |
| **Epitrochoid** | Shape of the stator of the Wankel engine. |
| **Equal weight techniques** | Digital samples based on equal area rather than equal intervals. |
| **Ergodic** | Statistical situation where a temporal average is equal to a spatial average. |
| **Error separation methods** | A technique in which the errors of a specimen and the instruments reference movement are isolated simultaneously. |
| **Errors** | Deviation from intended shape or size. |
| **Errors of form** | Long wavelength. Geometric deviation from the intended geometric shape. |
| **Expectation** | Average value statistical expectation. |
| **Extrapolation** | Technique to estimate value of a function outside its given range. |
| **f number** | Means of describing lens performance. Focal length divided by lens aperture. |
| *F* **test** | Fisher test for variance. |
| **Face motion** | Component of error motion of machine tool orthogonal to the axis of rotation. |
| **Factorial design** | Specialist form of experimentation designed to isolate the effect of the independent variables. |
| FECO **fringes** | Fringes of equal chromatic order. |
| FFT | Fast Fourier transform. |
| **Fibre optic** | Thin strand of optical transparent material along which information or electromagnetic energy can be transmitted. |
| **Fidelity** | The closeness of a measured signal to the original signal. |
| **Filter cut-off** | Frequency (or wavelength) corresponding to 50% attenuation. In the past this has been 75%. |
| **Fingerprint in manufacture** | A set of surface parameters completely defining the process and machine tool. |
| **Finish machining** | Final machining to achieve desired surface texture and dimension. |
| **Flash temperature** | Temperature of contacting asperities when scuffing occurs. |
| **Flatness** | Departure of workpiece from true flatness. |
| **Flaw** | Deviation from the expected statistics of the workpiece. |
| **Flicker noise** | Electrical noise inversely proportional to frequency. |
| **Flying spot microscope** | A microscope in which a small source of light and a small detector localized to it are synchronized. |
| **Follower** | A measuring system in which the surface geometry is followed by the measuring instrument. |
| **Force loop** | The path along which forces act in a metrology instrument or machine tool. Usually drive forces and inertial forces. |
| **Fourier transform** | Transformation depicting summation of sine waves. |
| **Fractal dimension** | The parameter depicting the scale of size. Given by 1/2(5−n) where n is power law of spectrum. |

| | |
|---|---|
| **Fractal property** | A multiscale property in which the parameters are the same for each scale of size. |
| **FRASTA** | Fracture surface topography analysis. Use of surface topography to assess fracture mechanics. |
| **Fraunhoffer diffraction** | Intensity pattern in the back focal plane of a lens. Usually the image of the source of light modulated by surface scatter. Typified by superposition of plane wavefront. |
| **Fresnel diffraction** | Intensity pattern produced by point source. Typified by spherical wavefronts. |
| **Fretting** | Wear produced by contact of bodies with lateral movement and vibration usually involving trapped debris. |
| **Function** | The application of a workpiece. |
| **Functional parameter** | Parameter which is important in a functional context. |
| **Fundamental motion** | Error motion of machine tool. |
| **Gabor transforms** | Space frequency function based on Gaussian weighting function. |
| **Gamma distribution** | The factorial distribution. |
| **Gaussian** | Distribution having an exponential with squared terms of the variable. |
| **Gaussian filter** | Filter having a Gaussian weighting function or frequency transmission characteristics. |
| **Geometric surface** | The ideal surface defined by the drawing or specification. |
| **Glossmeter** | Instrument for measuring light scatter from surface. |
| **Goodness of fit test** | Chi-squared test of the equivalence of hypotheses. |
| **Gram–Charlier series** | Method of characterizing a probability density function. |
| **Grubler's equation** | Mobility of linkages in terms of links and joints. |
| **Guard band** | Attenuation band in frequency characteristics isolating two blocks. |
| **Hadamard function** | Clipped signal transform. |
| **Harmonic weighting function** | Function describing the distortion of measured harmonic coefficients by a measuring technique. |
| **Hartley transform** | Transform similar to Fourier transform but having no complex components. |
| **Hatchet stylus** | Stylus used in measuring form. The hatchet shape integrates out the roughness. |
| **Hazard** | Rate of change of failure. |
| **Helical track** | Traverse technique sometimes used to measure cylindricity. Combines one linear and one rotary motion. |
| **Helmholtz–Kirchoff integral** | Equation which describes electromagnetic radiation properties. |
| **Hermite polynomials** | Related to differentials of Gaussian distribution. |
| **Heterodyne methods** | Techniques using two independent modes of measurement such as two frequencies or two polarizations. |
| **Hexagonal sampling** | Sampling pattern using seven points; usually one central, surrounded symmetrically by six. |
| **High pass filter** | Filter technique passing only small wavelengths. |
| **Hill climbing** | Technique of optimization. |
| **Hip prosthesis** | Replacement hip joint of partially spherical shape. |
| **Hole** | A closed contour on an areal map. |
| **Holography** | Photograph containing phase information trapped by means of reference beams. |
| **Homodyne methods** | Techniques using similar modes of measurement. |
| **Hybrid parameter** | A parameter such as slope derived from two or more independent variables. |
| **Hypergeometric function** | A three argument function used for characterization. |
| **Hypotrochoid** | A geometric figure of the family of trochoids. |
| **Impulse response** | The response of a system to an impulse. |
| **Inferior envelope** | Line through a succession of valleys. |
| **Instantaneous radius** | The radius of curvature at a point on the profile. |
| **Instantaneous axis** | Centre of rotation at a given time. |
| **Instrument capability** | Guaranteed performance specification of instrument. |
| **Integrated damping** | Damping coefficient optimized over a wide frequency band. |

| | |
|---|---|
| **Interpolation** | Derivation of the value of a signal between two known points. |
| **Intrinsic equation** | Equation developed from the profile itself. |
| **Inverse scatter problem** | The deduction of the surface statistics from the statistics of the scattered wavefront. |
| **Isotropy** | The uniformity of pattern of the surface. |
| **Iterative relation** | A relationship between values of system output and system input for various instances in time or space. |
| **Jacobian** | Relationship between the differential coefficients in different domains. |
| **JKR** | Model based on Johnson to explain adhesion between solids. |
| **Johnson noise** | Electronic noise dependent on input resistance. |
| **Joint probability density function** | Function describing statistical behaviour at an infinitesimal point. When integrated between two points gives probability. |
| **JPDF** | As above. |
| ***K* value** | Semi-empirical number found by Jakeman for fractal type surfaces. |
| **Kelvin clamp** | Method of achieving six constraints by means of kinematic location. |
| **Kinematics** | Laws of constraint of a free body. |
| **Kirchoff laws** | Laws of optical behaviour. |
| **Kitagawa plot** | Logarithmic plots of variables. |
| **Kurtosis** | Fourth central moment of a distribution. |
| **Lagrangian multipliers** | Method used to evaluate best fit coefficients usually with conditional differentials. |
| **Langmuir–Bloggett films** | Molecularly thin films. |
| **Laplace transform** | Transform used to determine the response of a system to generalized and impulsive inputs. |
| **Laser waist** | Minimum focused width of laser beam. |
| **Lateral roughness** | Roughness at right angles to movement. |
| **Lay** | Pattern of areal surface finish. Direction of the prevailing texture of the surface, usually determined by the method of manufacture. |
| **Least squares centre** | Centre position determined by a least squares method. |
| **Least squares cylinder** | A cylinder derived from the least squares straight line determined from the centres of the least squares circle of each measured profile. |
| **Length of profile** | Actual length of profile expressed in terms of differentials. Levelling depth $R_p$ maximum peak value in sampling length. |
| **Likelihood** | Expected value. |
| **Limaçon** | True equation of eccentric circular part when eccentric to rotate in axis of roundness instrument. |
| **Linear phase filter** | Filter having an impulse response with an axis of symmetry about a vertical axis. |
| **Linear programming** | Technique in which variables are related linearly. |
| **Lobing** | Undulations on a nominally round workpiece produced usually by clamping or chatter. |
| **Locus** | Line produced by some procedure. |
| **Log normal function** | Logarithm of the Gaussian distribution. Usually refers to distribution of extrema. |
| **Long crestedness** | Measurement of areal bandwidth of surface according to Longuett-Higgins. |
| **Longitudinal profile** | Profile resulting from the intersection of a surface by a plane parallel to the lay. |
| **Low pass filter** | Filter passing only long wavelengths. |
| **LSC** | Least squares circle. |
| **LVDT** | Linear voltage differential transformer. |
| **M system** | Measurement system based on mean line references. |
| **Map** | Areal coverage of surface. |

| | |
|---|---|
| **Markov process** | Sequence in which current value depends only on the previous one. |
| **Material ratio** | Ratio of material to air at any given level relative to mean line. |
| **Maximum material condition** | Basis of tolerancing system; minimizes the amount of machining. |
| MC | Minimum circumscribed circle. |
| MCE | Mean crest excursion. |
| **Mean line of roughness profile (M)** | Reference line in the evaluation length such that the area enclosed between it and the profile has equal positive and negative values. |
| **Mechanical loop** | Loop linking the mechanical elements usually of the metrology calliper. |
| **Meter cut-off** | Same as cut-off, filter cut-off. |
| **Method error or divergence** | Deviation in result obtained by using methods based on different principles. |
| **Method of exact fractions** | Method using multiple wavelengths in which, by measuring fracture of fringes, step heights can be measured. |
| **Metrology loop (measuring loop)** | The linking of all the mechanical elements making up the metrology calliper. |
| MFM | Magnetic force microscope. |
| MI | Minimum inscribed circle. |
| **Midpoint locus** | Locus produced by plotting outcome of running average procedure. |
| **Minimum phase** | System in which amplitude and phase are explicitly related. |
| **Minimum zone centre** | Centre position based on minimum zone procedure. |
| **Minimum zone references** | Zonal methods which have minimum separation. |
| MND | Multinormal distribution or multi-Gaussian distribution. |
| **Moiré fringes** | Fringes produced by overlaying of true structures of equal spacing at an angle. |
| **Motif** | Procedure for producing envelopes for use as reference lines. |
| **Motion copying** | Machining in which tool is forced to follow a predetermined geometric path. |
| **Motion error** | Error in machine tool. |
| **Multiprocess** | Manufacturing process such as plateau honing in which more than one process make up the final surface. |
| **Multivariate normal distributon** | Same as MND. |
| MZ | Minimum zone circle. |
| **Non-Newtonian** | System in lubrication in which viscosity changes with pressure. |
| **Normal equations** | Differential equations produced from the summation of the squared deviations of observed and ideal values. |
| **Numerical analysis** | Analysis of the discrete form of the surface by numerical rules. |
| **Numerical aperture** | Factor influencing resolution and depth of focus of optical element. |
| **Numerical model** | Discrete model of surface feature, i.e. three point model of peak. |
| **Nyquist criterion** | Frequency criterion advocating sampling at twice the highest frequency of interest. |
| **Objective speckle** | Intensity pattern produced when rough surface is illuminated by a laser. |
| **One number parameter** | The concept of attempting to judge the complete merit of a workpiece by one number. |
| OOR | Out of roundness. |
| **Optimized datum plane** | Reference plane obtained by optimized procedure. |
| **Ordinate** | Height measurement, usually in discrete form. |
| **Out of roundness** | OOR |
| **Ovality** | Maximum difference between length of two axes through object centre, usually but not necessarily at right angles. |
| **Parallax** | Produces distortion in scanning electron microscopes. |
| **Parameter** | Feature to be quantified. |
| **Parameterization** | The process of representing a function in terms of an indirect variable. |
| **Parasitic movement** | Spurious secondary movement or vibration. |
| **Pareto curve** | Most effort should be expended on most significant effects. Pareto's law. |
| **Partial arc** | An incomplete curve, usually circular. |
| PCF | Pressure concentration factor. |

| | |
|---|---|
| **Peak** | Maximum of profile between two adjacent valleys. |
| **Peak density** | Number of peaks per unit distance. |
| **Peak-to-valley (roundness)** | Radial separation of two concentric circles which are themselves concentric to the reference and totally enclose the measured profile. |
| **Periodic profile** | Profile which can be described by a periodic function, e.g. turned, milled. |
| **Perturbation methods** | Method for solving complex differential equations using small variations. |
| **Phase corrected filter** | Filter with no phase distortion. |
| **PHL** | Plastohydrodynamic lubrication. |
| **Physical metrology** | Metrology of physical parameters such as hardness and stress. |
| **Pick-up** | Device used to detect information. |
| **Planimeter** | Analogue device for measuring area. |
| **Plasticity index** | Numerical index for predicting plastic flow. |
| **Plug gauge centre** | Centre found by maximum inscribing circle algorithm. |
| **Poincaré section** | Two-dimensional section in 3D state space, used in chaos theory. |
| **Polar coordinates** | Use of radial variations as a function of angle. |
| **Polar distortion** | Deviations in shape of circle produced as a result of nature of roundness instruments. |
| **Polygonation** | The breakdown of a circular signal into a polygon form by means of triangles. |
| **Polynomial fit** | The fitting of a smooth polynomial through a set of discrete ordinals. |
| **Power spectral density** | The limiting value of the periodgram — the square of the Fourier coefficients. |
| **Pressure copying** | Topography generated as a result of pressure input. |
| **Pressure distribution** | Variation of pressure as a function of distance within the contact zone. |
| **Prestressing** | Application of fixed stress usually to position working zone in a more predictable region. |
| **Primal** | Conventional method of tackling exchange algorithms. |
| **Primary cutting edge** | Cutting along the track of cutting tool. |
| **Primary profile** | The profile resulting after the real profile has been investigated by a finite stylus. |
| **Process parameters** | Parameters such as depth of cut and feed typical of a manufacturing process. |
| **Profile** | A section taken through a workpiece. |
| **Profile length** | The true length of the profile as opposed to its $x$ dimension. |
| **Profile parameter** | An attempt to quantify the profile according to some procedure. |
| **Propagations of errors** | The resultant error in a system in terms of errors in its constituents. |
| **PSD** | Power spectral density. |
| **Pseudo-kinematics** | Kinematic design that allows a certain amount of elastic averaging. |
| **Pseudo-random sequence** | A random sequence generator that repeats after a limited number of operations. |
| **Quantization** | The breaking down of an analogue height into discrete binary units. |
| **Quill** | Housing for roundness instrument spindle. |
| **Radial** | Coordinate measured from a point irrespective of direction. |
| **Radial motion** | Movement of tool or workpiece in radial direction in error motion. |
| **Rain fall count** | Method of adding up effective count in fatigue. |
| **Raised cosine** | Unity added to cosine to make value always positive. Used in lag windows. |
| **Rake** | Angle of tool presented to workpiece. |
| **Raman scattering** | Non-elastic scattering—optical mode of vibration. |
| **Random error** | Error which cannot be predicted. Sometimes called stochastic error. |
| **Random process analysis** | Autocorrelation, power spectral density and probability density. |
| **Random profile** | Profile which can be described by a random function, e.g. ground, shot blasted. |
| **Range (measuring range)** | Usually taken to be the range over which a signal can be obtained as some function, not necessarily linear, of displacement. |
| **Rautiefe** | Original name for $R_t$. |
| **Ray tracing** | Optical behaviour based on geometrical optics. |
| **Rayleigh criterion** | Criterion of optical resolution. |

| | |
|---|---|
| **Reaction** | Force from surface responding to stylus, or other imposed force. |
| **Readability** | The ability to detect a change in value of a variable and represent it meaningfully. |
| **Real surface** | Surface limiting the body and separating it from the surrounding media. |
| **Real surface profile** | The line of intersection of a real surface and a plane. |
| **Recursive filters** | Filter whose current output can be expressed in terms of current inputs and past outputs. |
| **Reference circle** | Perfect circle of same size as workpiece from which deviations of roundness are measured. |
| **Reference cylinder** | Perfect cylinder of same size as workpiece from which deviations of cylindricity are measured. |
| **Reference line** | Line constructed from geometrical data from which certain features can be measured. Most often used to separate roughness and waviness. |
| **Reference line** | The line relative to which assessment of profile parameters are determined. |
| **Reference surface** | The surface relative to which roughness parameters are determined. |
| **Referred part** | A perfect part estimated from existing data. |
| **Regression** | The finding of a, usually linear, relationship between variables by least squares methods. |
| **Relaxation** | The tendency of a body to sink to a position of lower potential energy with time. |
| **Relocation profilometry** | Techniques used to ensure that a workpiece can be positioned time after time in the same place. |
| **Residual stress** | Stresses usually left in the subsurface of a workpiece as a result of machining. |
| **Resolution** | The ability to separate out two adjacent lateral details. |
| **Reversal methods** | Method used to separate out reference and workpiece values of geometry by changing the dependence of one on another. |
| **Reynold's equation** | Equation governing the pressure between two bodies in relative lateral motion separated by a fluid film. |
| **Reynold's roughness** | Roughness which is small compared with the film thickness. |
| **Ring gauge centre** | Centre found by means of the minimum circumscribing circle algorithm. |
| **RMSCE** | Root mean square crest excursion. |
| **Roughness** | Marks left on surface as a result of machining. |
| **Roughness evaluation length (1)** | Length of the reference line over which the mean values of the parameters are determined. |
| **Roughness profile** | A real profile after modification by filtering, e.g. elec filter, stylus. |
| **Roundness typology** | Characterization of a workpiece according to the coefficients of the Fourier components of the roundness signal. |
| **Running-in** | The process of easing a workpiece into a steady state wear regime (shakedown). |
| **Runout** | Total deviation as seen by instrument probe. |
| **_S_ number** | Sputtering rate. |
| **_s_ plane** | Mapping of complex plane. |
| **Sampling** | The process of taking discrete values of an analogue signal. |
| **Sampling length (1) (roughness)** | The basic length over which roughness is measured and is equal to the cut-off length. |
| **Sampling rate** | The rate at which discrete measurements are made. |
| **Scalar theory** | Theory of electromagnetic radiation in which phase and amplitude of the wavefront is important. |
| **Scatterometer** | Instrument for measuring the light scattered from a surface. |
| **SCM** | Scanning confocal microscope. |
| **Scuffing wear** | Wear produced by failure due to thermal runaway. |
| **Secondary cutting edge** | Cutting produced in the axial direction. |
| **Seidal aberration** | Errors in optical wavefront produced by second-order faults in the optical system. |

| | |
|---|---|
| **Self-affine** | Property which allows equivalence if one axis is changed. |
| **Self-similarity** | Property which allows equivalence of properties at all scales of size. |
| SEM | Scanning electron microscope. |
| **Sensitive direction** | Direction in which the rate of change of energy conversion is a maximum. |
| **Sensitivity** | The energy or information conversion rate with displacement. |
| **Shadowing** | Produced by detail on the surface being washed out by oblique angle illumination. |
| **Shakedown** | Settling to final condition (e.g. running-in). |
| **Sheppard's correction** | Numerical noise produced as a result of finite quantization. |
| **Shot noise** | Noise produced by variation in standing electric currents akin to Brownian motion. |
| **Simplex method** | An exchange algorithm used in linear programming. |
| **Simpson's rule** | Rule for numerical integration based on fitting parabolas between points. |
| **Singular value decomposition** | Method of rearranging matrices so as to be easy to solve. |
| **Skew** | Third central moment of a distribution. |
| **Skew limaçon** | Limaçon shape found in cross-section of tilted cylinder. |
| **Skid** | Mechanical filter, used as reference. |
| **Space-frequency functions** | Functions such as Wigner which have two arguments, one in space and one in frequency. |
| **Spalling** | Another word for pitting. |
| **Spark-out** | A means of ensuring that all possible elements on a grinding wheel actually cut. |
| **Spatial wavelength** | Wavelength on surface expressed in millimetres not seconds. |
| **Spanzipfal** | Triangular piece of material left on surface during cutting. |
| **Speckle** | Grainy appearance of object produced in laser tight by aperture limited optics. |
| **Spectroscopy** | The breaking down of a signal into frequency components. |
| **Specular reflection** | Light scattered at the reflected angle. |
| **Speed ratio** | Ratio of workpiece circumference speed to that of wheel circumference speed. |
| **Sphericity** | Departure of geometry from a true sphere. |
| **Spline function** | Mathematical function usually of cubic nature derived from use of elastic wooden spline. |
| **Squareness** | Errors in position of two planes nominally at right angles. Expressed as a zone rather than as an angle. |
| **Squeeze film** | Lubrication film in the presence of normal vibration. |
| **Straightness** | Departure of surface geometry from a straight line. |
| **Stribeck diagram** | Shows the degree of lubrication. |
| **Standard cut-off (wavelength)** | Nominally equal to 0.8 mm or 0.03″. |
| **Standard sampling length** | As above. |
| **Standardization** | Limitation of procedures to those agreed. |
| **State-space** | Velocity displacement display of system output. |
| **Static interaction** | Two body contact without lateral movement. |
| **Stiefel exchange mechanism** | Mechanism used in linear programming. |
| **Stieltjes correlator** | Correlator using one channel of clipped signals working on reduced band width. |
| **Stiffness** | Slope of force displacement curve of contact. |
| STM | Scanning tunnelling microscope. |
| **Stoke's roughness** | Roughness which is significant proportion of film thickness. |
| **Strange attractor** | Word given to apparent focus in space state diagram. |
| **Stratified process** | Finishing produced by more than one process. Processes usually present at different heights in the topography. |
| **Structure function** | Expected values of square of difference between two values of waveform. |
| **Stylus** | Mechanical implement used to contact workpiece and to communicate displacement to transducer. |

| | |
|---|---|
| **Stylus tree** | Rod containing a number of styli at different axial positions. Used to measure spindle concentricity with rotating spindle instrument. |
| **Subjective speckle** | Speckle effect produced when surface illuminated by a laser is imaged. Speckle is in image plane. |
| **Subsurface characteristics** | Physical characteristics immediately below geometric boundary of surface. |
| **Summit** | Areal peak. |
| **Superior envelope** | Line connecting suitable peaks. |
| **Superposition integral** | Linear summation of properties. |
| **Suppressed radius** | Workpiece size removed in order to enable roundness instrument to measure the surface skin with high accuracy. |
| **Surface integrity** | Properties taken sometimes to mean physical and geometrical properties in the US. Usually stress related properties only. |
| **Surface roughness** | Irregularities in the surface texture inherent in the manufacturing process but excluding waviness and errors of form. |
| **Surface texture** | The geometry imparted. |
| **Symbolic notation of errors** | Symbolic method of writing error propagation equations. |
| **System response** | Output of the system to a given input. |
| **Systematic error** | Error attributable to deterministic cause and hence can be predicted. |
| **Tactile instrument** | Instrument which uses contact as a means of collecting information. |
| **Taper** | Cylindricity in which the radial value varies with height linearly. |
| **TEM** | Transmission electron microscope. |
| **Tetragonal sampling** | Sampling with five point analysis. |
| **Tetrahedron ball frame** | Frame built in form of tetrahedron to calibrate coordinate measuring machine. |
| **Theoretical surface finish** | Surface finish determined totally from process parameters such as feed, depth of cut, tool radius. |
| **TIS** | Total integrated scatter. |
| **Tolerance** | Range of allowable dimensional or other parameter variation of workpiece. |
| **Top down mechanics** | Mechanisms simulating contact. |
| **Topography** | The study of surface features. |
| **Topothesy** | Length of chord on surface having average slope of one radian. |
| **Traceability** | The linking of measurement procedures to national standards. |
| **Trackability** | The ability of stylus system to follow geometry without losing contact. |
| **Transducer** | Energy conversion system. |
| **Transfer function** | Dynamic relationship between input and output of system. |
| **Transmission characteristics** | Plot of transfer function. |
| **Transverse profile** | Profile resulting from the intersection of a surface by a plane normal to the surface lay. |
| **Transverse roughness** | Roughness at right angles to direction of motion. |
| **Trapezoidal rule** | Method of numerical integration linking adjacent points linearly. |
| **Tribology** | Science of rubbing parts. |
| **Trigonal sampling** | Sampling pattern using 4 point analysis of height in a profile. |
| **Truncation** | The curtailing of aspects of a function. In particular the cutting off of height in a profile. |
| **Typology** | A characterization. |
| **Union Jack pattern** | A pattern of measurements used to cover surface plates in an economic and self regulating way. |
| **Unit event of machining** | Basic unit relating finish to average grain thickness or cutting particle size. |
| **UPL, LPL, LVL, UVL** | Upper peak limit, etc. |
| **Valley** | Minimum of the profile. |
| **Van Der Waal's forces** | Atomic forces. |
| **Vector theory** | Electromagnetic theory incorporation polarization. |
| **Vee block method** | Roundness assessment by chordal measurement. |

| | |
|---|---|
| **Vernier fringes** | Fringes produced by overlay of two scales of slightly different spacing. |
| **Vibrating table method** | Method of calibrating surface texture instruments. |
| **Walsh function** | Transform using square waves rather than sine waves. |
| **Wankel engine** | Rotary engine with stator shaped as epitrochoid. |
| **Wavelet transform** | Space frequency function. |
| **Waviness** | Geometric features imparted to workpiece usually by errors in the machine tool motion. |
| **Weibull distribution** | Graph of hazard with time. |
| **Weierstrass function** | Fractal function. |
| **Weighting function** | Set of weighting factors. |
| **Wein's law** | Law relating maximum radiation to temperature |
| **Whole field technique** | Method of showing up distortion. |
| **Wigner distribution function** | Space frequency function based on quantum mechanics. |
| **Work function** | Energy required to release electrons from surfaces. |
| **Wringing** | Sticking of smooth surfaces in contact. Presence of fluid film necessary. |
| **z transform** | Discrete transformation technique for dealing with the dynamics of discrete systems. |
| **Zeeman splitting** | Splitting of laser frequency by magnetic field. |
| **Zernicke polynomials** | Polynomials used in specifying optical performance. |
| **Zero crossing** | Crossing of the mean line. |

Direction notation:

| | |
|---|---|
| $x$ | Along the surface in direction of motion. |
| $y$ | Lateral motion relative to direction of measurement. |
| $z$ | Direction of measurement normal to surface. |

Surface roughness parameters associated with the properties of irregularities in the direction of the roughness profile *height*:

| | |
|---|---|
| $z_p$ | Peak height. |
| $z_v$ | Valley depth. |
| $R_p$ | Maximum peak height. |
| $R_v$ | Maximum valley depth. |
| $R_z$ | Maximum height of roughness profile. |
| $R_{zs}$ | Ten point height of irregularities. |
| $R_c$ | Mean height of roughness profile irregularities. |
| $\bar{R}$ | Average value of roughness parameter. |
| $R_a$ | Arithmetic mean deviation. |
| $R_q$ | Root-mean-square deviation. |

Surface roughness parameters associated with irregularity properties in the profile *length* direction:

| | |
|---|---|
| $\lambda_q$ | Root-mean-square wavelength. |
| $\lambda_a$ | Average wavelength. |
| $S_m$ | Mean spacing of irregularities. |
| $S$ | Mean spacing of local peaks. |
| $l$ | Roughness profile length. |

Surface roughness parameters associated with the *form* of the profile irregularity:

| | |
|---|---|
| $S_k$ | Skewness. |
| $\Delta_q$ | Root-mean-square slope. |

| $\Delta_a$ | Arithmetic mean slope. |
| $\eta_p$ | Bearing length. |
| $t_p$ | Bearing length ratio. |

**Radial measurements**

*Least squares circle (LS)*

The least squares reference figure is calculated to be a circle fitted through the data such that the sum of the squares of the deviations inside and outside that circle is at a minimum

*Minimum zone circles (MZ)*

The minimum zone reference figure is calculated as the two concentric circles which totally enclose the data, such that the radial separation between them is a minimum.

*Minimum circumscribed circles (MC)*

The minimum circumscribed reference figure is calculated to be the smallest circle that totally encloses the data.

*Maximum inscribed circle (MI)*

The maximum inscribed reference figure is calculated to be the largest circle that is totally enclosed by the data.