

T H E

MEASUREMENT,
INSTRUMENTATION,
AND
SENSORS

H A N D B O O K

Editor-in-Chief

John G. Webster



CRC PRESS



Springer



IEEE PRESS

A CRC Handbook Published in Cooperation with IEEE Press

T H E

MEASUREMENT,
INSTRUMENTATION,
AND
SENSORS

H A N D B O O K

Editor-in-Chief
John G. Webster



A CRC Handbook Published in Cooperation with IEEE Press

This One



KSHZ-R7P-WNBW

Library of Congress Cataloging-in-Publication Data

The measurement, instrumentation, and sensors handbook ; John G. Webster, editor-in-chief.

p. cm. — (Electrical engineering handbook series)

Includes bibliographical references and index.

ISBN 3-540-64830-5

1. Physical measurements—Handbooks, manuals, etc.
 2. Mensuration—Handbooks, manuals, etc.
 3. Scientific apparatus and instruments—Handbooks, manuals, etc.
- I. Webster, John G., 1932-
II. Series.

QC39.M393 1999

530.8—dc21

98-31681
CIP

Co-published by

CRC Press LLC

2000 Corporate Blvd., N.W.

Boca Raton, FL 33431, U.S.A

(Orders from the U.S.A. and Canada (only) to CRC Press LLC)

and by

Springer-Verlag GmbH & Co. KG

Tiergartenstraße 17

D-69121 Heidelberg

Germany

(Orders from outside the U.S.A. and Canada to Springer-Verlag)

ISBN 3-540-64830-5

This book contains information obtained from authentic and highly regarded sources. Reprinted material is quoted with permission, and sources are indicated. A wide variety of references are listed. Reasonable efforts have been made to publish reliable data and information, but the author and the publisher cannot assume responsibility for the validity of all materials or for the consequences of their use.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation, without intent to infringe.

© 1999 by CRC Press LLC

No claim to original U.S. Government works

International Standard Book Number 3-540-64830-5

Library of Congress Card Number 98-31681

Printed in the United States of America 1 2 3 4 5 6 7 8 9 0

Printed on acid-free paper

Contents

Section I Measurement Characteristics

- 1 Characteristics of Instrumentation *R. John Hansman, Jr.* 1-1
- 2 Operational Modes of Instrumentation *Richard S. Figliola* 2-1
- 3 Static and Dynamic Characteristics of Instrumentation
Peter H. Sydenham..... 3-1
- 4 Measurement Accuracy *Ronald H. Dieck* 4-1
- 5 Measurement Standards *DeWayne B. Sharp* 5-1

Section II Spatial Variables Measurement

- 6 Displacement Measurement, Linear and Angular
 - 6.1 Resistive Displacement Sensors *Keith Antonelli, James Ko, and Shyan Ku*.....6-2
 - 6.2 Inductive Displacement Sensors *Halit Eren*6-15
 - 6.3 Capacitive Sensors—Displacement *Halit Eren and Wei Ling Kong*6-37
 - 6.4 Piezoelectric Transducers and Sensors *Ahmad Safari, Victor F. Janas, Amit Bandyopadhyay, and Andrei Kholkin*.....6-53
 - 6.5 Laser Interferometer Displacement Sensors *Bernhard Günther Zagar*6-65
 - 6.6 Bore Gaging Displacement Sensors *Viktor P. Astakhov*.....6-78
 - 6.7 Time-of-Flight Ultrasonic Displacement Sensors
Teklic Ole Pedersen and Nils Karlsson.....6-92
 - 6.8 Optical Encoder Displacement Sensors *J. R. René Mayer*.....6-98
 - 6.9 Magnetic Displacement Sensors *David S. Nyce*.....6-119
 - 6.10 Synchro/Resolver Displacement Sensors
Robert M. Hyatt, Jr. and David Dayton6-128
 - 6.11 Optical Fiber Displacement Sensors
Richard O. Claus, Vikram Bhatia, and Anbo Wang6-141
 - 6.12 Optical Beam Deflection Sensing *Grover C. Wetsel*.....6-156

7	Thickness Measurement	<i>John C. Brasunas, G. Mark Cushman, and Brook Lakew</i>	7-1
8	Proximity Sensing for Robotics	<i>R.E. Saad, A. Bonen, K.C. Smith, and B. Benhabib</i>	8-1
9	Distance	<i>W. John Ballantyne</i>	9-1
10	Position, Location, Altitude Measurement		
10.1	Altitude Measurement	<i>Dimitris E. Manolakis</i>	10-1
10.2	Attitude Measurement	<i>Mark A. Stedham, Partha B. Banerjee, Seiji Nishfuji, and Shogo Tanaka</i>	10-17
10.3	Inertial Navigation	<i>Halit Eren and C.C. Fung</i>	10-34
10.4	Satellite Navigation and Radiolocation	<i>Halit Eren and C.C. Fung</i>	10-48
10.5	Occupancy Detection	<i>Jacob Fraden</i>	10-62
11	Level Measurement	<i>Detlef Brumbi</i>	11-1
12	Area Measurement	<i>Charles B. Coulbourn and Wolfgang P. Buerner</i>	12-1
13	Volume Measurement	<i>René G. Aarnink and Hessel Wijkstra</i>	13-1
14	Angle Measurement	<i>Robert J. Sandberg</i>	14-1
15	Tilt Measurement	<i>Adam Chrzanowski and James M. Secord</i>	15-1
16	Velocity Measurement	<i>Charles P. Pinney and William E. Baker</i>	16-1
17	Acceleration, Vibration, and Shock Measurement	<i>Halit Eren</i>	17-1
Section III Time and Frequency Measurement			
18	Time Measurement	<i>Michael A. Lombardi</i>	18-1
19	Frequency Measurement	<i>Michael A. Lombardi</i>	19-1
Section IV Mechanical Variables Measurement — Solid			
20	Mass and Weight Measurement	<i>Mark Fritz and Emil Hazarian</i>	20-1
21	Density measurement	<i>Halit Eren</i>	21-1

59	Optical Loss	<i>Halit Eren</i>	59-1
60	Polarization Measurement	<i>Soe-Mie F. Nee</i>	60-1
61	Refractive Index Measurement	<i>G. H. Meeten</i>	61-1
62	Turbidity Measurement	<i>Daniel Harrison and Michael Fisch</i>	62-1
63	Laser Output Measurement	<i>Haiyin Sun</i>	63-1
64	Vision and Image Sensors	<i>Stanley S. Ipson and Chima Okereke</i>	64-1

IX Radiation Measurement

65	Radioactivity Measurement	<i>Bert M. Coursey</i>	65-1
66	Radioactivity Measurement	<i>Larry A. Franks, Ralph B. James, and Larry S. Darken</i>	66-1
67	Charged Particle Measurement	<i>John C. Armitage, Madhu S. Dixit, Jacques Dubeau, Hans Mes, and F. Gerald Oakham</i>	67-1
68	Neutron Measurement	<i>Steven M. Grimes</i>	68-1
69	Dosimetry Measurement	<i>Brian L. Justus, Mark A. Miller, and Alan L. Huston</i>	69-1

X Chemical Variables Measurement

70	Composition Measurement		
70.1	Electrochemical Composition Measurement	<i>Michael J. Schöning, Olaf Glück, and Marion Thust</i>	70-1
70.2	Thermal Composition Measurement	<i>Mushtaq Ali, Behrooz Pahlavanpour, and Maria Eklund</i>	70-49
70.3	Kinetic Methods	<i>E.E. Uzgiris and J.Y. Gui</i>	70-61
70.4	Chromatography Composition Measurement	<i>Behrooz Pahlavanpour, Mushtaq Ali, and C.K. Laird</i>	70-74
71	pH Measurement	<i>Norman F. Sheppard, Jr. and Anthony Guiseppi-Elie</i>	71-1
72	Humidity and Moisture Measurement	<i>Gert J.W. Vischer</i>	72-1

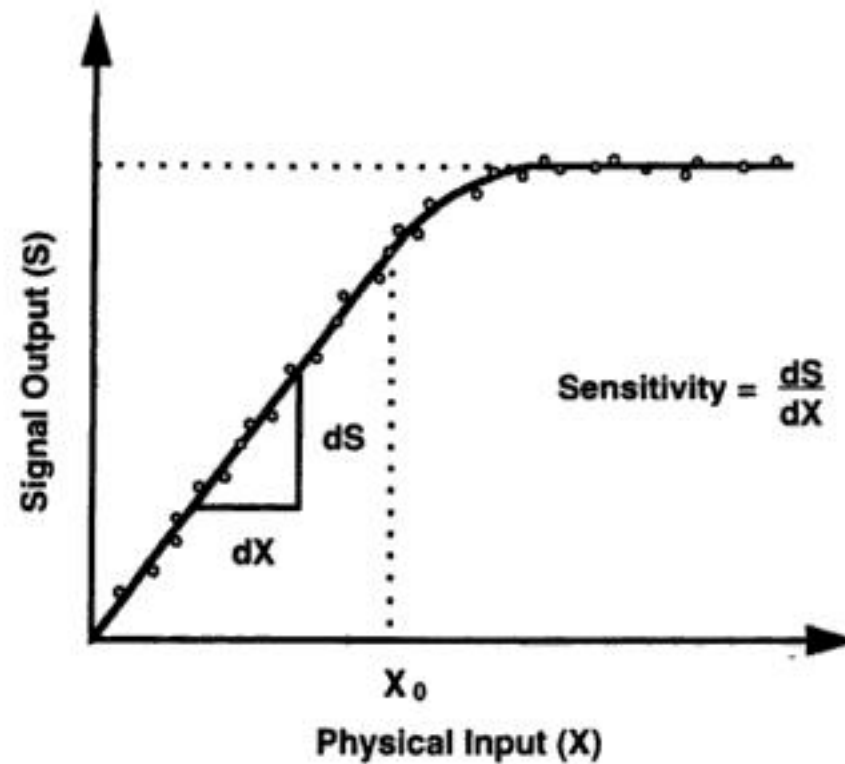


FIGURE 1.3 Calibration curve example.

can be categorized into two broad classes depending on how they interact with the environment they are measuring. *Passive sensors* do not add energy as part of the measurement process but may remove energy in their operation. One example of a passive sensor is a thermocouple, which converts a physical temperature into a voltage signal. In this case, the temperature gradient in the environment generates a thermoelectric voltage that becomes the signal variable. Another passive transducer is a pressure gage where the pressure being measured exerts a force on a mechanical system (diaphragm, aneroid or Borden pressure gage) that converts the pressure force into a displacement, which can be used as a signal variable. For example, the displacement of the diaphragm can be transmitted through a mechanical gearing system to the displacement of an indicating needle on the display of the gage.

Active sensors add energy to the measurement environment as part of the measurement process. An example of an active sensor is a radar or sonar system, where the distance to some object is measured by actively sending out a radio (radar) or acoustic (sonar) wave to reflect off of some object and measure its range from the sensor.

Calibration

The relationship between the physical measurement variable input and the signal variable (output) for a specific sensor is known as the *calibration* of the sensor. Typically, a sensor (or an entire instrument system) is calibrated by providing a known physical input to the system and recording the output. The data are plotted on a calibration curve such as the example shown in Figure 1.3. In this example, the sensor has a linear response for values of the physical input less than X_0 . The *sensitivity* of the device is determined by the slope of the calibration curve. In this example, for values of the physical input greater than X_0 , the calibration curve becomes less sensitive until it reaches a limiting value of the output signal. This behavior is referred to as *saturation*, and the sensor cannot be used for measurements greater than its saturation value. In some cases, the sensor will not respond to very small values of the physical input variable. The difference between the smallest and largest physical inputs that can reliably be measured by an instrument determines the *dynamic range* of the device.

Modifying and Interfering Inputs

In some cases, the sensor output will be influenced by physical variables other than the intended measurand. In Figure 1.4, X is the intended measurand, Y is an *interfering input*, and Z is a *modifying input*. The interfering input Y causes the sensor to respond in the same manner as the linear superposition of Y and the intended measurand X . The measured signal output is therefore a combination of X and Y ,

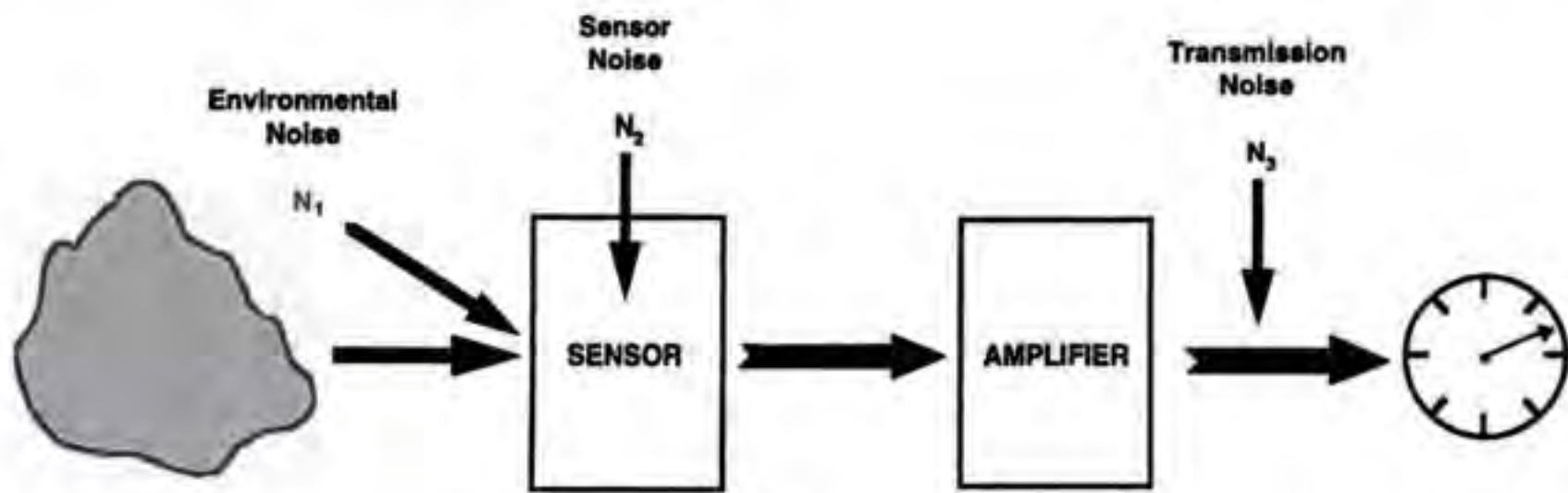


FIGURE 1.8 Instrument model with noise sources.

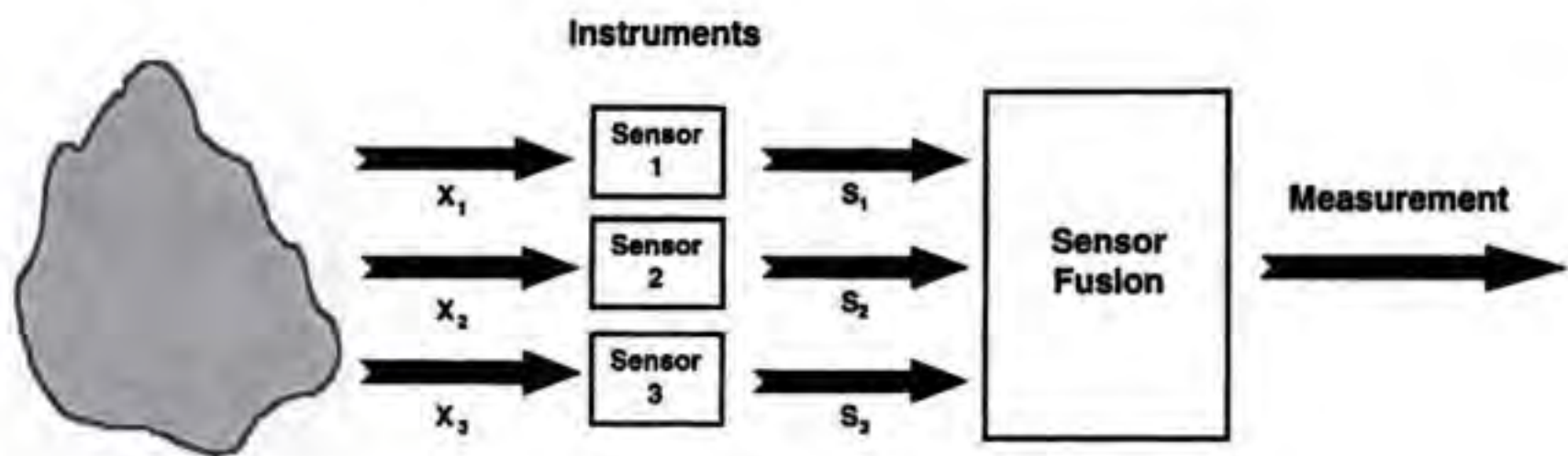


FIGURE 1.9 Example of sensor fusion.

Random error generating noise can also be introduced at each stage in the measurement process, as shown schematically in Figure 1.8. Random interfering inputs will result in noise from the measurement environment N_1 that are introduced before the sensor, as shown in the figure. An example would be background noise received by a microphone. Sensor noise N_2 can also be introduced within the sensor. An example of this would be thermal noise within a sensitive transducer, such as an infrared sensor. Random motion of electrons, due to temperature, appear as voltage signals, which are apparently due to the high sensitivity of the device. For very sensitive measurements with transducers of this type (e.g., infrared detectors), it is common to cool the detector to minimize this noise source.

Noise N_3 can also be introduced in the transmission path between the transducer and the amplifier. A common example of transmission noise in the U.S. is 60 Hz interference from the electric power grid that is introduced if the transmission path is not well grounded, or if an inadvertent electric ground loop causes the wiring to act as an antenna.

It is important to note that the noise will be amplified along with the signal as it passes through the amplifier in Figure 1.8. As a consequence, the figure of merit when analyzing noise is not the level of the combined noise sources, but the *signal to noise ratio (SNR)*, defined as the ratio of the signal power to the power in the combined noise sources. It is common to report SNR in decibel units.

The SNR is ideally much greater than 1 (0 dB). However, it is sometimes possible to interpret a signal that is lower than the noise level if some identifying characteristics of that signal are known and sufficient signal processing power is available. The human ability to hear a voice in a loud noise environment is an example of this signal processing capability.

Sensor Fusion

The process of *sensor fusion* is modeled in Figure 1.9. In this case, two or more sensors are used to observe the environment and their output signals are combined in some manner (typically in a processor) to

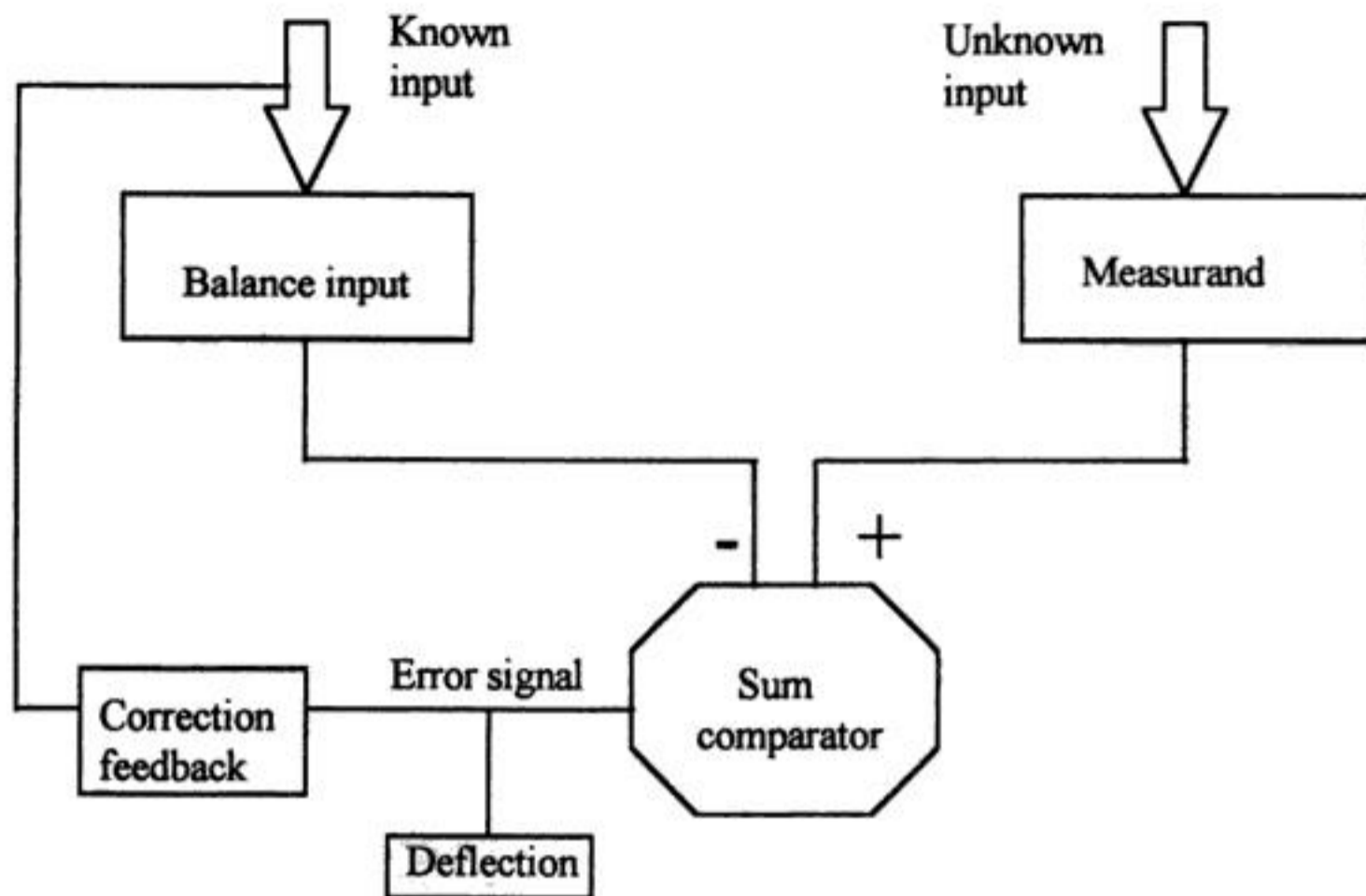


FIGURE 2.2 A null instrument requires input from two sources for comparison.

Deflection instruments are the most common of measuring instruments. The relationship between the measurand and the prime element or measuring circuit can be a direct one, with no balancing mechanism or comparator circuits used. The proportional response can be manipulated through signal conditioning methods between the prime element and the output scale so that the output reading is a direct indication of the measurand. Effective designs can achieve a high accuracy, yet sufficient accuracy for less demanding uses can be achieved at moderate costs.

An attractive feature of the deflection instrument is that it can be designed for either static or dynamic measurements or both. An advantage to deflection design for dynamic measurements is in the high dynamic response that can be achieved. A disadvantage of deflection instruments is that by deriving its energy from the measurand, the act of measurement will influence the measurand and change the value of the variable being measured. This change is called a loading error. Hence, the user must ensure that the resulting error is acceptable. This usually involves a careful look at the instrument input impedance for the intended measurement.

A spring scale is a good, simple example of a deflection instrument. As shown in Figure 2.3, the input weight or measurand acts on a plate-spring. The plate-spring serves as a prime element. The original position of the spring is influenced by the applied weight and responds with a translational displacement, a deflection x . The final value of this deflection is a position that is at equilibrium between the downward force of the weight, W , and the upward restoring force of the spring, kx . That is, the input force is balanced against the restoring force. A mechanical coupler is connected directly or by linkage to a pointer. The pointer position is mapped out on a corresponding scale that serves as the readout scale. For example, at equilibrium $W = kx$ or by measuring the deflection of the pointer the weight is deduced by $x = W/k$.

The flow diagram logic for a deflection instrument is rather linear, as shown in Figure 2.4. The input signal is sensed by the prime element or primary circuit and thereby deflected from its initial setting. The deflection signal is transmitted to signal conditioners that act to condition the signal into a desired form. Examples of signal conditioning are to multiply the deflection signal by some scaler magnitude, such as in amplification or filtering, or to transform the signal by some arithmetic function. The conditioned signal is then transferred to the output scale, which provides the indicated value corresponding to the measurand value.

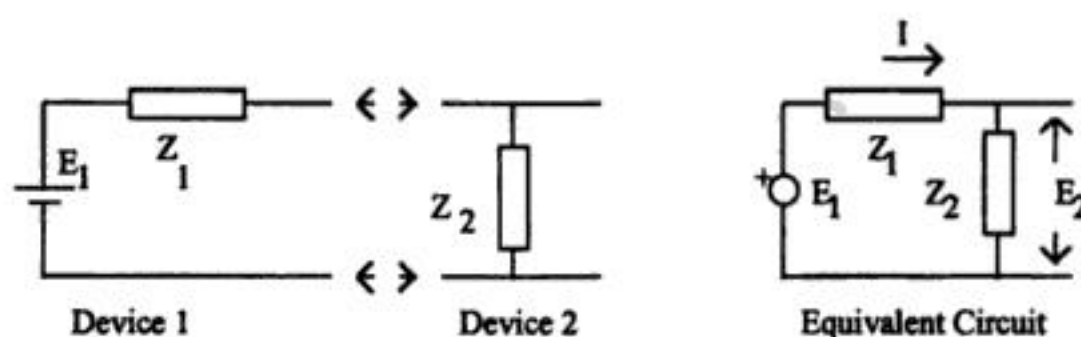


FIGURE 2.8 An equivalent circuit is formed by applying a measuring instrument to the output terminals of an instrument.

the output terminals of device 1 to the input terminals of device 2 creates the equivalent circuit also shown in Figure 2.7. The potential actually sensed by device 2 will be

$$E_2 = E_1 \frac{1}{1 + Z_1/Z_2}$$

The difference between the actual potential E_1 at the output terminals of device 1 and the measured potential E_2 is a **loading error** brought on by the input impedance of measuring device 2. It is clear that a high input impedance Z_2 relative to Z_1 minimizes this error. A general rule is for the input impedance to be at least 100 times the source impedance to reduce the loading error to 1%.

In general, null instruments and null methods will minimize loading errors. They provide the equivalent of a very high input impedance to the measurement, minimizing energy drain from the measured system. Deflection instruments and deflection measuring techniques will derive energy from the process being measured and therefore require attention to proper selection of input impedance.

Defining Terms

Analog sensor: Sensors that output a signal that is continuous in both magnitude and time (or space).

Deflection instrument: A measuring device whose output deflects proportional to the magnitude of the measurand.

Digital sensor: Sensors that output a signal that is discrete (noncontinuous) in time and/or magnitude.

Input impedance: The impedance measured across the input terminals of a device.

Loading error: That difference between the measurand and the measuring system output attributed to the act of measuring the measurand.

Measurand: A physical quantity, property, or condition being measured. Often, it is referred to as a measured value.

Null instrument: A measuring device that balances the measurand against a known value, thus achieving a null condition. A null instrument minimizes measurement loading errors.

Readout: This is the display of a measuring system.

Resolution: This is the least count or smallest detectable change in measurand capable.

Sensor: The portion of a measurement system that responds directly to the physical variable being measured.

Further Information

E. O. Doebelin, *Measurement Systems, 4th ed.*, New York: McGraw-Hill, 1990.

R. S. Figliola and D. E. Beasley, *Theory and Design for Mechanical Measurements, 2nd ed.*, New York: Wiley, 1995.

D. Wobschall, *Circuit Design for Electronic Instrumentation: Analog and Digital Devices from Sensor to Display, 2nd ed.*, New York: McGraw-Hill, 1987.

100	Explosion-Proof Instruments	Sam S. Khalilieh	100-1
99	Electropneumatic and Electrohydraulic Instruments: Modeling of Electrohydraulic and Electrostatic Actuators	M. Pachter and C. H. Houpis	99-1
98	Optimal Control	Halit Eren	98-1
97	PID Control	F. Greg Shinsky	97-1

XIV Control

96	Reading/Recording Devices	96.1 Graphic Recorders Herman Vermairen	96-1
		96.2 Data Acquisition Systems Edward McConnell	96-10
		96.3 Magnetic and Optical Recorders Yufeng Li	96-22
95	Light-Emitting Diode Displays	Mohammad A. Karim	95-1
94	Electroluminescent Displays	William A. Barrow	94-1
93	Plasma-Driven Flat Panel Displays	Robert T. McGrath, Ramanapathy Verasingam, William C. Moffatt, and Robert B. Campbell	93-1
92	Liquid Crystal Displays	Kalluri R. Sarma	92-1
91	Cathode Ray Tube Displays	Christopher J. Sherman	91-1
90	Human Factors in Displays	Steven A. Murray, Barrett S. Caldwell	90-1

XIII Displays

89	Electromagnetic Compatibility	89.1 Grounding and Shielding in Industrial Electronic Systems	
		Daryl Gerke, and William Kimmel	89-1
		89.2 EMI and EMC Test Methods Jeffrey P. Mills	89-12
88	Sensor Networks and Communication	Robert M. Crovella	88-1
87	Telemetry	Albert Lozano-Nieto	87-1

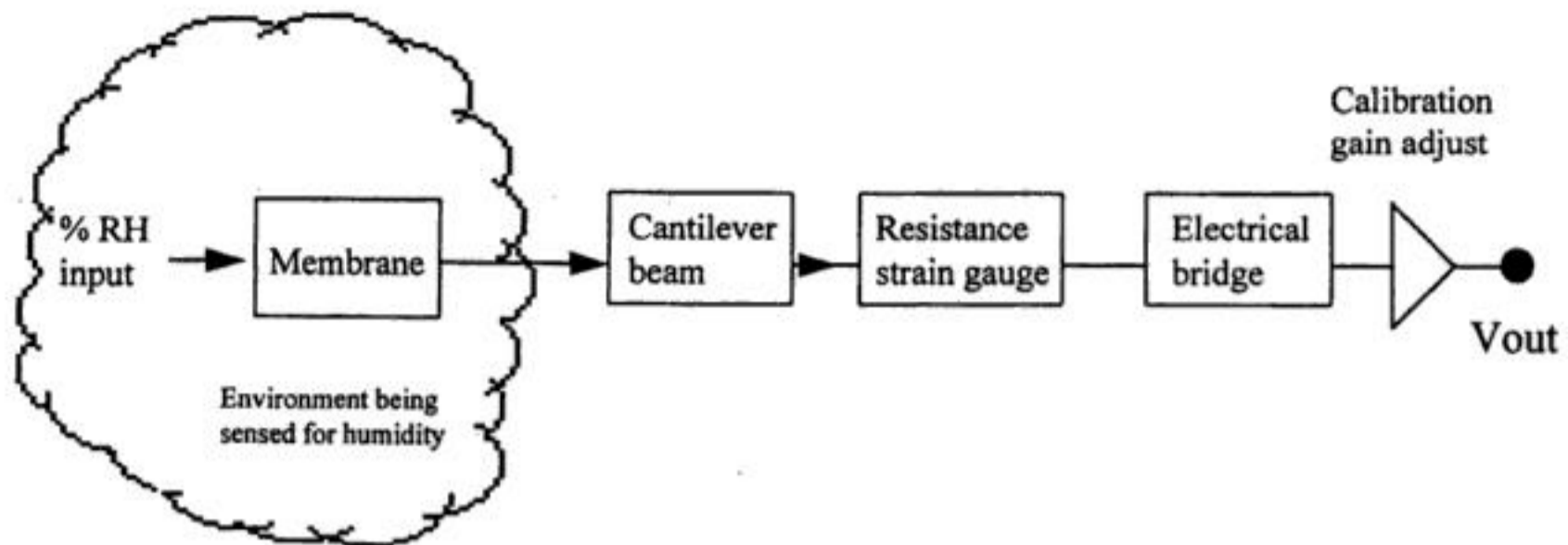


FIGURE 3.4 Instruments are formed from a connection of blocks. Each block can be represented by a conceptual and mathematical model. This example is of one type of humidity sensor.

3.1 Static Characteristics of Instrument Systems

Output/Input Relationship

Instrument systems are usually built up from a serial linkage of distinguishable building blocks. The actual physical assembly may not appear to be so but it can be broken down into a representative diagram of connected blocks. Figure 3.4 shows the block diagram representation of a humidity sensor. The sensor is activated by an input physical parameter and provides an output signal to the next block that processes the signal into a more appropriate state.

A key generic entity is, therefore, the relationship between the input and output of the block. As was pointed out earlier, all signals have a time characteristic, so we must consider the behavior of a block in terms of both the static and dynamic states.

The behavior of the static regime alone and the combined static and dynamic regime can be found through use of an appropriate mathematical model of each block. The mathematical description of system responses is easy to set up and use if the elements all act as linear systems and where addition of signals can be carried out in a linear additive manner. If nonlinearity exists in elements, then it becomes considerably more difficult — perhaps even quite impractical — to provide an easy to follow mathematical explanation. Fortunately, general description of instrument systems responses can be usually be adequately covered using the linear treatment.

The output/input ratio of the whole cascaded chain of blocks 1, 2, 3, etc. is given as:

$$[\text{output/input}]_{\text{total}} = [\text{output/input}]_1 \times [\text{output/input}]_2 \times [\text{output/input}]_3 \dots$$

The output/input ratio of a block that includes both the static and dynamic characteristics is called the *transfer function* and is given the symbol G .

The equation for G can be written as two parts multiplied together. One expresses the static behavior of the block, that is, the value it has after all transient (time varying) effects have settled to their final state. The other part tells us how that value responds when the block is in its dynamic state. The static part is known as the *transfer characteristic* and is often all that is needed to be known for block description.

The static and dynamic response of the cascade of blocks is simply the multiplication of all individual blocks. As each block has its own part for the static and dynamic behavior, the cascade equations can be rearranged to separate the static from the dynamic parts and then by multiplying the static set and the dynamic set we get the overall response in the static and dynamic states. This is shown by the sequence of Equations 3.1 to 3.4.

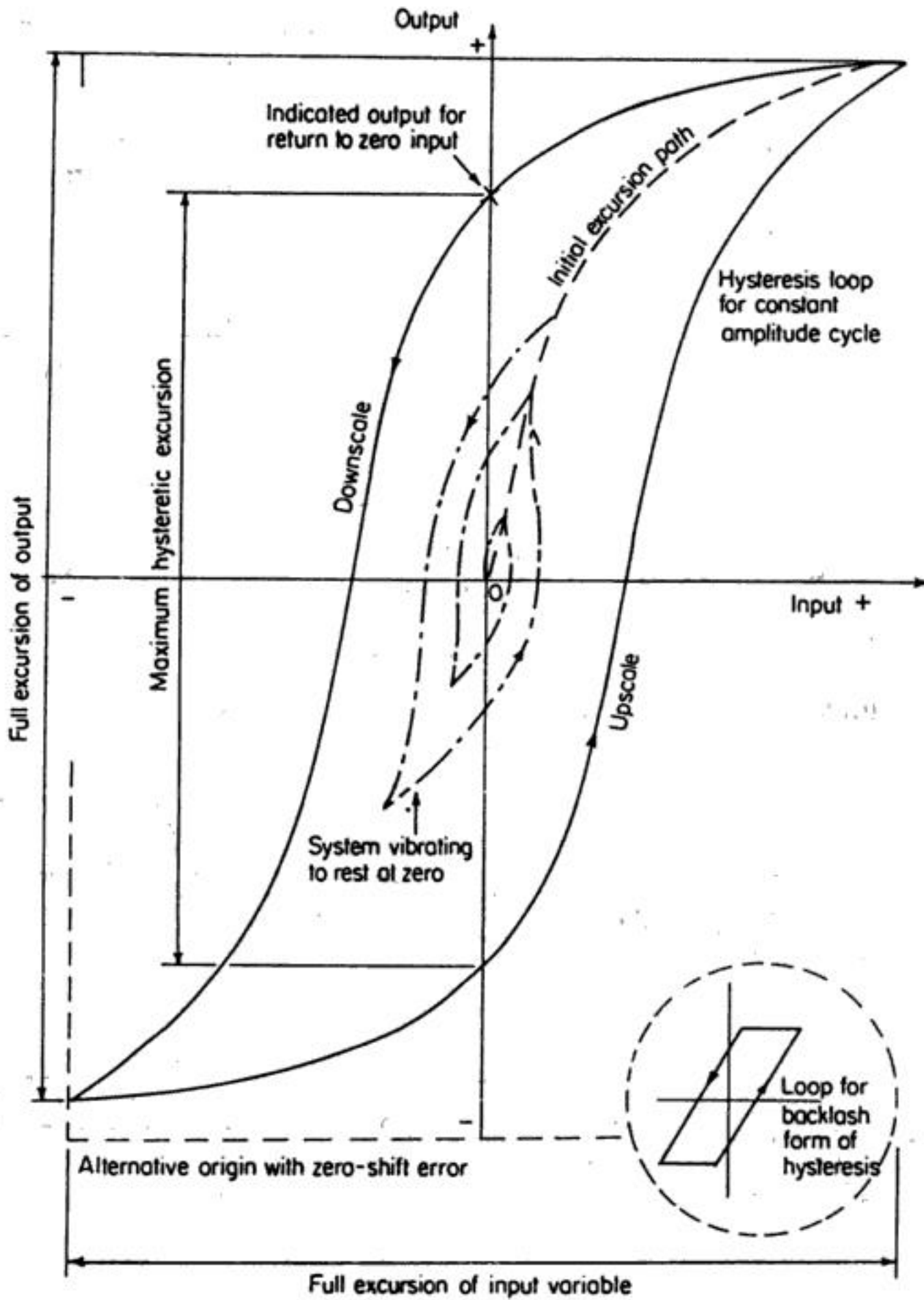


FIGURE 3.7 Generalized graph of output/input relationship where hysteresis is present. (From P. H. Sydenham, *Handbook of Measurement Science*, Vol. 2, Chichester, U.K., John Wiley & Sons, 1983. With permission.)

can be amplified by a single semiconductor element. Raising the level of all of the waveform equally takes all parts into the reasonably linear zone of an amplifier, allowing more faithful replication. If bias were not used here, then the lower half cycle would not be amplified, resulting in only the top half appearing in the output.

Error of Nonlinearity

Ideally, it is often desired that a strictly linear relationship exists between input and output signals in amplifiers. Practical units, however, will always have some degree of nonconformity, which is called the *nonlinearity*. If an instrument block has constant gain for all input signal levels, then the relationship graphing the input against the output will be a straight line; the relationship is then said to be linear.

To properly appreciate instrumentation design and its use, it is now necessary to develop insight into the most commonly encountered types of dynamic response and to develop the mathematical modeling basis that allows us to make concise statements about responses.

If the transfer relationship for a block follows linear laws of performance, then a generic mathematical method of dynamic description can be used. Unfortunately, simple mathematical methods have not been found that can describe all types of instrument responses in a simplistic and uniform manner. If the behavior is nonlinear, then description with mathematical models becomes very difficult and might be impracticable. The behavior of nonlinear systems can, however, be studied as segments of linear behavior joined end to end. Here, digital computers are effectively used to model systems of any kind provided the user is prepared to spend time setting up an adequate model.

Now the mathematics used to describe linear dynamic systems can be introduced. This gives valuable insight into the expected behavior of instrumentation, and it is usually found that the response can be approximated as linear.

The modeled response at the output of a block G_{result} is obtained by multiplying the mathematical expression for the input signal G_{input} by the transfer function of the block under investigation $G_{response}$, as shown in Equation 3.5.

$$G_{result} = G_{input} \times G_{response} \tag{3.5}$$

To proceed, one needs to understand commonly encountered input functions and the various types of block characteristics. We begin with the former set: the so-called *forcing functions*.

Forcing Functions

Let us first develop an understanding of the various types of input signal used to perform tests. The most commonly used signals are shown in Figure 3.12. These each possess different valuable test features. For example, the sine-wave is the basis of analysis of all complex wave-shapes because they can be formed as a combination of various sine-waves, each having individual responses that add to give all other wave-shapes. The step function has intuitively obvious uses because input transients of this kind are commonly encountered. The ramp test function is used to present a more realistic input for those systems where it is not possible to obtain instantaneous step input changes, such as attempting to move a large mass by a limited size of force. Forcing functions are also chosen because they can be easily described by a simple mathematical expression, thus making mathematical analysis relatively straightforward.

Characteristic Equation Development

The behavior of a block that exhibits linear behavior is mathematically represented in the general form of expression given as Equation 3.6.

$$\dots\dots\dots a_2 d^2 y / dt^2 + a_1 dy / dt + a_0 y = x(t) \tag{3.6}$$

Here, the coefficients a_2 , a_1 , and a_0 are constants dependent on the particular block of interest. The left-hand side of the equation is known as the *characteristic equation*. It is specific to the internal properties of the block and is not altered by the way the block is used.

The specific combination of forcing function input and block characteristic equation collectively decides the combined output response. Connections around the block, such as feedback from the output to the input, can alter the overall behavior significantly: such systems, however, are not dealt with in this section being in the domain of feedback control systems.

Solution of the combined behavior is obtained using Laplace transform methods to obtain the output responses in the time or the complex frequency domain. These mathematical methods might not be familiar to the reader, but this is not a serious difficulty for the cases most encountered in practice are

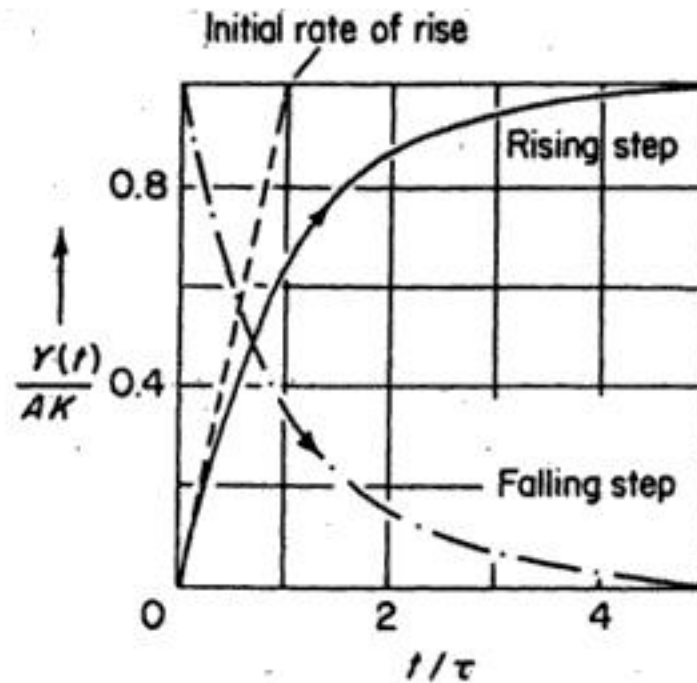


FIGURE 3.14 The step response for all first-order systems is covered by these two normalized graphs. (From P. H. Sydenham, *Handbook of Measurement Science*, Vol. 2, Chichester, U.K., John Wiley & Sons, 1983. With permission.)

If the input is a sine-wave, the output response is quite different; but again, it will be found that there is a general solution for all situations of this kind. As before, the input forcing equation is multiplied by the characteristic equation for the first-order block and Laplace transformation is used to get back to the time domain response. After rearrangement into two parts, this yields:

$$y(t) = \left[AK\tau\omega e^{-t/\tau} / (\tau^2\omega^2 + 1) \right] + \left[AK / (\tau^2\omega^2 + 1)^{1/2} \cdot \sin(\omega t + \phi) \right] \quad (3.16)$$

where ω is the signal frequency in angular radians, $\phi = \tan^{-1}(-\omega\tau)$, A the amplitude of the sine-wave input, K the gain of the first-order block, t the time in consistent units, and τ the time constant associated with the block.

The left side of the right-hand bracketed part is a short-lived, normally ignored, time transient that rapidly decays to zero, leaving a steady-state output that is the parameter of usual interest. Study of the steady-state part is best done by plotting it in a normalized way, as has been done in Figure 3.15.

These plots show that the amplitude of the output is always reduced as the frequency of the input signal rises and that there is always a phase lag action between the input and the output that can range from 0 to 90° but never be more than 90°. The extent of these effects depends on the particular coefficients of the block and input signal. These effects must be well understood when interpreting measurement results because substantial errors can arise with using first-order systems in an instrument chain.

Second-Order Blocks

If the second-order differential term is present, the response of a block is quite different, again responding in quite a spectacular manner with features that can either be wanted or unwanted.

As before, to obtain the output response, the block's characteristic function is multiplied by the chosen forcing function. However, to make the results more meaningful, we first carry out some simple substitution transformations.

The steps begin by transforming the second-order differential Equation 3.6 into its Laplace form to obtain:

$$X(s) = a_2s^2Y(s) + a_1sY(s) + a_0Y(s) \quad (3.17)$$

This is then rearranged to yield:

$$G(s) = Y(s)/X(s) = 1/a_0 \cdot 1/\left\{ (a_2/a_0)s^2 + (a_1/a_0)s + 1 \right\} \quad (3.18)$$

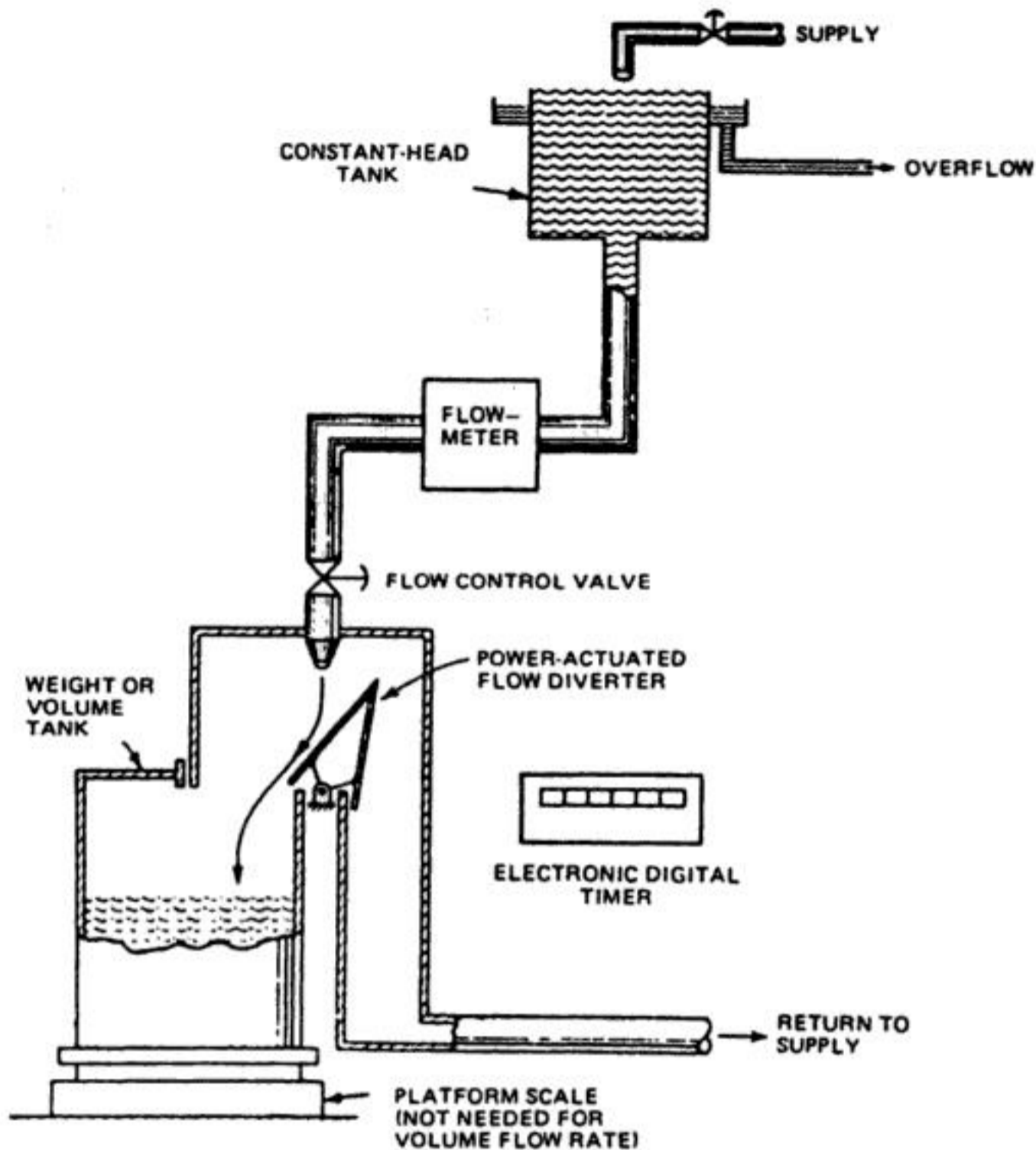


FIGURE 3.18 This practical example illustrates how flow meters are calibrated by passing a known quantity of fluid through the meter over a given time. (Originally published in P. H. Sydenham, *Transducers in Measurement and Control*, Adam Hilger, Bristol, IOP Publishing, Bristol, 1984. Copyright P. H. Sydenham.)

of comparison must necessarily be less than that required. This means that calibration is often an expensive process. Conducting a good calibration requires specialist expertise.

The method and apparatus for performing measurement instrumentation calibrations vary widely. An illustrative example of the comparison concept underlying them all is given in the calibration of flow meters, shown diagrammatically in Figure 3.18.

By the use of an overflowing vessel, the top tank provides a flow of water that remains constant because it comes from a constant height. The meter to be calibrated is placed in the downstream pipe.

The downstream is either deflected into the weigh tank or back to the supply. To make a measurement, the water is first set to flow to the supply. At the start of a test period, the water is rapidly and precisely deflected into the tank. After a given period, the water is again sent back to the supply. This then has filled the tank with a given amount of water for a given time period of flow. Calculations are then undertaken to work out the quantity of water flowing per unit time period, which is the *flow rate*. The meter was already registering a flow rate as a constant value. This is then compared with the weighed method to yield the error. Some thought will soon reveal many sources of error in the test apparatus, such as that the temperature of the water decides the volume that flows through and thus this must be allowed for in the calculations.

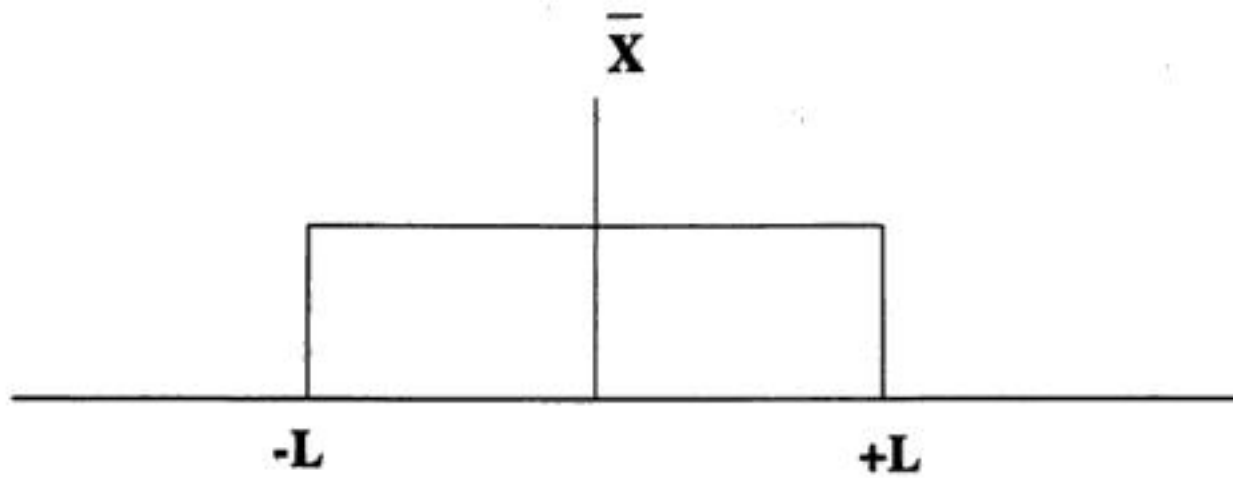


FIGURE 4.2

$$\bar{X} = \sum_{i=1}^M (X_i / N) \quad (4.5)$$

X_i = the i^{th} data point used to calculate the sample standard deviation and the average, \bar{X} , from the data

N = the number of data points used to calculate the standard deviation

$(N - 1)$ = the degrees of freedom of S_x and $S_{\bar{x}}$

M = the number of data points in the reported average test result

Note in Equation 4.4 that N does not necessarily equal M . It is possible to obtain S_x from historical data with many degrees of freedom ($[N - 1]$ greater than 30) and to run the test only M times. The test result, or average, would therefore be based on M measurements, and the standard deviation of the average would still be calculated with Equation 4.4. In that case, there would be two averages, \bar{X} . One \bar{X} would be from the historical data used to calculate the sample standard deviation, and the other \bar{X} , the average test result for M measurements.

Note that the sample standard deviation, S_x , is simply:

$$S_x = \sqrt{\frac{\sum_{i=1}^N (X_i - \bar{X})^2}{N - 1}} \quad (4.6)$$

In some cases, a particular error distribution may be assumed or known to be a *uniform or rectangular distribution*, Figure 4.2, instead of a normal distribution. For those cases, the sample standard deviation of the data is calculated as:

$$S_x = L / \sqrt{3} \quad (4.7)$$

where L = the plus/minus limits of the uniform distribution for a particular error [3].

For those cases, the standard deviation of the average is written as:

$$S_{\bar{x}} = \frac{L / \sqrt{3}}{\sqrt{M}} \quad (4.8)$$

nonindependent error source. These are sometimes called *dependent error sources*. Their degree of dependence may be estimated with the linear correlation coefficient. If they are nonindependent, whether Type A or Type B, Equation 4.13 becomes [7]:

$$U_{R,ISO} = t_{95} \left\{ \sum_{T=A}^B \sum_{i=1}^{N_{i,T}} \left[(\theta_i U_{i,T})^2 + \sum_{j=1}^{N_{j,T}} \theta_i \theta_j U_{(i,T),(j,T)} (1 - \delta_{i,j}) \right] \right\}^{1/2} \quad (4.14)$$

where: $U_{i,T}$ = the i^{th} elemental uncertainty of Type T (can be Type A or B)
 $U_{R,ISO}$ = the total uncertainty of the measurement or test result
 θ_i = the sensitivity of the test or measurement result to the i^{th} Type T uncertainty
 θ_j = the sensitivity of the test or measurement result to the j^{th} Type T uncertainty
 $U_{(i,T),(j,T)}$ = the covariance of $U_{i,T}$ on $U_{j,T}$

$$= \sum_{l=1}^K U_{i,T}(l) U_{j,T}(l) \quad (4.15)$$

= the sum of the products of the elemental systematic uncertainties that arise from a common source (l)

l = an index or counter for common uncertainty sources

K = the number of common source pairs of uncertainties

$\delta_{i,j}$ = the Kronecker delta. $\delta_{i,j} = 1$ if $i = j$, and $\delta_{i,j} = 0$ if not [7]

T = an index or counter for the ISO uncertainty type, A or B

This ISO classification equation will yield the same total uncertainty as the engineering classification, but the ISO classification does not provide insight into how to improve an experiment's or test's uncertainty. That is, whether to possibly take more data because the random uncertainties are too high or calibrate better because the systematic uncertainties are too large. The engineering classification now presented is therefore the preferred approach.

Engineering Classification

The engineering classification recognizes that experiments and tests have two major types of errors whose limits are estimated with uncertainties at some chosen confidence. These error types may be grouped as *random* and *systematic*. Their corresponding limit estimators are the random uncertainty and systematic uncertainties, respectively.

Random

The general expression for random uncertainty is the ($1S_{\bar{X}}$) standard deviation of the average [6]:

$$S_{\bar{X},R} = \left[\sum_{T=A}^B \sum_{i=1}^{N_{i,T}} (\theta_i S_{\bar{X}_{i,T}})^2 \right]^{1/2} = \left[\sum_{T=A}^B \sum_{i=1}^{N_{i,T}} (\theta_i S_{X_{i,T}} / \sqrt{M_{i,T}})^2 \right]^{1/2} \quad (4.16)$$

where: $S_{X_{i,T}}$ = the sample standard deviation of the i^{th} random error source of Type T
 $S_{\bar{X}_{i,T}}$ = the random uncertainty (standard deviation of the average) of the i^{th} parameter random error source of Type T
 $S_{\bar{X},R}$ = the random uncertainty of the measurement or test result
 $N_{i,T}$ = the total number of random uncertainties, Types A and B, combined
 $M_{i,T}$ = the number of data points averaged for the i^{th} error source, Type A or B
 θ_i = the sensitivity of the test or measurement result to the i^{th} random uncertainty

$$U_B = \left[\left(\frac{0.06}{2} \right)^2 + \left(\frac{0.10}{2} \right)^2 \right]^{1/2} = 0.058 \quad (4.28)$$

$$U_{R,ISO} = \pm K \left[(U_A)^2 + (U_B)^2 \right]^{1/2} = \pm K \left[(0.21)^2 + (0.058)^2 \right]^{1/2} \quad (4.29)$$

Here, remember that the 0.21 is the root sum square of the $1S_{\bar{x}}$ Type A uncertainties in Table 4.2, and 0.058 that for the $1S_{\bar{x}}$ Type B uncertainties. Also note that in most cases, the Type B uncertainties have infinite degrees of freedom and represent an equivalent $2S_{\bar{x}}$. That is why they are divided by 2 — to get an equivalent $1S_{\bar{x}}$. Where there are less than 30 degrees of freedom, one needs to divide by the appropriate Student's t that gave the 95% confidence interval. For the reference junction systematic uncertainty above, that was 2.18.

If “ K ” is taken as Student's t_{95} , the degrees of freedom must first be calculated. Remember that all the systematic components of Type “B” have infinite degrees of freedom except for the 0.07, which has 12 degrees of freedom. Also, all the B_i in Table 4.1 represent an equivalent $2S_{\bar{x}}$ except for 0.07, which represents $2.18S_{\bar{x}}$, as its degrees of freedom are 12 and not infinity. To use their data here, divide them all but the 0.07 by 2 and the 0.07 by 2.18 so they all now represent $1S_{\bar{x}}$, as do the random components. All Type A uncertainties, whether systematic or random in Table 4.1, have degrees of freedom as noted in the table. The degrees of freedom for U_{ISO} is then:

$$df_R = \nu_R = \frac{\left[(0.095)^2 + (0.045)^2 + (0.173)^2 + (0.06/2)^2 + (0.07/2.18)^2 + (0.10/2)^2 \right]^2}{\frac{(0.095)^4}{9} + \frac{(0.045)^4}{4} + \frac{(0.173)^4}{11} + \frac{(0.06)^4}{\infty} + \frac{(0.07/2.18)^4}{12} + \frac{(0.10/2)^4}{\infty}} = 22.51 \approx 22 \quad (4.30)$$

t_{95} is therefore 2.07. $U_{R,ISO}$ is then:

$$U_{R,ISO} = \pm 2.07 \left[(0.21)^2 + (0.058)^2 \right]^{1/2} = 0.45 \text{ for 95\% confidence} \quad (4.31)$$

For a detailed comparison to the engineering system, here denoted as the $U_{R,ENG}$ model, three significant figures are carried so as not to be affected by round-off errors. Then:

$$U_{R,ISO} = \pm 2.074 \left[(0.205)^2 + (0.0583)^2 \right]^{1/2} = 0.442 \text{ for 95\% confidence} \quad (4.32)$$

For the engineering system, $U_{R,ENG}$ model, Equation 4.23, one obtains the expression:

$$U_{R,ENG} = \pm t_{95} \left[(0.13/2)^2 + (0.20)^2 \right]^{1/2} \quad (4.33)$$

Here, the $(0.13/2)$ is the $B_r/2$ and the 0.20 is as before the random component. To obtain the proper t_{95} , the degrees of freedom need to be calculated just as in Equation 4.30. There, the degrees of freedom were 22 and t_{95} equals 2.07. $U_{R,ENG}$ is then:

to serve the parochial needs of commerce, trade, land division, and taxation. Because the standards were defined by local or regional authorities, differences arose that often caused problems in commerce and early scientific investigation. The rapid growth of science in the late 17th century highlighted a number of serious deficiencies in the system of units then in use and, in 1790, led the French National Assembly to direct the French Academy of Sciences to “deduce an invariable standard for all measures and all the weights.” The Academy proposed a system of units, the metric system, to define the unit of length in terms of the earth’s circumference, with the units of volume and mass being derived from the unit of length. Additionally, they proposed that all multiples of each unit be a multiple of 10.

In 1875, the U.S. and 16 other countries signed the “Treaty of the Meter,” establishing a common set of units of measure. It also established an International Bureau of Weights and Measures (called the BIPM). That bureau is located in the Parisian suburb of Sèvres. It serves as the worldwide repository of all the units that maintain our complex international system of weights and measures. It is through this system that compatibility between measurements made thousands of miles apart is currently maintained.

The system of units set up by the BIPM is based on the meter and kilogram instead of the yard and the pound. It is called the *Système International d’Unités* (SI) or the International System of Units. It is used in almost all scientific work in the U.S. and is the only system of measurement units in most countries of the world today.

Even a common system of units does not guarantee measurement agreement, however. Therein lies the crux of the problem. We must make measurements, and we must know how accurately (or, to be more correct, with what uncertainty) we made those measurements. In order to know that, there must be standards. Even more important, everyone must agree on the values of those standards and use the same standards.

As the level of scientific sophistication improved, the basis for the measurement system changed dramatically. The earliest standards were based on the human body, and then attempts were made to base them on “natural” phenomena. At one time, the basis for length was supposed to be a fraction of the circumference of the earth but it was “maintained” by the use of a platinum/iridium bar. Time was maintained by a pendulum clock but was defined as a fraction of the day and so on. Today, the meter is no longer defined by an artifact. Now, the meter is the distance that light travels in an exactly defined fraction of a second. Since the speed of light in a vacuum is now defined as a constant of nature with a specified numerical value (299, 792, 458 m/s), the definition of the unit of length is no longer independent of the definition of the unit of time.

Prior to 1960, the second was defined as 1/86,400th of a mean solar day. Between 1960 and 1967, the second was defined in terms of the unit of time implicit in the calculation of the ephemerides: “The second is the fraction 1/31, 556, 925.9747 of the tropical year for January 0 at 12 hours of ephemeris time.” With the advent of crystal oscillators and, later, atomic clocks, better ways were found of defining the second. This, in turn, allowed a better understanding of things about natural phenomena that would not have been possible before. For example, it is now known that the earth does not rotate on its axis in a uniform manner. In fact, it is erratically slowing down. Since the second is maintained by atomic clocks it is necessary to add “leap seconds” periodically so that the solar day does not gradually change with respect to the time used every day. It was decided that a constant frequency standard was preferred over a constant length of the day.

5.2 What Are Standards?

One problem with standards is that there are several kinds. In addition to “measurement standards,” there are “standards of practice or protocol standards” that are produced by the various standards bodies such as the International Organization for Standardization (ISO), the International Electrotechnical Commission (IEC), the American National Standards Institute (ANSI), and the Standards Council of Canada (SCC). See Figure 5.1.

5.5 Types of Standards

Basic or Fundamental Standards

In the SI system, there are seven basic measurement units from which all other units are derived. All of the units except one are defined in terms of their unitary value. The one exception is the unit of mass. It is defined as 1000 grams (g) or 1 kilogram (kg). It is also unique in that it is the only unit currently based on an artifact. The U.S. kilogram and hence all other standards of mass are based on one particular platinum/iridium cylinder kept at the BIPM in France. If that International Prototype Kilogram were to change, all other mass standards throughout the world would be wrong.

The seven basic units are listed in Appendix 1, Table 1. Their definitions are listed in Appendix 1, Table 2.

Derived Standards

All of the other units are derived from the seven basic units described in Appendix 1, Table 1. Measurement standards are devices that represent the SI standard unit in a measurement. (For example, one might use a zener diode together with a reference amplifier and a power source to supply a known voltage to calibrate a digital voltmeter. This could serve as a measurement standard for voltage and be used as a reference in a measurement.)

Appendix 1, Table 3 lists the most common derived SI units, together with the base units that are used to define the derived unit. For example, the unit of frequency is the hertz; it is defined as the reciprocal of time. That is, 1 hertz (1 Hz) is one cycle per second.

The Measurement Assurance System

Figure 5.3 illustrates the interrelationship of the various categories of standards throughout the world. While it gives more detail to U.S. structure, similar structures exist in other nations. Indeed, a variety of regional organizations exist that help relate measurements made in different parts of the world to each other.

5.6 Numbers, Dimensions, and Units

A measurement is always expressed as a multiple (or submultiple) of some unit quantity. That is, both a numeric value and a unit are required. If electric current were the measured quantity, it might be expressed as some number of milliamperes or even microamperes. It is easy to take for granted the existence of the units used, because their names form an indispensable part of the vocabulary.

5.7 Multiplication Factors

Since it is inconvenient to use whole units in many cases, a set of multiplication factors has been defined that can be used in conjunction with the units to bring a value being measured to a more reasonable size. It would be difficult to have to refer to large distances in terms of the meter; thus, one defines longer distances in terms of kilometers. Short distances are stated in terms of millimeters, micrometers, nanometers, etc. See Appendix 1, Table 4.

Defining Terms

Most of the definitions in this listing were taken from the *International Vocabulary of Basic and General Terms in Metrology*, published by the ISO, 1993 (VIM) [7]. They are indicated by the inclusion (in brackets) of their number designation in the VIM. The remainder of the definitions are not intended to

Reference Material [6.13]: A material or substance, one or more of whose property values are sufficiently homogeneous and well established to be used for the calibration of an apparatus, the assessment of a measurement method, or for assigning values to materials.

NOTE: A reference material can be in the form of a pure or mixed gas, liquid or solid. Examples are water for the calibration of viscometers, sapphire as a heat-capacity calibrant in calorimetry, and solutions used for calibration in chemical analysis.

This definition, including the Note, is taken from ISO Guide 30:1992.

Repeatability (of results of measurements) [3.6]: The closeness of the agreement between the results of successive measurements of the same measurand carried out under the same conditions of measurement.

NOTES:

1. These conditions are called *repeatability conditions*.
2. Repeatability conditions include:
 - a. The same measurement process
 - b. The same observer
 - c. The same measuring instrument, used under the same conditions
 - d. The same location
 - e. Repetition over a short period of time
3. Repeatability can be expressed quantitatively in terms of the dispersion of characteristics of the results.

Reproducibility (of results of measurements) [3.7]: The closeness of the agreement between the results of measurements of the same measurand carried out under changed conditions of measurement.

NOTES:

1. A valid statement of reproducibility requires specification of the conditions changed.
2. The changed conditions include:
 - a. Principle of measurement
 - b. Method of measurement
 - c. Observer
 - d. Measuring instrument
 - e. Reference standard
 - f. Location
 - g. Condition of use
 - h. Time
3. Reproducibility can be expressed quantitatively in terms of the dispersion characteristics of the results.
4. Results here are usually understood to be corrected results.

Secondary Standard [6.5]: A standard whose value is assigned by comparison with a primary standard of the same quantity.

Standards Laboratory: A work space, provided with equipment and standards, a properly controlled environment, and trained personnel, established for the purpose of maintaining traceability of standards and measuring equipment used by the organization it supports. Standards laboratories typically perform fewer, more specialized and higher accuracy measurements than Calibration Laboratories.

Tolerance: In metrology, the limits of the range of values (the uncertainty) that apply to a properly functioning measuring instrument.

- 15 Tilt Measurement** *Adam Chrzanowski, James M. Secord*..... 15-1
Tiltmeters or Inclinometers • Geodetic Leveling • Hydrostatic Leveling • Suspended and Inverted Plumb Lines • Integration of Observations
- 16 Velocity Measurement** *Charles P. Pinney, William E. Baker* 16-1
Introduction • Measurement of Linear Velocity • Velocity: Angular • Conclusion
- 17 Acceleration, Vibration, and Shock Measurement** *Halit Eren* 17-1
Accelerometer Dynamics: Frequency Response, Damping, Damping Ratio, and Linearity • Electromechanical Force-Balance (Servo) Accelerometers • Piezoelectric Accelerometers • Piezoresistive Accelerometers • Differential-Capacitance Accelerometers • Strain-Gage Accelerometers • Seismic Accelerometers • Inertial Types, Cantilever, and Suspended-Mass Configuration • Electrostatic Force Feedback Accelerometers • Microaccelerometers • Cross-Axis Sensitivity • Selection, Full-Scale Range, and Overload Capability • Signal Conditioning

TABLE 6.2 Characteristics of Conductive Plastic, Wirewound, and Hybrid Resistive Elements

	Conductive plastic	Wirewound	Hybrid
Resolution	Infinitesimal	Quantized	Infinitesimal
Power rating	Low	High	Low
Temperature stability	Poor	Excellent	Very good
Noise	Very low	Low, but degrades with time	Low
Life	10^6 – 10^8 cycles	10^5 – 10^6 cycles	10^6 – 10^7 cycles

TABLE 6.3 Potentiometer Terminal Markings

Terminal	Possible color codings			Rotary pot	Linear-motion pot
1	Yellow	Red	Black	CCW limit	Fully retracted limit
2	Red	Green	White	Wiper	Wiper
3	Green	Black	Red	CW limit	Fully extended limit

Hybrid elements feature a wirewound core with a conductive plastic coating, combining wirewound and conductive plastic technologies to realize some of the more desirable attributes of both. The plastic limits power dissipation abilities in exchange for low noise, long life, and unlimited resolution. Like wirewounds, hybrids offer excellent temperature stability. They make an excellent choice for precision measurement.

Cermet elements, made from a ceramic-metal alloy, offer unlimited resolution and reasonable noise levels. Their advantages include high power dissipation abilities and excellent stability in adverse conditions. Cermet elements are rarely applied to precision measurement because conductive plastic elements offer lower noise, lower friction, and longer life.

Carbon composition elements, molded under pressure from a carbon-plastic mixture, are inexpensive and very popular for general use, but not for precision measurement. They offer unlimited resolution and low noise, but are sensitive to environmental stresses (e.g., temperature, humidity) and are subject to wear.

Table 6.2 summarizes the distinguishing characteristics of the preferred resistive elements for precision measurement.

Electrical Characteristics

Before selecting a pot and integrating it into a measurement system, the following electrical characteristics should be considered.

Terminals and Taps

Table 6.3 shows the conventional markings found on the pot housing [4, 5]; CW and CCW indicate clockwise and counter-clockwise rotation as seen from the front end. Soldering studs and eyelets, integral connectors, and flying leads are common means for electrical connection. In addition to the wiper and end terminals, a pot may possess one or more terminals for *taps*. A tap enables an electrical connection to be made with a particular point along the resistive element. Sometimes, a *shunt resistor* is connected to a tap in order to modify the output function. End terminations and taps can exhibit different electrical characteristics depending on how they are manufactured. See [2] for more details.

Taper

Pots are available in a variety of different tapers that determine the shape of the output function. With a linear-taper pot, the output varies linearly with wiper motion, as shown in Figure 6.2. (Note that a pot with a linear taper should not be confused with a linear-motion pot, which is sometimes called a “linear pot.”) Linear-taper pots are the most commonly available, and are widely used in sensing and control

to as high as 150°C. Operating outside specified limits can cause material failure, either directly from temperature or from thermally induced misalignment.

Vibration, Shock, and Acceleration

Vibration, shock, and acceleration are all potential sources of contact discontinuities between the wiper and the resistive element. In general, a contact failure is considered to be a discontinuity equal to or greater than 0.1 ms [2]. The values quoted in specification sheets are in gs and depend greatly on the particular laboratory test. Some characterization tests use sinusoidal vibration, random vibration, sinusoidal shock, sawtooth shock, or acceleration to excite the pot. Manufacturers use mechanical design strategies to eliminate weaknesses in a pot's dynamic response. For example, one technique minimizes vibration-induced contact discontinuities using multiple wipers of differing resonant frequencies.

Speed

Exceeding a pot's specified maximum speed can cause premature wear or discontinuous values through effects such as wiper bounce. As a general rule, the slower the shaft motion, the longer the unit will last (in total number of cycles). Speed limitations depend on the materials involved. For rotary pots, wirewound models have preferred maximum speeds on the order of 100 rpm, while conductive plastic models have allowable speeds as high as 2000 rpm. Linear-motion pots have preferred maximum velocities up to 10 m s⁻¹.

Life

Despite constant mechanical wear, a pot's expected lifetime is on the order of a million cycles when used under proper conditions. A quality film pot can last into the hundreds of millions of cycles. Of wirewound, hybrid, and conductive plastic pots, the uneven surface of a wirewound resistive element inherently experiences the most wear and thus has the shortest expected operating life. Hybrids improve on this by using a wirewound construction in combination with a smooth conductive film coating. Conductive plastic pots generally have the longest life expectancy due to the smooth surface of their resistive element.

Contamination and Seals

Foreign material contaminating pots can promote wear and increase friction between the wiper and the resistive element. Consequences range from increased mechanical loading to outright failure (e.g., seizing, contact discontinuity). Fortunately, sealed pots are available from most manufacturers for industrial applications where dirt and liquids are often unavoidable. To aid selection, specifications often include the type of *case sealing* (i.e., mechanisms and materials) and the *seal resistance* to cleaning solvents and other commonly encountered fluids.

Misalignment

Shaft misalignment in a pot can prematurely wear its bearing surfaces and increase its mechanical loading effects. A good design minimizes misalignment. (See *Implementation*, below.) Manufacturers list a number of alignment tolerances. In linear-motion pots, *shaft misalignment* is the maximum amount a shaft can deviate from its axis. The degree to which a shaft can rotate around its axis is listed under *shaft rotation*. In rotary pots, *shaft end play* and *shaft radial play* both describe the amount of shaft deflection due to a radial load. *Shaft runout* denotes the shaft diameter eccentricity when a shaft is rotated under a radial load.

Mechanical Mounting Methods

Hardware features on a pot's housing determine the mounting method. Options vary with manufacturer, and among rotary, linear-motion, and string pots. Offerings include custom bases, holes, tabs, flanges, and brackets — all of which secure with machine screws — and threaded studs, which secure with nuts. Linear-motion pots are available with rod or slider actuation, some with internal or external return springs. Mounting is typically accomplished by movable clamps, often supplied by the pot manufacturer. Other linear-motion pots mount via a threaded housing. For rotary pots, the two most popular mounting methods are the *bushing mount* and the *servo mount*. See Figure 6.5.

In a typical measurement application, the pot shaft couples to a mechanical component (e.g., a gear, a pulley), or to another shaft of the same or different diameter. Successful couplings provide a positive link to the shaft without stressing the pot's mechanics. Satisfying these objectives with rotary and linear-motion pots requires a balance between careful alignment and compliant couplings. Alignment is not as critical with a string pot. Useful coupling methods include the following.

Compliant couplings. It is generally wise to put a compliant coupling between a pot's shaft and any other shafting. A compliant coupling joins two misaligned shafts of the same or different diameter. Offerings from the companies in Table 6.4 include bellows couplings, flex couplings, spring couplings, spider couplings, Oldham couplings, wafer spring couplings, flexible shafts, and universal joints. Each type has idiosyncrasies that impact measurement error; manufacturer catalogs provide details.

Sleeve couplings. Less expensive than a compliant coupling, a rigid sleeve coupling joins two shafts of the same or different diameter with the requirement that the shafts be perfectly aligned. Perfect alignment is difficult to achieve initially, and impossible to maintain as the system ages. Imperfect alignment accelerates wear and risks damaging the pot. Sleeve couplings are available from the companies listed in Table 6.4.

Press fits. A press fit is particularly convenient when the bore of a small plastic part is nominally the same as the shaft diameter. Carefully force the part onto the shaft. Friction holds the part in place, but repeated reassembly will compromise the fit.

Shrink fits. Components with a bore slightly under the shaft diameter can be heated to expand sufficiently to slip over the shaft. A firm grip results as the part cools and the bore contracts.

Pinning. Small hubbed components can be pinned to a shaft. The pin should extend through the hub partway into the shaft, and the component should fit on the shaft without play. Use roll pins or spiral pins combined with a thread-locking compound (e.g., Loctite 242).

Set-screws. Small components are available with hubs that secure with set-screws. The component should fit on the shaft without play. For best results, use two set-screws against a shaft with perpendicular flats. Dimple a plain shaft using the component's screw hole(s) as a drill guide. Apply a thread-locking compound (e.g., Loctite 242) to prevent the set-screws from working loose.

Clamping. Small components are also available with split hubs that grip a shaft when squeezed by a matching hub clamp. Clamping results in a secure fit without marring the shaft.

Adhesives. Retaining compounds (e.g., Loctite 609) can secure small components to a shaft. Follow manufacturer's instructions for best results.

Spring-loaded contact. A spring-loaded shaft will maintain positive contact against a surface that moves at reasonable speeds and without sudden acceleration.

Costs and Sources

Precision pots are inexpensive compared to other displacement measurement technologies. Table 6.5 lists approximate costs for off-the-shelf units in single quantity. Higher quality generally commands a higher price; however, excellent pots are often available at bargain prices due to volume production or surplus conditions. Electronic supply houses offer low-cost pots (i.e., under \$20) that can suffice for short-term projects. Regardless of price, always check the manufacturer's specifications to confirm a pot's suitability for a given application.

Table 6.6 lists several sources of precision pots. Most manufacturers publish catalogs, and many have Web sites. In addition to a standard product line, most manufacturers will custom-build pots for high-volume applications.

TABLE 6.5 Typical Single-quantity Prices (\$US) for Commercially Available Pots

Potentiometer type	Approximate price range
Rotary	\$10–\$350
Linear-motion	\$20–\$2000
String	\$250–\$1000

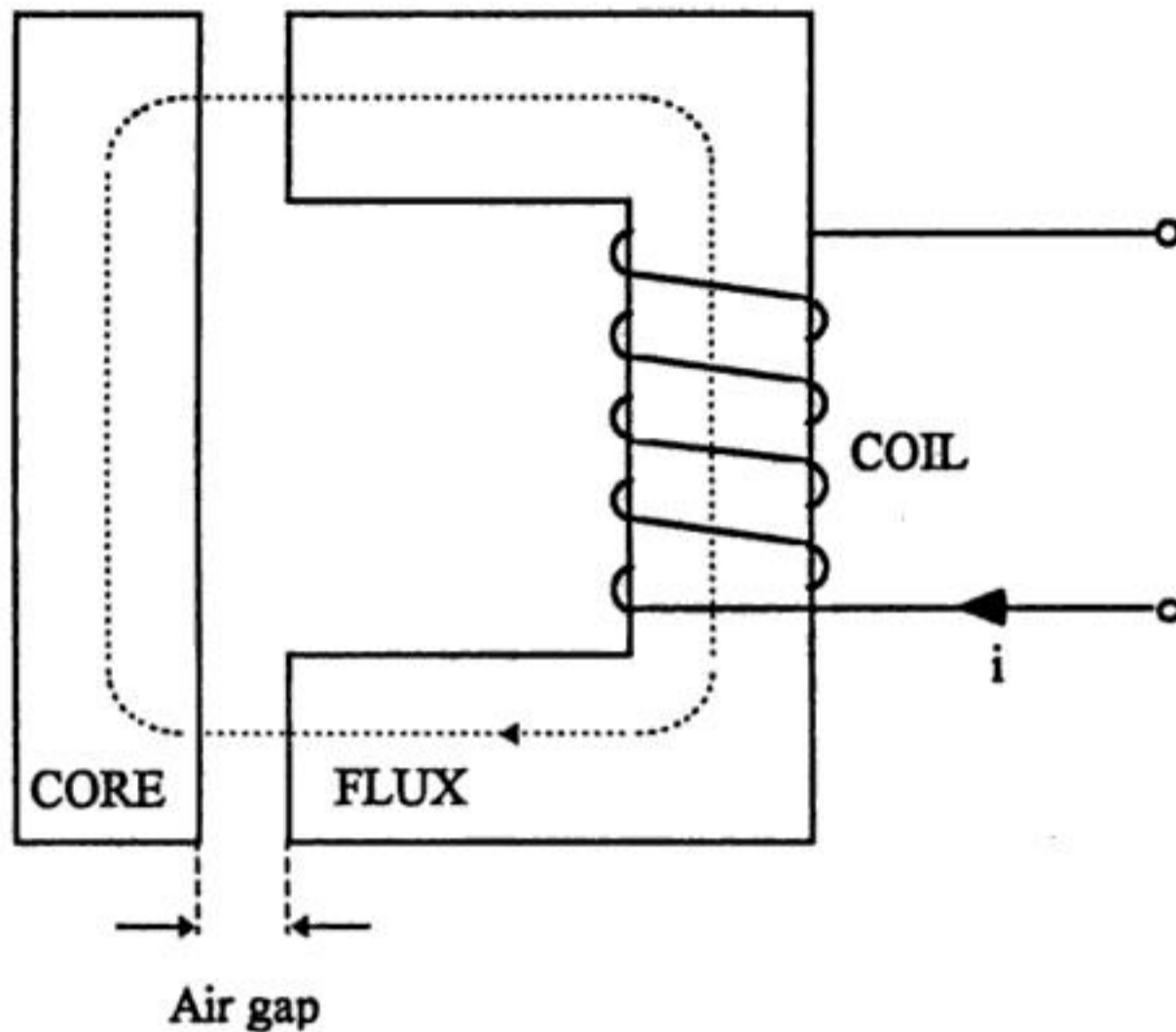


FIGURE 6.7 A basic inductive sensor consists of a magnetic circuit made from a ferromagnetic core with a coil wound on it. The coil acts as a source of magnetomotive force (mmf) that drives the flux through the magnetic circuit and the air gap. The presence of the air gap causes a large increase in circuit reluctance and a corresponding decrease in the flux. Hence, a small variation in the air gap results in a measurable change in inductance.

$$\Psi = n\Phi = n^2 i / \mathcal{R} \quad \text{weber} \quad (6.4)$$

Equation 6.4 leads to self inductance L of the coil, which is described as the total flux (Ψ weber) per unit current for that particular coil; that is

$$L = \Psi / I = n^2 / \mathcal{R} \quad (6.5)$$

This indicates that the self inductance of an inductive element can be calculated by magnetic circuit properties. Expressing \mathcal{R} in terms of dimensions as:

$$\mathcal{R} = l / \mu \mu_0 A \quad (6.6)$$

where l = the total length of the flux path

μ = the relative permeability of the magnetic circuit material

μ_0 = the permeability of free space ($= 4\pi \times 10^{-7}$ H/m)

A = the cross-sectional area of the flux path

The arrangement illustrated in Figure 6.7 becomes a basic inductive sensor if the air gap is allowed to vary. In this case, the ferromagnetic core is separated into two parts by the air gap. The total reluctance of the circuit now is the addition of the reluctance of core and the reluctance of air gap. The relative permeability of air is close to unity, and the relative permeability of the ferromagnetic material is of the order of a few thousand, indicating that the presence of the air gap causes a large increase in circuit reluctance and a corresponding decrease in the flux. Hence, a small variation in the air gap causes a measurable change in inductance. Most of the inductive transducers are based on these principles and are discussed below in greater detail.

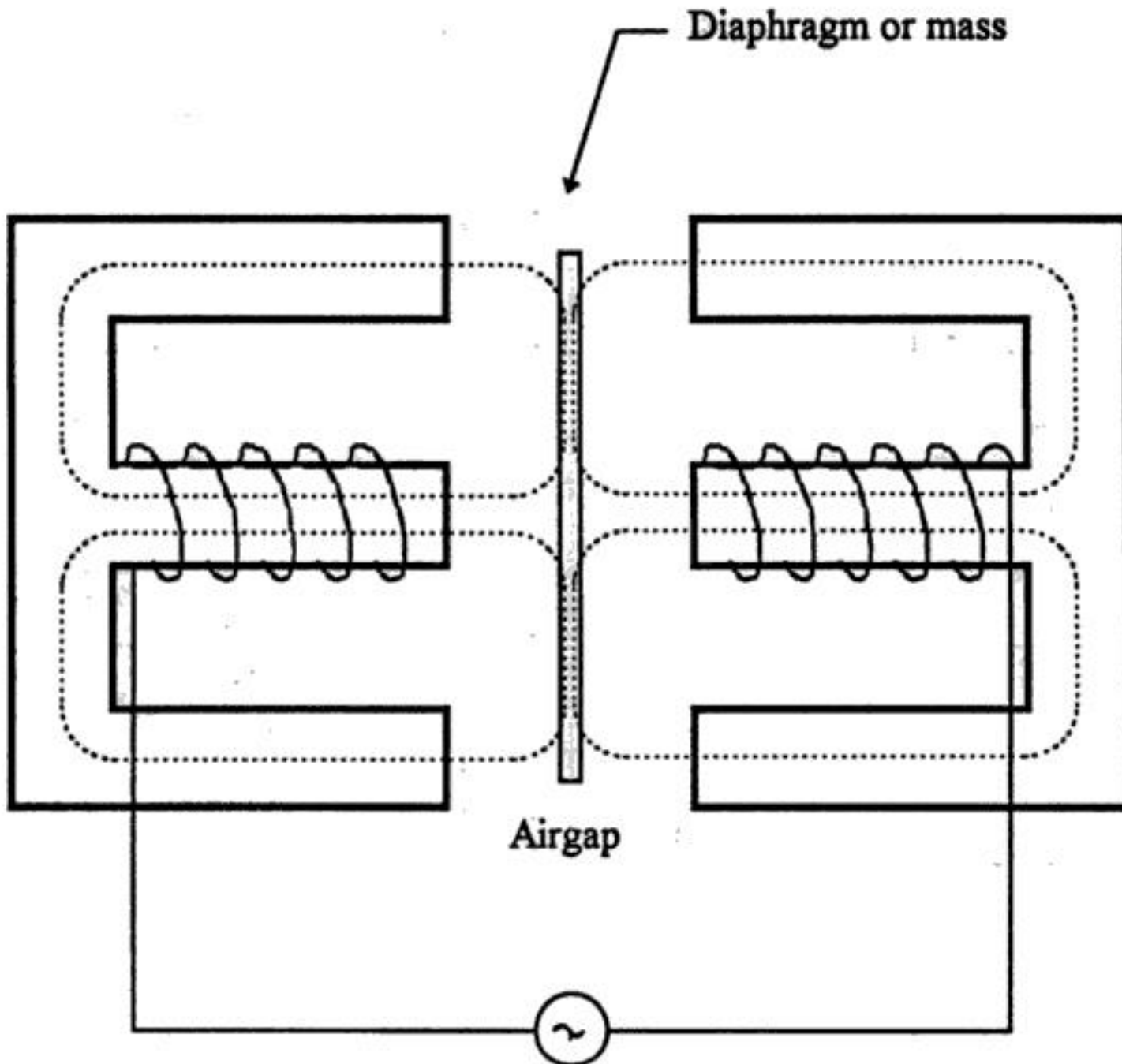


FIGURE 6.10 A typical commercial variable differential sensor. The iron core is located half-way between the two E frames. Motion of the core increases the air gap for one of the E frames while decreasing the other side. This causes reluctances to change, thus inducing more voltage on one side than the other. Motion in the other direction reverses the action, with a 180° phase shift occurring at null. The output voltage can be processed, depending on the requirements, by means of rectification, demodulation, or filtering. The full-scale motion may be extremely small, on the order of few thousandths of a centimeter.

increases as the tooth moves away from the pole. When the wheel rotates with a velocity ω , the flux may mathematically be expressed as:

$$\Psi(\theta) = A + B \cos m\theta \quad (6.12)$$

where A = the mean flux
 B = the amplitude of the flux variation
 m = the number of teeth

The induced emf is given by:

$$E = -d\Psi(\theta)/dt = -(d\Psi(\theta)/d\theta) \times (d\theta/dt) \quad (6.13)$$

or

$$E = bm\omega \sin m\omega t \quad (6.14)$$

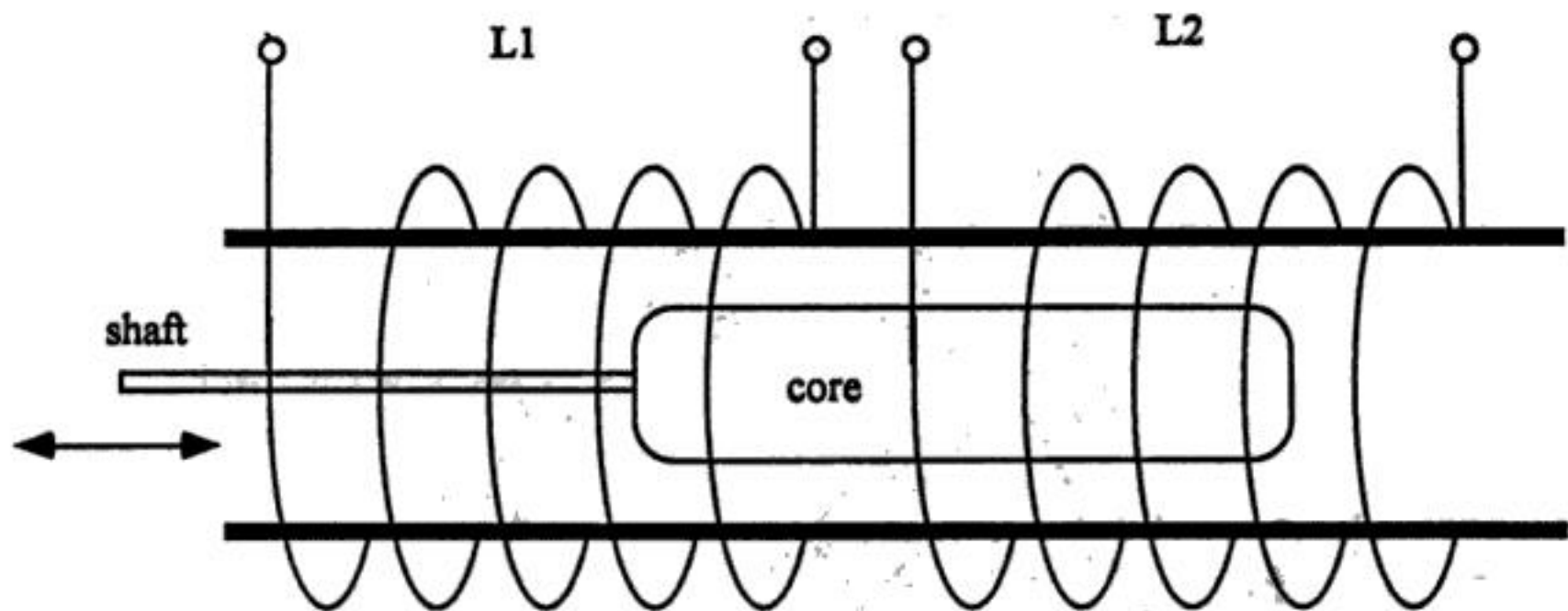


FIGURE 6.14 A typical linear-variable inductor consists of a movable iron core inside a former holding a center-tapped coil. The core and both coils have the same length l . When the core is in the reference position, each coil will have equal inductances of value L . As the core moves by δl , changes in inductances $+\delta L$ and $-\delta L$ create voltage outputs from the coils.

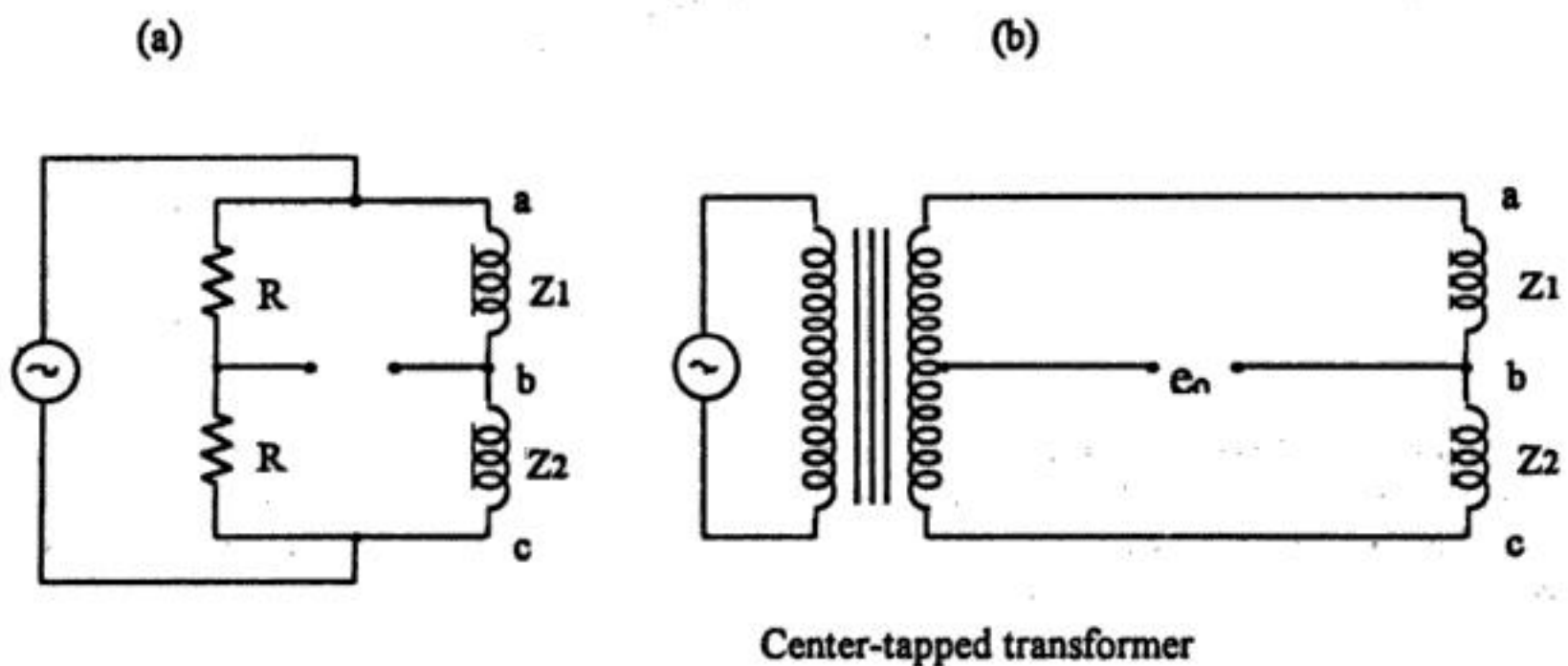


FIGURE 6.15 The two coils of a linear-variable inductor are usually placed to form two arms of a bridge circuit, also having two equal balancing resistors as in circuit (a). The bridge is excited with ac of 5 V to 25 V with a frequency of 50 Hz to 5 kHz. At a selected excitation frequency, the total transducer impedance at null conditions is set in the 100 Ω to 1000 Ω range. By careful construction of the bridge, the output voltage can be made a linear function displacement of the core within a limited range. In some cases, in order to reduce power losses due to heating of resistors, center-tapped transformers may be used as a part of the bridge network (b).

It is particularly easy to construct transducers of this type, by simply winding a center-tapped coil on a suitable former. The variable-inductance transducers are commercially available in strokes from about 2 mm to 500 cm. The sensitivity ranges between 1% full scale to 0.02% in long stroke special constructions. These devices are also known as linear displacement transducers or LDTs, and they are available in various shape and sizes.

Apart from linear-variable inductors, there are rotary types available too. Their cores are specially shaped for rotational applications. Their nonlinearity can vary between 0.5% to 1% full scale over a range of 90° rotation. Their sensitivity can be up to 100 mV per degree of rotation.

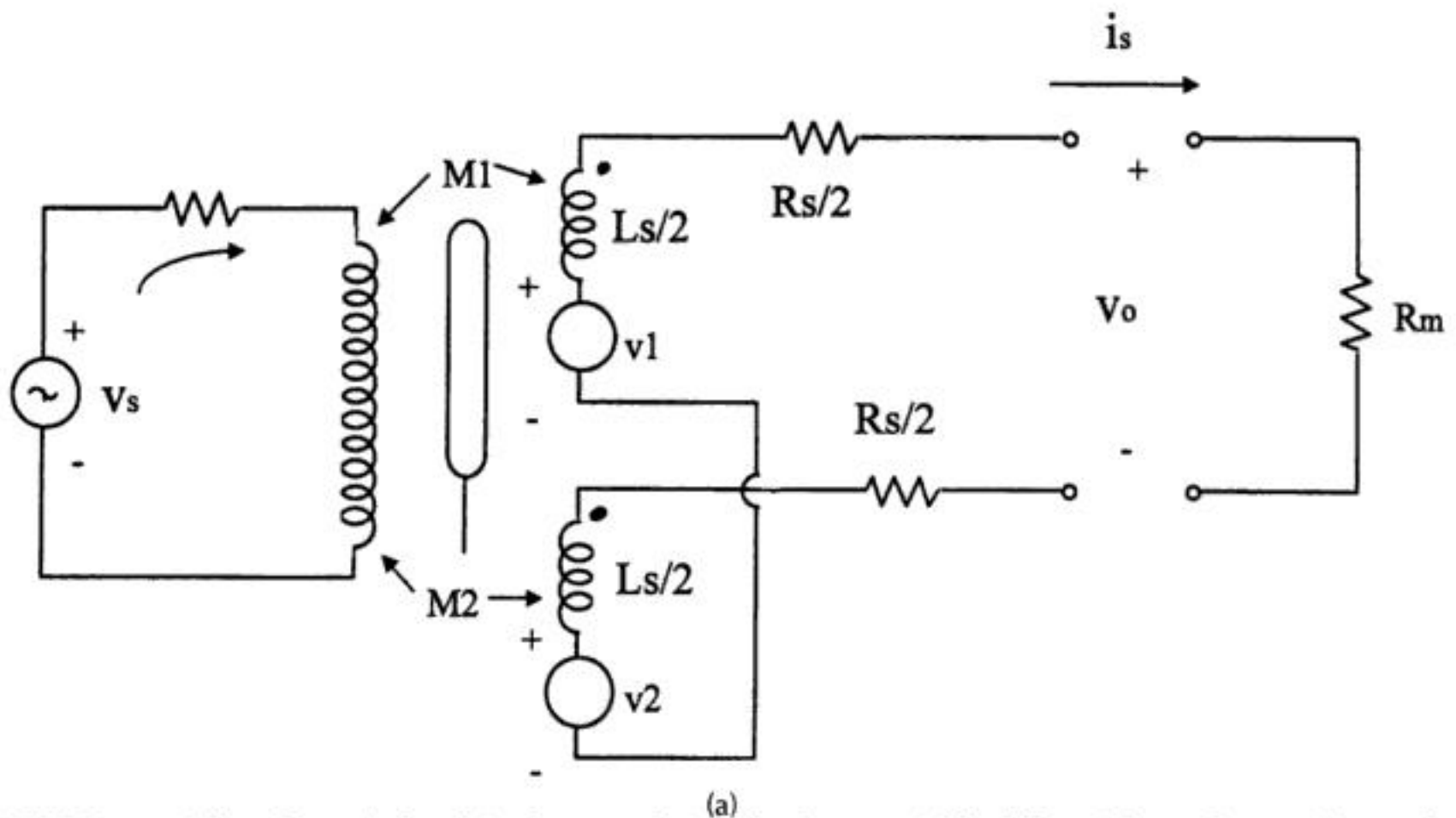


FIGURE 6.18 The voltages induced in the secondaries of a linear-variable differential-transformer (a) may be processed in a number of ways. The output voltages of individual secondaries v_1 and v_2 at null position are illustrated in (b). In this case, the voltages of individual coils are equal and in phase with each other. Sometimes, the outputs are connected opposing each other, and the output waveform v_o becomes a function of core position x and phase angle α , as in (c). Note the phase shift of 180° as the core position changes above and below the null position.

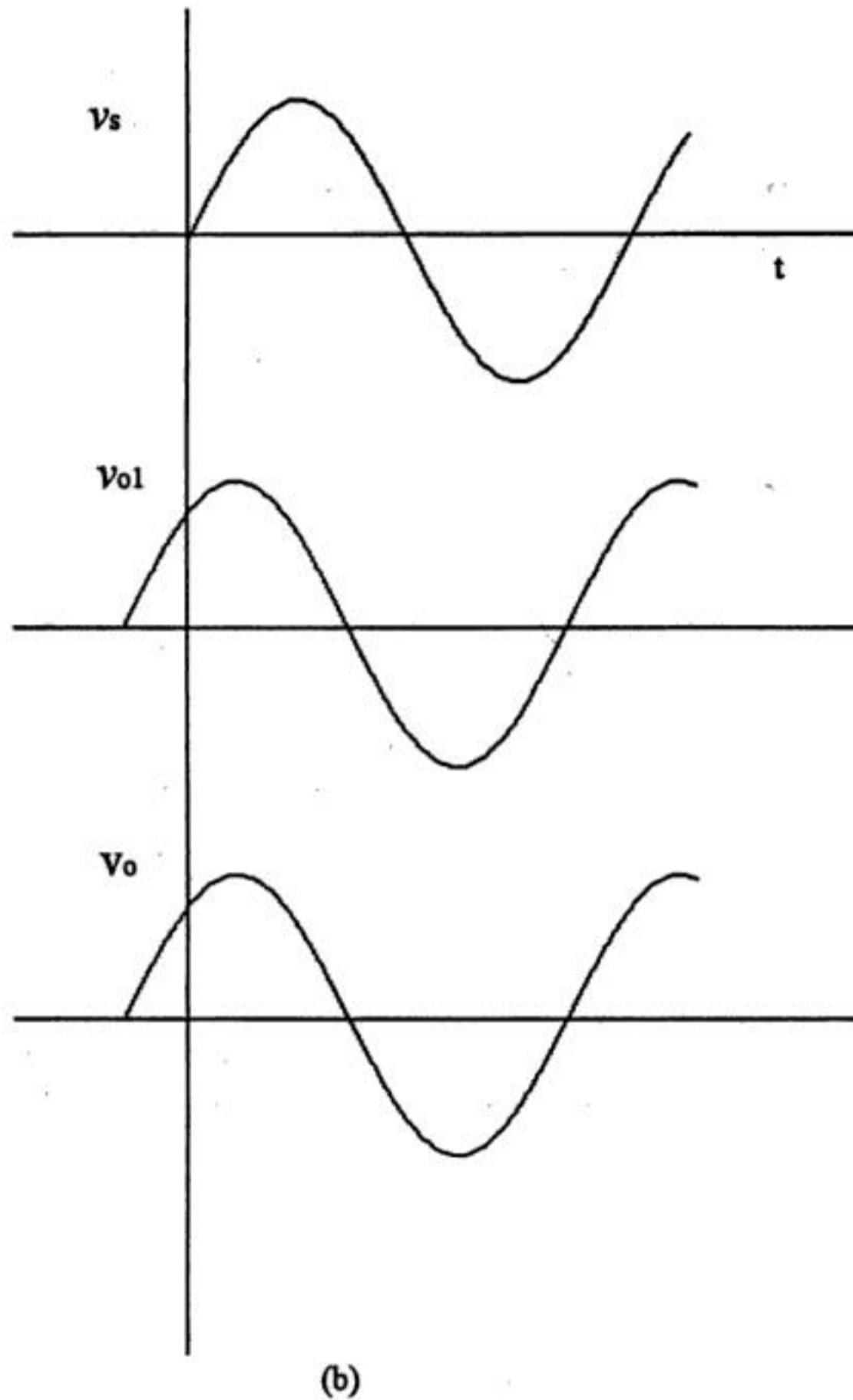
an equal coupling between the primary and secondary windings, thus giving a null point or reference point of the sensor. As long as the core remains near the center of the coil arrangement, output is very linear. The linear ranges of commercial differential transformers are clearly specified, and the devices are seldom used outside this linear range.

The ferromagnetic core or plunger moves freely inside the former, thus altering the mutual inductance between the primary and secondaries. With the core in the center, or at the reference position, the induced emfs in the secondaries are equal; and since they oppose each other, the output voltage is zero. When the core moves, say to the left, from the center, more magnetic flux links with the left-hand coil than with the right-hand coil. The voltage induced in the left-hand coil is therefore larger than the induced voltage on the right-hand coil. The magnitude of the output voltage is then larger than at the null position and is equal to the difference between the two secondary voltages. The net output voltage is in phase with the voltage of the left-hand coil. The output of the device is then an indication of the displacement of the core. Similarly, movement in the opposite direction to the right from the center reverses this effect, and the output voltage is now in phase with the emf of the right-hand coil.

For mathematical analysis of the operation of LVDTs, Figure 6.18(a) can be used. The voltages induced in the secondary coils are dependent on the mutual inductance between the primary and individual secondary coils. Assuming that there is no cross-coupling between the secondaries, the induced voltages may be written as:

$$v_1 = M_1 s i_p \quad \text{and} \quad v_2 = M_2 s i_p \tag{6.15}$$

where M_1 and M_2 are the mutual inductances between primary and secondary coils for a fixed core position; s is the Laplace operator; and i_p is the primary current



(b)
FIGURE 6.18 (continued)

In the case of opposing connection, no load output voltage v_o without any secondary current may be written as:

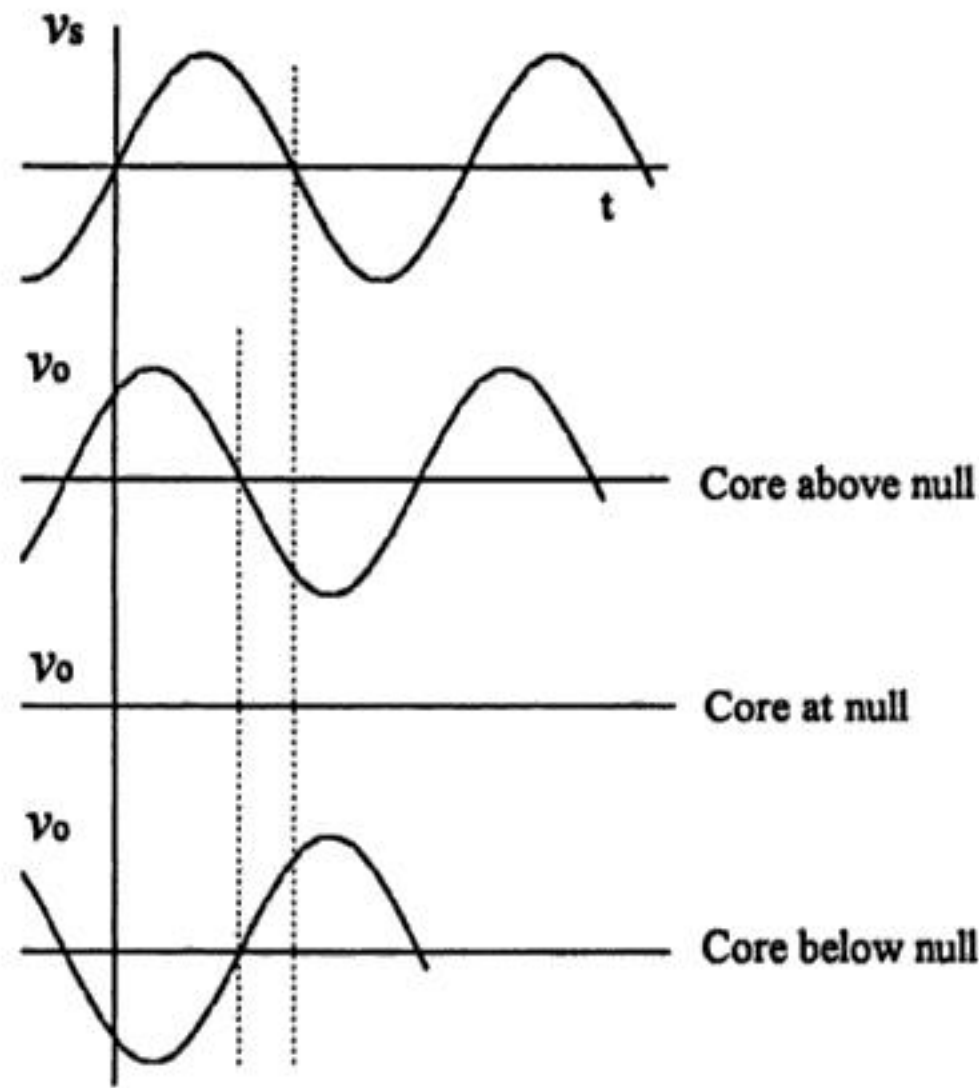
$$v_o = v_1 - v_2 = (M_1 - M_2) s i_p \quad (6.16)$$

writing

$$v_s = i_p (R + sL_p) \quad (6.17)$$

Substituting i_p in Equation 6.15 gives the transfer function of the transducer as:

$$v_o/v_s = (M_1 - M_2) s / (R + sL_p) \quad (6.18)$$



(c)

FIGURE 6.18 (continued)

However, if there is a current due to output signal processing, then describing equations may be modified as:

$$v_o = R_m i_s \tag{6.19}$$

where $i_s = (M_1 - M_2) si_p / (R_s + R_m + sL_s)$
and

$$v_s = i_p (R + sL_p) - (M_1 - M_2) si_s \tag{6.20}$$

Eliminating i_p and i_s from Equations 4.19 and 4.20 results in a transfer function as:

$$v_o/v_s = R_m (M_1 - M_2) s / \left\{ \left[(M_1 - M_2)^2 + L_s L_p \right] s^2 + \left[L_p (R + R_m) + R L_s \right] s + (R_s + R_m) + R \right\} \tag{6.21}$$

This is a second-order system, which indicates that due to the effect of the numerator of Eq. 6.21, the phase angle of the system changes from +90° at low frequencies to -90° at high frequencies. In practical applications, the supply frequency is selected such that at the null position of the core, the phase angle of the system is 0°.

The amplitudes of the output voltages of secondary coils are dependent on the position of the core. These outputs may directly be processed from each individual secondary coil for slow movements of the core, and when the direction of the movement of the core does not bear any importance. However, for fast movements of the core, the signals can be converted to dc and the direction of the movement from the null position can be detected. There are many options to do this; however, a *phase-sensitive demodulator*

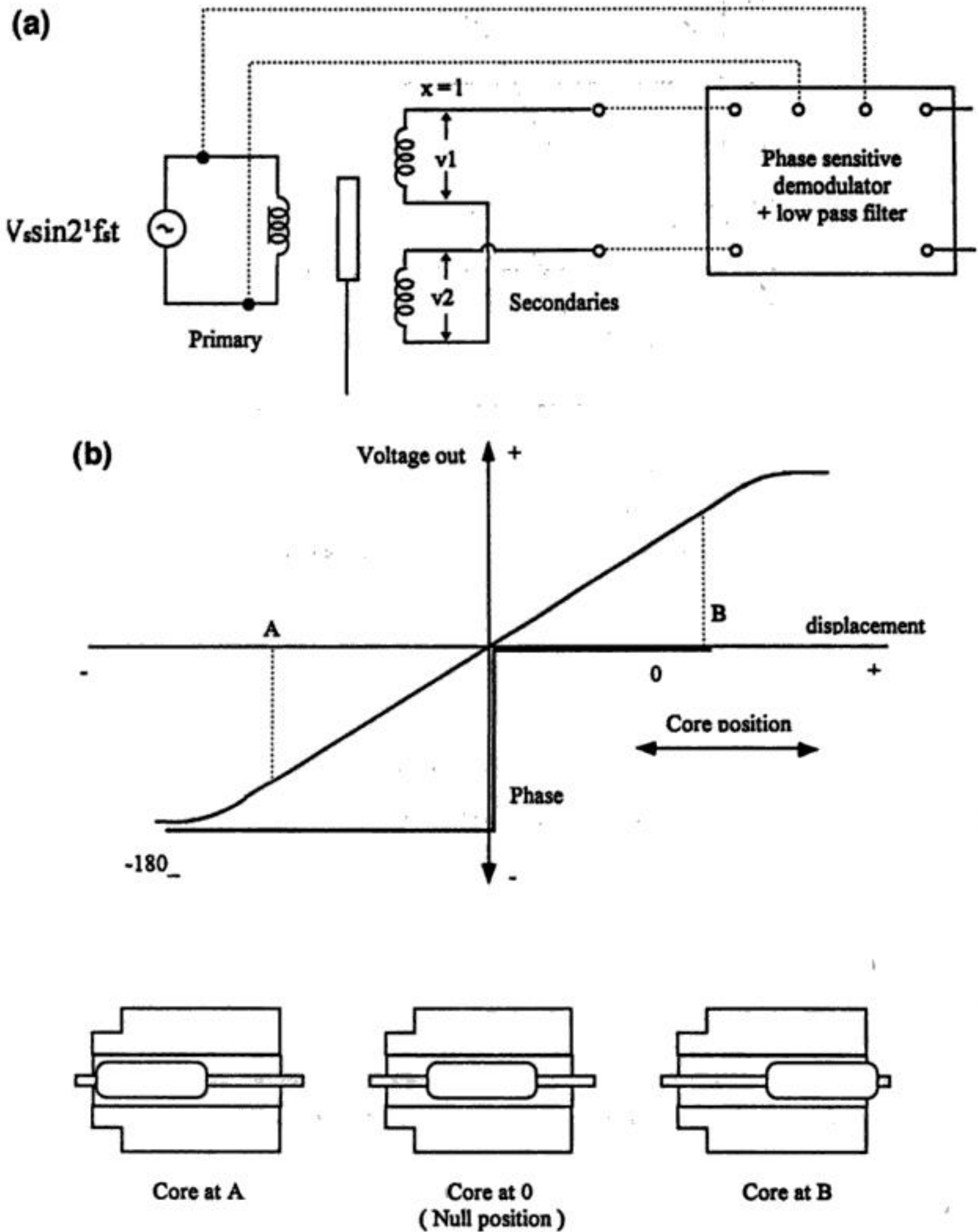


FIGURE 6.19 Phase-sensitive demodulator and (a) are commonly used to obtain displacement proportional signals from LVDTs and other differential type inductive sensors. They convert the ac outputs from the sensors into dc values and also indicate the direction of movement of the core from the null position. A typical output of the phase-sensitive demodulator is shown in (b). The relationship between output voltage v_o and phase angle α is also shown against core position x .

and filter arrangement is commonly used, as shown in Figure 6.19(a). A typical output of the phase-sensitive demodulator is illustrated in Figure 6.19(b), in relation to output voltage v_o , displacement x , and phase angle α .

The phase-sensitive demodulators are used extensively in differential type inductive sensors. They basically convert the ac outputs to dc values and also indicate the direction of movement of the core

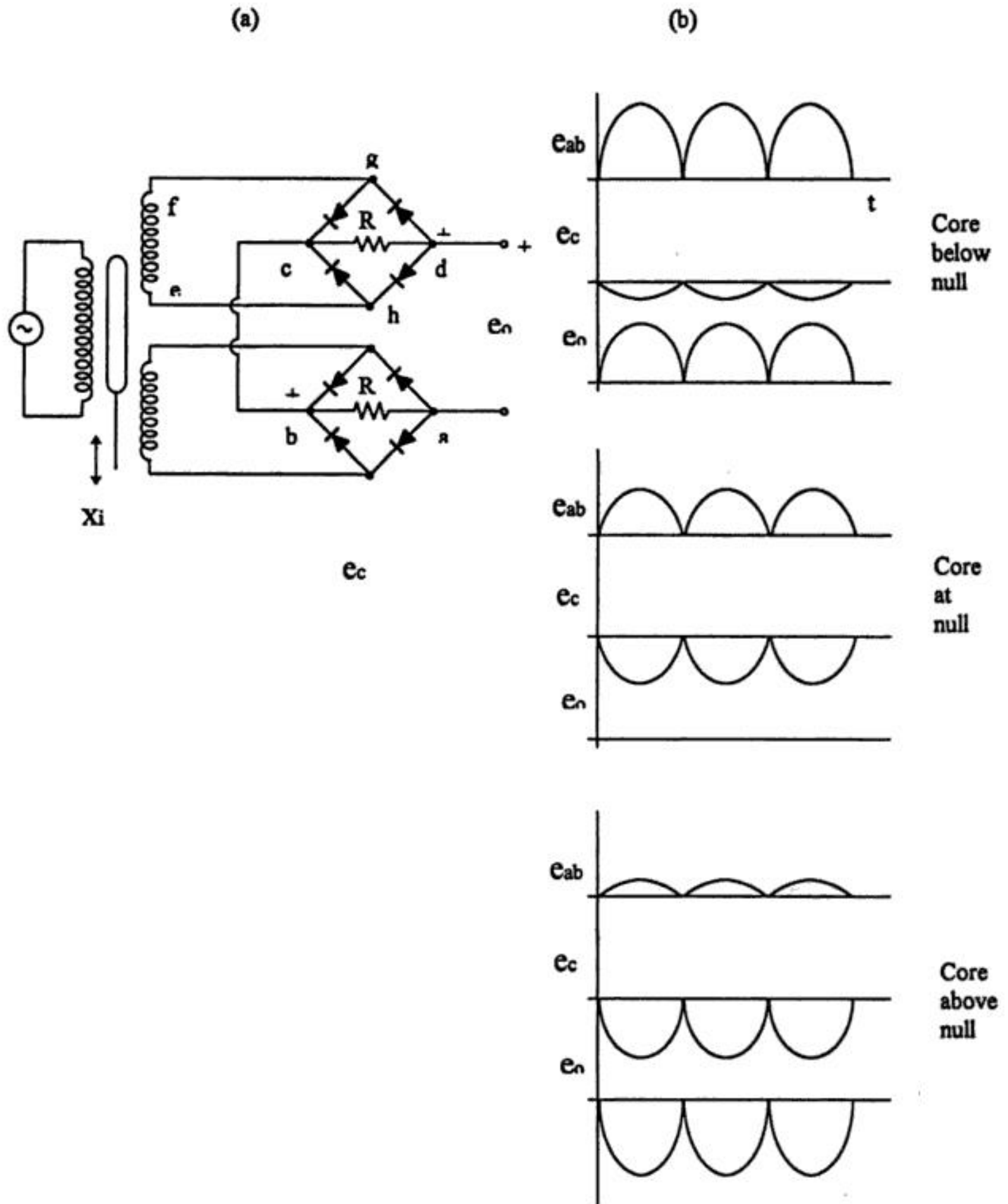


FIGURE 6.20 A typical phase-sensitive demodulation circuit based on diode bridges as in (a). Bridge 1 acts as a rectification circuit for secondary 1, and bridge 2 acts as a rectifier for secondary 2 where the net output voltage is the difference between the two bridges, as in (b). The position of the core can be determined from the amplitude of the dc output, and the direction of the movement of the core can be determined from the polarity of the voltage. For rapid movements of the core, the output of the diode bridges must be filtered, for this, a suitably designed simple RC filter may be sufficient.

from the null position. A typical phase-sensitive demodulation circuit may be constructed, based on diodes shown in Figure 6.20(a). This arrangement is useful for very slow displacements, usually less than 1 or 2 Hz. In this figure, bridge 1 acts as a rectification circuit for secondary 1, and bridge 2 acts as a rectifier for secondary 2. The net output voltage is the difference between the outputs of two bridges, as

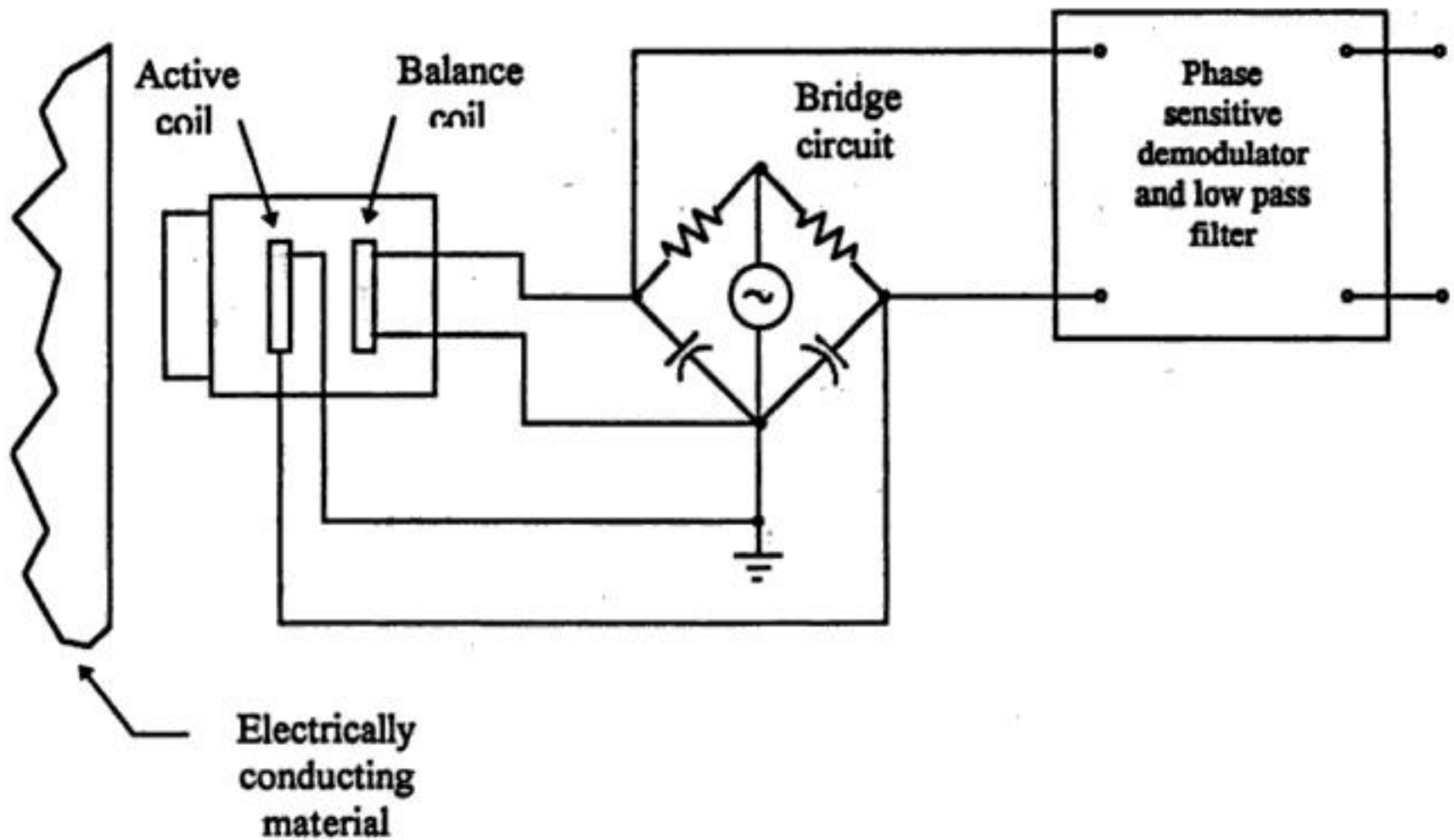


FIGURE 6.22 Eddy current transducers are inductive transducers using probes. The probes contain one active and one balance coil. The active coil responds to the presence of a conducting target, while the balance coil completes a bridge circuit and provides temperature compensation. When the probe is brought close to the target, the flux from the probe links with the target, producing eddy currents within the target that alter the inductance of the active coil. This change in inductance is detected by a bridge circuit.

In most rotary linear-variable differential transformers, the rotor mass is very small, usually less than 5 g. The nonlinearity in the output ranges between $\pm 1\%$ and $\pm 3\%$, depending on the angle of rotation. The motion in the radial direction produces a small output signal that can affect the overall sensitivity. However, this transverse sensitivity is usually kept below 1% of the longitudinal sensitivity.

Eddy Current

Inductive transducers based on eddy currents are mainly probe types, containing two coils as shown in Figure 6.22. One of the coils, known as the active coil, is influenced by the presence of the conducting target. The second coil, known as the balance coil, serves to complete the bridge circuit and provides temperature compensation. The magnetic flux from the active coil passes into the conductive target by means of a probe. When the probe is brought close to the target, the flux from the probe links with the target, producing eddy currents within the target.

The eddy current density is greatest at the target surface and become negligibly small, about three skin depths below the surface. The skin depth depends on the type of material used and the excitation frequency. While thinner targets can be used, a minimum of three skin depths is often necessary to minimize the temperature effects. As the target comes closer to the probe, the eddy currents become stronger, causing the impedance of the active coil to change and altering the balance of the bridge in relation to the target position. This unbalance voltage of the bridge may be demodulated, filtered, and linearized to produce a dc output proportional to target displacement. The bridge oscillation may be as high as 1 MHz. High frequencies allow the use of thin targets and provide good system frequency response.

Probes are commercially available with full-scale diameter ranging from 0.25 to 30 mm with a nonlinearity of 0.5% and a maximum resolution of 0.0001 mm. Targets are usually supplied by the clients, involving noncontact measurements of machine parts. For nonconductive targets, conductive materials of sufficient thickness must be attached to the surface by means of commercially available adhesives. Since the target material, shape, etc. influence the output, it is necessary to calibrate the system statistically

for a specific target. The recommended measuring range of a given probe begins at a standoff distance equal to about 20% of the stated range of the probe. In some cases, a standoff distance of 10% is recommended for which the system is calibrated as standard. A distance greater than 10% of the measuring range can be used as long as the calibrated measuring range is reduced by the same amount.

Flat targets must be of the same diameter as the probe or larger. If the target diameter is smaller than the probe diameter, the output drops considerably, thus becoming unreliable. Curved-surface targets can behave similar to flat surfaces if the diameter exceeds about three or four times the diameter of the probe. In this case, the target essentially becomes an infinite plane. This also allows some cross-axis movement without affecting the system output. Target diameters comparable to the sensor could result in detrimental effects in measurements due to cross-axis movements.

For curved or irregular shaped targets, the system must be calibrated using the exact target that is seen in the operation. This tends to eliminate any errors caused by the curved surfaces during application. However, special multiprobe systems are available for orbital motions of rotating shafts. If the curved (shaft) target is about 10 times larger than the sensor diameter, it acts as an infinite plane and does not need special calibrations. Care must be exercised to deal with electrical runout due to factors such as inhomogeneities in hardness, etc., particularly valid for ferrous targets. However, nonferrous targets are free from electrical runout concerns.

Shielding and Sensitivity of Inductive Sensors to Electromagnetic Interference

Magnetic fields are produced by currents in wires and more strongly by the coils. The fields due to coils are important due to magnetic coupling, particularly when there are two or more coils in the circuit. The magnetic coupling between coils may be controlled by large spacing between coils, orientation of coils, the shape of the coils, and by shielding.

Inductive sensors come in different shape and sizes. While some sensors have closed cores such as toroidal shapes, others have open cores and air gaps between cores and coils. Closed cores can have practically zero external fields, except for small leakage fluxes. Even if the sensors do not have closed cores, most variable inductor sensors have rather limited external fields, due to two neighboring sets of coils connected in opposite directions that minimize the external fields.

Inductive sensors are made from closed conductors. This implies that if the conductor moves in a magnetic field, a current will flow. Alternatively, a magnetic change produces current in a stationary closed conductor. Unless adequate measures are taken, there may be external magnetic fields linking (interference) with the sensor coils, thus producing currents and unwanted responses.

Due to inherent operations, inductive sensors are designed to have high sensitivity to magnetic flux changes. External electromagnetic interference and external fields can seriously affect the performance of the sensors. It is known that moderate magnetic fields are found near power transformers, electrical motors, and power lines. These small fields produce current in the inductive sensors elements. One way of eliminating external effects is accomplished by magnetic shielding of the sensors and by grounding appropriately. In magnetic shielding, one or more shells of high-permeability magnetic materials surround the part to be shielded. Multiple shells may be used to obtain very complete shielding. The ends of each individual shell are separated by insulation so that the shell does not act as a single shorted turn, thus accommodating high current flows. Similarly, in the case of multiple shielding, shells are isolated from each other by proper insulation.

Alternating magnetic fields are also screened by interposing highly conductive metal sheets such as copper or aluminum on the path of the magnetic flux. The eddy currents induced in the shield give a counter mmf that tends to cancel the interfering magnetic field. This type of shielding is particularly effective at high frequencies. Nevertheless, appropriate grounding must still be observed.

In many inductive sensors, stray capacitances can be a problem, especially at the null position of the moving core. If the capacitive effect is greater than a certain value, say 1% of the full-scale output, this effect may be reduced by the use of center-tapped supply and appropriate grounding.

References

1. J. P. Bentley, *Principles of Measurement Systems*, 2nd ed., United Kingdom: Longman Scientific and Technical, 1988.
2. E. O. Doebelin, *Measurement Systems: Application and Design*, 4th ed., New York: McGraw-Hill, 1990.
3. J. P. Holman, *Experimental Methods for Engineers*, 5th ed., New York: McGraw-Hill, 1989.
4. W. J. Tompkins and J. G. Webster, *Interfacing Sensors to the IBM PC*, Englewood Cliffs, NJ: Prentice-Hall, 1988.

Appendix to Section 6.2

LIST OF MANUFACTURERS

Adsen Tech. Inc.

18310 Bedford Circle
La Puente, CA 91744
Fax: (818) 854-2776

Dynalco Controls

3690 N.W. 53rd Street
Ft. Lauderdale, FL 33309
Tel: (954) 739-4300 & (800) 368-6666
Fax: (954) 484-3376

Electro Corporation

1845 57th Street
Sarasato, FL 34243
Tel: (813) 355-8411 & (800) 446-5762
Fax: (813) 355-3120

Honeywell

Dept 722
11 W. Spring Street
Freeport, IL 61032
Tel: (800) 537-6945
Fax: (815) 235-5988

Kaman Inst. Co.

1500 Garden of the Gods Rd.
Colorado Springs, CO 80907
Tel: (719) 599-1132 & (800) 552-6267
Fax: (719) 599-1823

Kavlico Corporation

14501 Los Angeles Avenue
Moorpark, CA 93021
Tel: (805) 523-2000
Fax: (805) 523-7125

Lucas

1000 Lucas Way
Hampton, VA 23666
Tel: (800) 745-8008
Fax: (800) 745-8004

Motion Sensors Inc.

786 Pitts Chapel Rd.
Alizabath City, NC 27909
Tel: (919) 331-2080
Fax: (919) 331-1666

Rechner Electronics Ind. Inc.

8651 Buffalo Ave.
Niagara Falls, NY 14304
Tel: (800) 544-4106
Fax: (716) 283-2127

Reed Switch Developments Co. Inc.

P.O. Drawer 085297
Racine, WI 53408
Tel: (414) 637-8848
Fax: (414) 637-8861

Smith Research & Technology Inc.

205 Sutton Lane, Dept. TR-95
Colorado Springs, CO 80907
Tel: (719) 634-2259
Fax: (719) 634-2601

Smith Systems Inc.

6 Mill Creek Dr.
Box 667
Brevard, NC 28712
Tel: (704) 884-3490
Fax: (704) 877-3100

Standex Electronics

4538 Camberwell Rd.
Dept. 301L
Cincinnati, OH 45209
Tel: (513) 871-3777
Fax: (513) 871-3779

Turck Inc.

3000 Campus Drive
Minneapolis, MN 55441
Tel: (612) 553-7300 & (800) 544-7769
Fax: (612) 553-0708

Xolox Sensor Products

6932 Gettysburg Pike
Ft. Wayne, IN 46804
Tel: (800) 348-0744
Fax: (219) 432-0828

6.3 Capacitive Sensors—Displacement

Halit Eren and Wei Ling Kong

Capacitive sensors are extensively used in industrial and scientific applications. They are based on changes in capacitance in response to physical variations. These sensors find many diverse applications — from humidity and moisture measurements to displacement sensing. In some cases, the basic operational and sensing principles are common in dissimilar applications; and in other cases, different principles can be used for the same applications. For example, capacitive microphones are based on variations of spacing between plates in response to acoustical pressure, thus turning audio signals to variations in capacitance. On the other hand, a capacitive level indicator makes use of the changes in the relative permittivity between the plates. However, capacitive sensors are best known to be associated with displacement measurements for rotational or translational motions, as will be described next. Other applications of capacitance sensors such as humidity and moisture will be discussed.

Capacitive Displacement Sensors

The measurement of distances or displacements is an important aspect of many industrial, scientific, and engineering systems. The displacement is basically the vector representing a change in position of a body or point with respect to a reference point. Capacitive displacement sensors satisfy the requirements of applications where high linearity and wide ranges (from a few centimeters to a couple of nanometers) are needed.

The basic sensing element of a typical displacement sensor consists of two simple electrodes with capacitance C . The capacitance is a function of the distance d (cm) between the electrodes of a structure, the surface area A (cm²) of the electrodes, and the permittivity ϵ (8.85×10^{-12} F m⁻¹ for air) of the dielectric between the electrodes; therefore:

$$C = f(d, A, \epsilon) \quad (6.22)$$

There are three basic methods for realizing a capacitive displacement sensor: by varying d , A , or ϵ , as discussed below.

Variable Distance Displacement Sensors

A capacitor displacement sensor, made from two flat coplanar plates with a variable distance x apart, is illustrated in Figure 6.23. Ignoring fringe effects, the capacitance of this arrangement can be expressed by:

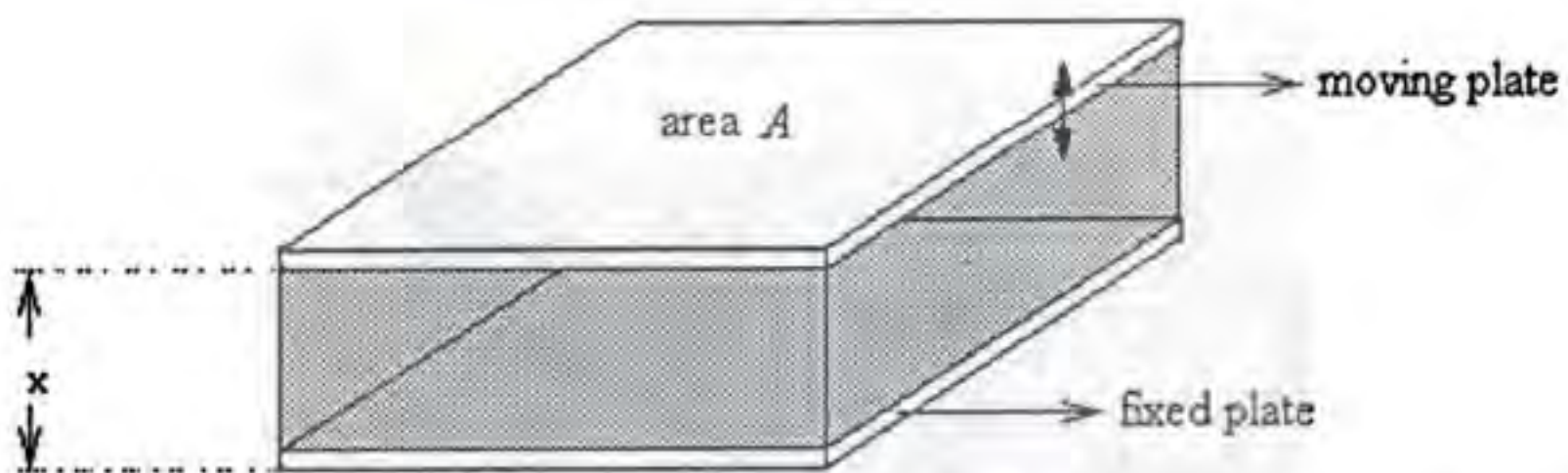


FIGURE 6.23 A variable distance capacitive displacement sensor. One of the plates of the capacitor moves to vary the distance between plates in response to changes in a physical variable. The outputs of these transducers are nonlinear with respect to distance x having a hyperbolic transfer function characteristic. Appropriate signal processing must be employed for linearization.

$$C(x) = \epsilon A/x = \epsilon_r \epsilon_0 A/x \quad (6.23)$$

where ϵ = the dielectric constant or permittivity

ϵ_r = the relative dielectric constant (in air and vacuum $\epsilon_r = 1$)

$\epsilon_0 = 8.854188 \times 10^{-12}$ F/m⁻¹, the dielectric constant of vacuum

x = the distance of the plates in m

A = the effective area of the plates in m²

The capacitance of this transducer is nonlinear with respect to distance x , having a hyperbolic transfer function characteristic. The sensitivity of capacitance to changes in plate separation is

$$dC/dx = -\epsilon_r \epsilon_0 A/x^2 \quad (6.24)$$

Equation 6.24 indicates that the sensitivity increases as x decreases. Nevertheless, from Equations 6.23 and 6.24, it follows that the percent change in C is proportional to the percent change in x . This can be expressed as:

$$dC/C = -dx/x \quad (6.25)$$

This type of sensor is often used for measuring small incremental displacements without making contact with the object.

Variable Area Displacement Sensors

Alternatively, the displacements may be sensed by varying the surface area of the electrodes of a flat plate capacitor, as illustrated in Figure 6.24. In this case, the capacitance would be:

$$C = \epsilon_r \epsilon_0 (A - wx)/d \quad (6.26)$$

where w = the width

wx = the reduction in the area due to movement of the plate

Then, the transducer output is linear with displacement x . This type of sensor is normally implemented as a rotating capacitor for measuring angular displacement. The rotating capacitor structures are also used as an output transducer for measuring electric voltages as capacitive voltmeters.

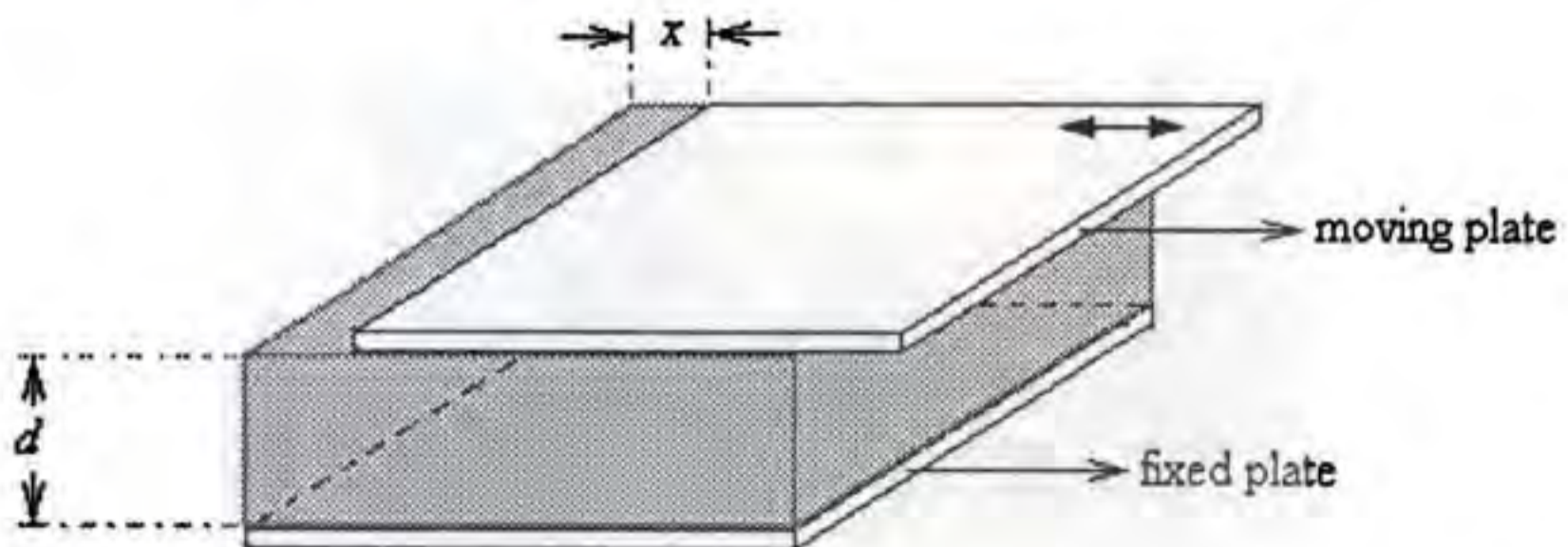


FIGURE 6.24 A variable area capacitive displacement sensor. The sensor operates on the variation in the effective area between plates of a flat-plate capacitor. The transducer output is linear with respect to displacement x . This type of sensor is normally implemented as a rotating capacitor for measuring angular displacement.

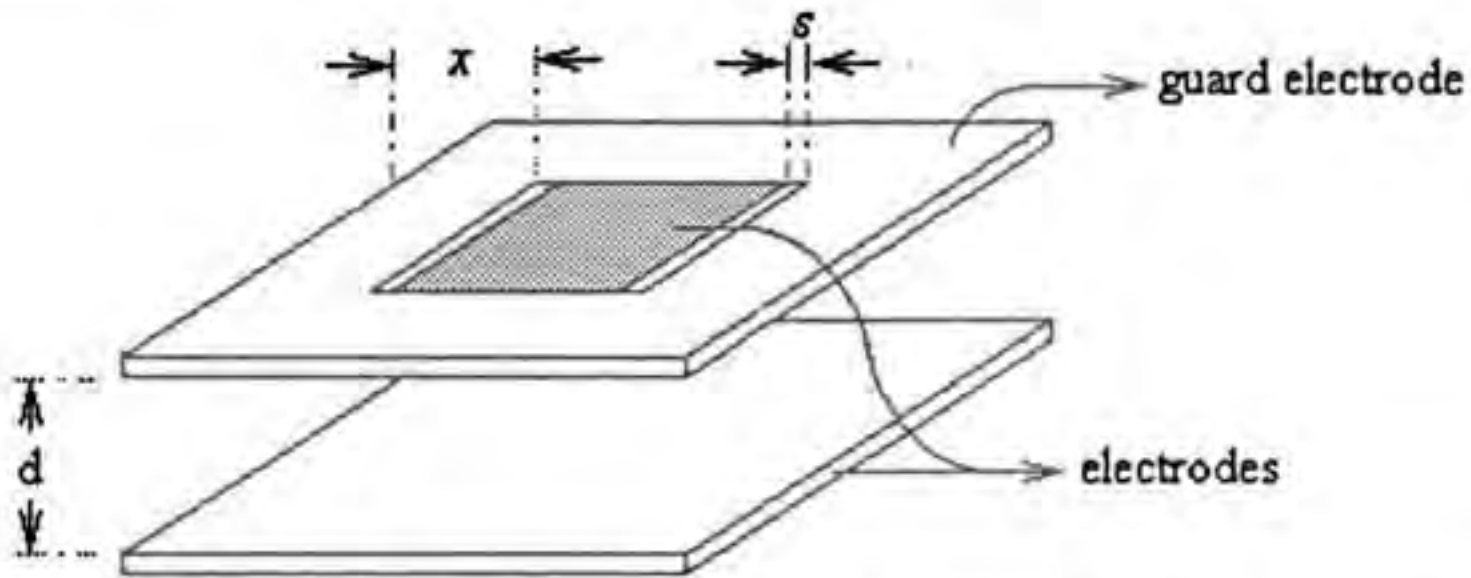


FIGURE 6.27 A typical smart capacitive position sensor. This type of microstructure position sensor contains three electrodes, two of which are fixed and the third electrode moves infinitesimally relative to the others. Although the response is highly nonlinear the integrated chip contains linearization circuits. They feature a 0 mm to 1 mm measuring range with 1 μ m accuracy.

$$\delta < \exp(-\pi x/d) \tag{6.31}$$

where x is the width of the guard and d is the distance between the electrodes. Since this deviation introduces nonlinearity, δ is required to be less than 100 ppm. Another form of deviation also exists between the small electrode and the surrounding guard, particularly for gaps

$$\delta < \exp(-\pi d/s) \tag{6.32}$$

where s is the width of the gap. When the gap width, s , is less than 1/3 of the distance between electrodes, this deviation is negligible.

For signal processing, the system uses the three-signal concept. The capacitor C_x is connected to an inverting operational amplifier and oscillator. If the external movements are linear, by taking into account the parasitic capacitors and offsetting effects, the following equation can be written:

$$M_x = mC_x + M_{off} \tag{6.33}$$

where m is the unknown gain and M_{off} is the unknown offset. By performing the measurement of a reference C_{ref} , by measuring the offset, M_{off} , and by making $m = 0$, the parameters m and M_{off} can be eliminated. The final measurement result for the position, P_{os} , can be defined as:

$$P_{os} = \frac{M_{ref} - M_{off}}{M_x - M_{off}} \tag{6.34}$$

In this case, the sensor capacitance C_x can be simplified to:

$$C_x = \frac{\epsilon A_x}{d_0 + \Delta d} \tag{6.35}$$

where A_x is the area of the electrode, d_0 is the initial distance between them, ϵ is the dielectric constant, and Δd is the displacement to be measured. For the reference electrodes, the reference capacitance may be found by:

$$C_{\text{ref}} = \frac{\epsilon A_{\text{ref}}}{d_{\text{ref}}} \quad (6.36)$$

with A_{ref} the area and d_{ref} the distance. Substitution of Equations 6.35 and 6.36 into Equations 6.33 and 6.34 yields:

$$P_{\text{os}} = \frac{A_{\text{ref}}(d_0 + \Delta d)}{A_x d_{\text{ref}}} = a_1 \frac{\Delta d}{d_{\text{ref}}} + a_0 \quad (6.37)$$

P_{os} is a value representing the position if the stable constants a_1 and a_0 are unknown. The constant $a_1 = A_{\text{ref}}/A_x$ becomes a stable constant so long as there is good mechanical matching between the electrode areas. The constant $a_0 = (A_{\text{ref}} d_0)/(A_x d_{\text{ref}})$ is also a stable constant for fixed d_0 and d_{ref} . These constants are usually determined by calibration repeated over a certain time span. In many applications, these calibrations are omitted if the displacement sensor is part of a larger system where an overall calibration is necessary. This overall calibration usually eliminates the requirement for a separate determination of a_1 and a_0 .

The accuracy of this type of system could be as small as 1 μm over a 1 mm range. The total measuring time is better than 0.1 s. The capacitance range is from 1 pF to 50 fF. Interested readers should refer to [4] at the end of this chapter.

Capacitive Pressure Sensors

A commonly used two-plate capacitive pressure sensor is made from one fixed metal plate and one flexible diaphragm, as shown in Figure 6.28. The flat circular diaphragm is clamped around its circumference and bent into a curve by an applied pressure P . The vertical displacement y of this system at any radius r is given by:

$$y = 3(1 - \nu^2)(a^2 - r^2)P/16Et^3 \quad (6.38)$$

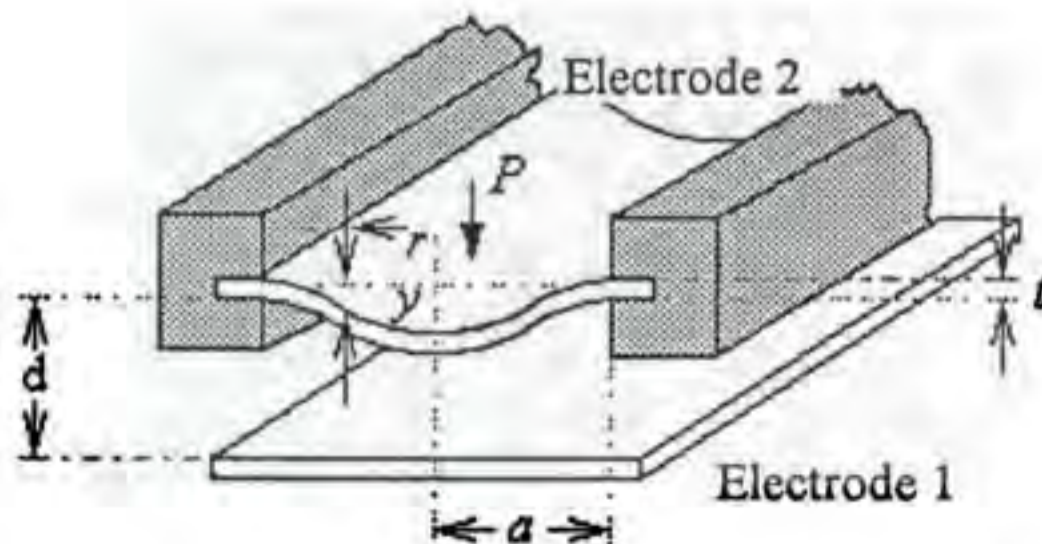


FIGURE 6.28 A capacitive pressure sensor. These pressure sensors are made from a fixed metal plate and a flexible diaphragm. The flat flexible diaphragm is clamped around its circumference. The bending of the flexible plate is proportional to the applied pressure P . The deformation of the diaphragm results in changes in capacitance.

where a = the radius of diaphragm
 t = the thickness of diaphragm
 E = Young's modulus
 ν = Poisson's ratio

Deformation of the diaphragm means that the average separation of the plates is reduced. Hence, the resulting increase in the capacitance ΔC can be calculated by:

$$\Delta C/C = (1 - \nu^2) a^4 P / 16 E t^3 \quad (6.39)$$

where d is the initial separation of the plates and C is the capacitance at zero pressure.

Another type of sensor is the differential capacitance pressure sensor shown in Figure 6.29. The capacitances C_1 and C_2 of the sensor change with respect to the fixed central plate in response to the applied pressures P_1 and P_2 . Hence, the output of the sensor is proportional to $(P_1 - P_2)$. The signals are processed using one of the techniques described in the "Signal Processing" section of this chapter.

Capacitive Accelerometers and Force Transducers

In recent years, capacitive-type micromachined accelerometers, as illustrated in Figure 6.30, are gaining popularity. These accelerometers use the proof mass as one plate of the capacitor and use the other plate as the base. When the sensor is accelerated, the proof mass tends to move; thus, the voltage across the capacitor changes. This change in voltage corresponds to the applied acceleration.

In Figure 6.30, let $F(x)$ be the positive force in the direction in which x increases. Neglecting all losses (due to friction, resistance, etc.), the energy balance of the system can be written for an infinitesimally small displacement dx , electrical energy dE_e , and field energy dE_f of the electrical field between the electrodes as:

$$dE_m + dE_e = dE_f \quad (6.40)$$

in which:

$$dE_m = F(x) dx \quad (6.41)$$

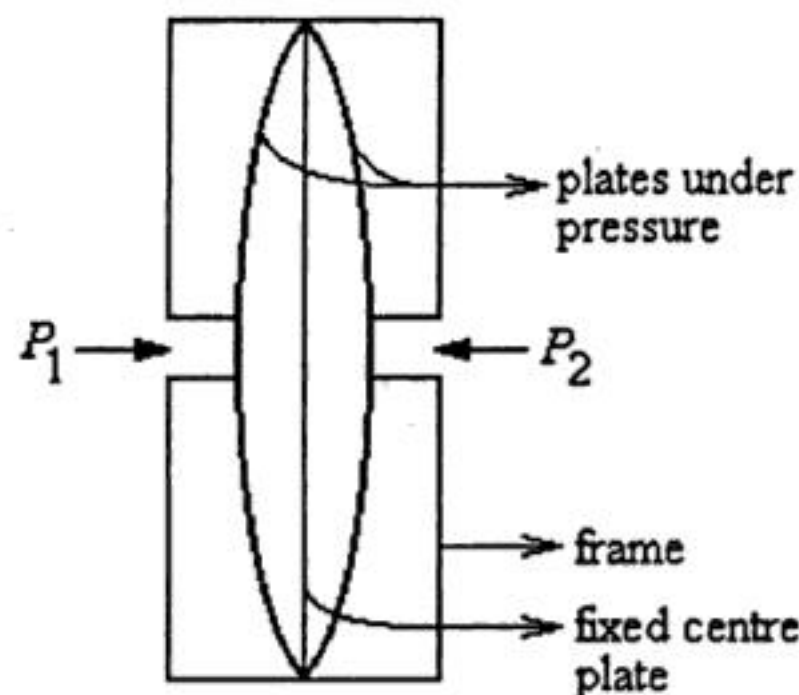


FIGURE 6.29 A differential capacitive pressure sensor. The capacitances C_1 and C_2 of the sensor changes due to deformation in the outer plates, with respect to the fixed central plate in response to the applied pressures P_1 and P_2 .

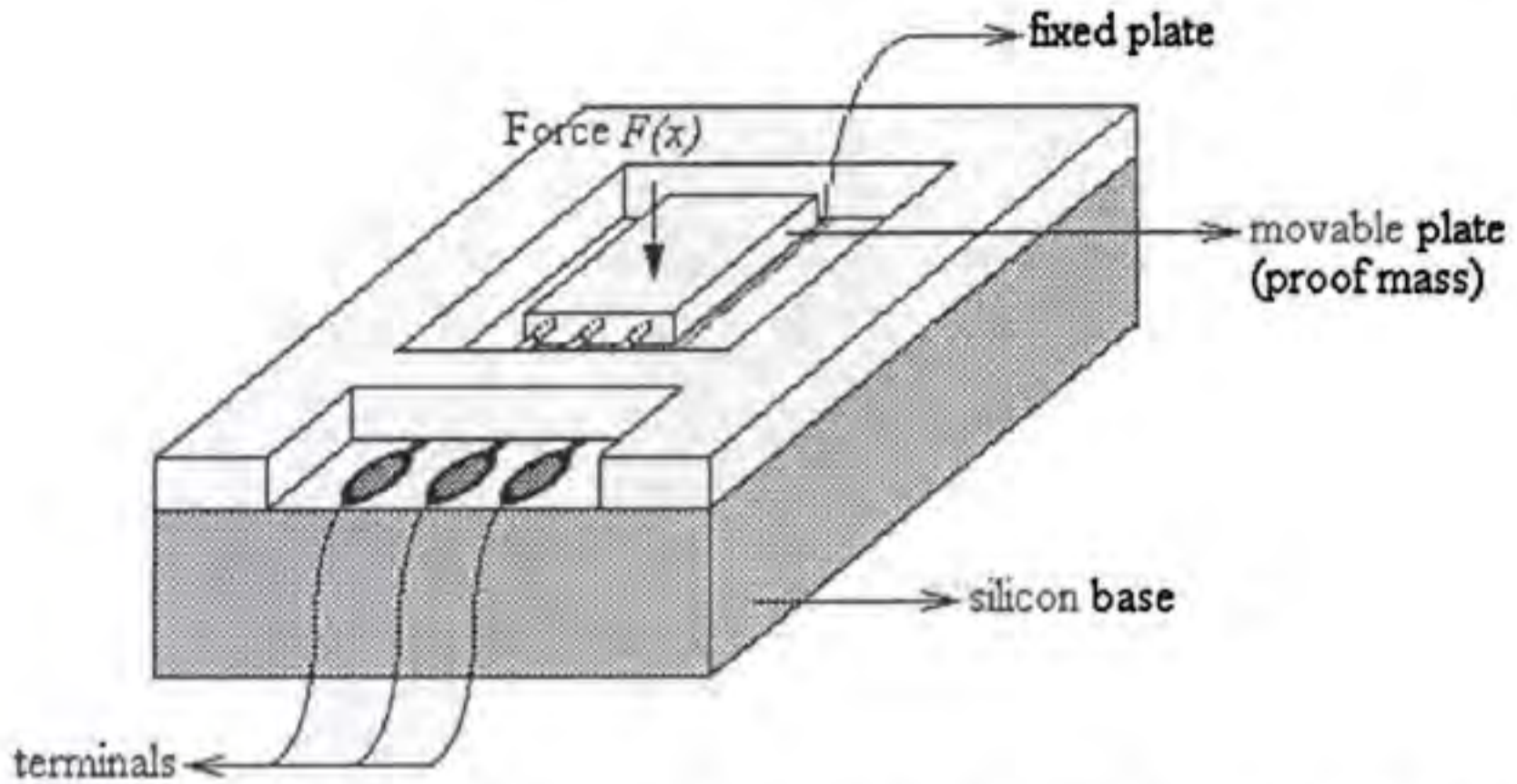


FIGURE 6.30 A capacitive force transducer. A typical capacitive micromachined accelerometer has one of the plates as the proof mass. The other plate is fixed, thus forming the base. When the sensor is accelerated, the proof mass tends to move, thus varying the distance between the plates and altering the voltage across the capacitor. This change in voltage is made to be directly proportional to the applied acceleration.

Also,

$$dE_m = d(QV) = Q dV + V dQ \quad (6.42)$$

If the supply voltage V across the capacitor is kept constant, it follows that $dV = 0$. Since $Q = VC(x)$, the Coulomb force is given by:

$$F(x) = -V^2 \frac{dC(x)}{dx} \quad (6.43)$$

Thus, if the movable electrode has complete freedom of motion, it will have assumed a position in which the capacitance is maximal; also, if C is a linear function of x , the force $F(x)$ becomes independent of x .

Capacitive silicon accelerometers are available in a wide range of specifications. A typical lightweight sensor will have a frequency range of 0 to 1000 Hz, and a dynamic range of acceleration of $\pm 2 g$ to $\pm 500 g$.

Capacitive Liquid Level Measurement

The level of a nonconducting liquid can be determined by capacitive techniques. The method is generally based on the difference between the dielectric constant of the liquid and that of the gas or air above it. Two concentric metal cylinders are used for capacitance, as shown in Figure 6.31. The height of the liquid, h , is measured relative to the total height, l . Appropriate provision is made to ensure that the space between the cylindrical electrodes is filled by the liquid to the same height as the rest of the container. The usual operational conditions dictate that the spacing between the electrodes, $s = r_2 - r_1$, should be much less than the radius of the inner electrode, r_1 . Furthermore, the tank height should be much greater than r_2 . When these conditions apply, the capacitance is approximated by:

$$C = \frac{\epsilon_l(l) + \epsilon_g(h-l)}{4.6 \log \left[1 - \left(\frac{s}{r} \right) \right]} \quad (6.44)$$

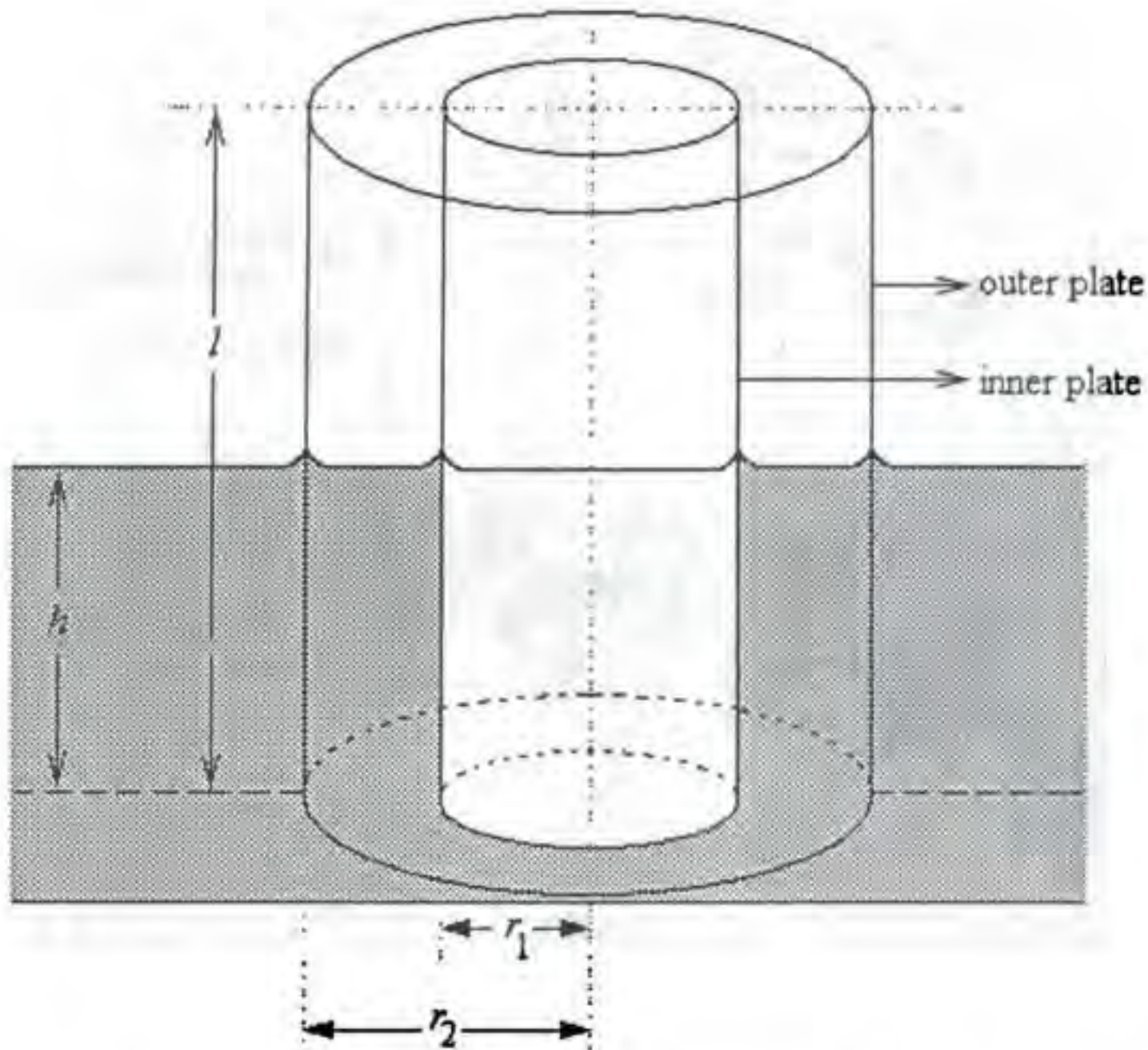


FIGURE 6.31 A capacitive liquid level sensor. Two concentric metal cylinders are used as electrodes of a capacitor. The value of the capacitance depends on the permittivity of the liquid and that of the gas or air above it. The total permittivity changes depending on the liquid level. These devices are usually applied in nonconducting liquid applications.

where ϵ_l and ϵ_g are the dielectric constants of the liquid and gas (or air), respectively. The denominator of the above equation contains only terms that relate to the fixed system. Therefore, they become a single constant. A typical application is the measurement of the amount of gasoline in a tank in airplanes. The dielectric constant for most compounds commonly found in gasoline is approximately equal to 2, while that of air is approximately unity. A linear change in capacitance with gasoline level is expected for this situation. Quite high accuracy can be achieved if the denominator is kept quite small, thus accentuating the level differences. These sensors often incorporate an ac deflection bridge.

Capacitive Humidity and Moisture Sensors

The permittivities of atmospheric air, of some gases, and of many solid materials are functions of moisture content and temperature. Capacitive humidity devices are based on the changes in the permittivity of the dielectric material between plates of capacitors. The main disadvantage of this type sensor is that a relatively small change in humidity results in a capacitance large enough for a sensitive detection.

Capacitive humidity sensors enjoy wide dynamic ranges, from 0.1 ppm to saturation points. They can function in saturated environments for long periods of time, a characteristic that would adversely affect many other humidity sensors. Their ability to function accurately and reliably extends over a wide range of temperatures and pressures. Capacitive humidity sensors also exhibit low hysteresis and high stability with minimal maintenance requirements. These features make capacitive humidity sensors viable for many specific operating conditions and ideally suitable for a system where uncertainty of unaccounted conditions exists during operations.

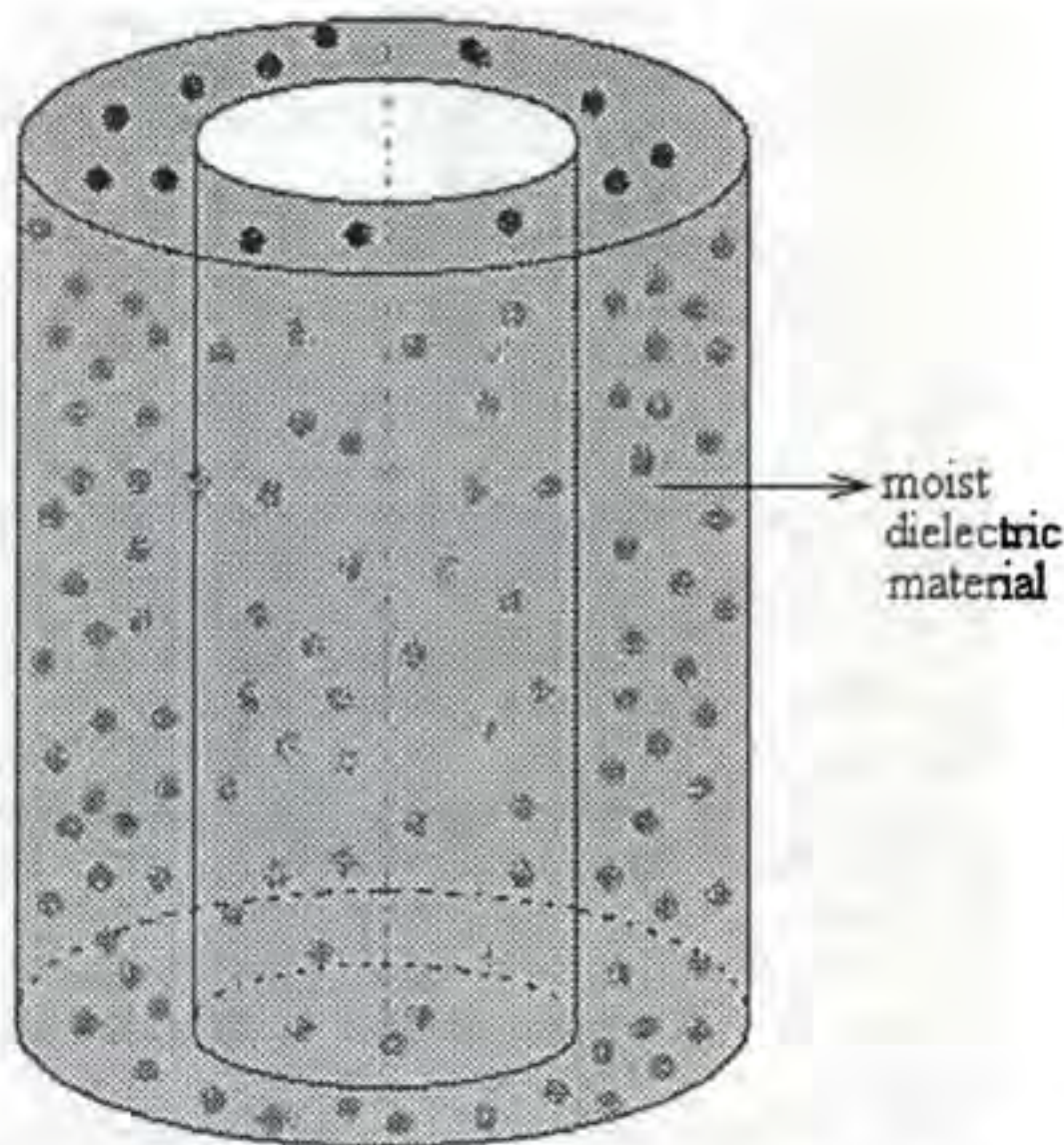


FIGURE 6.33 A capacitive moisture sensor. The permittivity of material between two cylindrical or parallel plates with fixed dimensions changes, depending on the moisture level of the materials in the chamber. The variations in capacitance values with respect to water content is processed. The capacitor is incorporated as a part of an oscillatory circuit operating at a suitable frequency, usually at radio frequencies.

Capacitive moisture sensors must be calibrated for samples made from different materials, as the materials themselves demonstrate different permittivities. Accurate temperature is necessary as the dielectric constant may be highly dependent on temperature. Most of these devices are built to operate at temperature ranges of 0°C to 50°C, supported by tight temperature compensation circuits. Once calibrated for a specific application, they are suitable for measuring moisture in the range of 0% to 40%.

Signal Processing

Generally, capacitive type pickups require relatively complex circuitry in comparison to many other sensor types, but they have the advantage of mechanical simplicity. They are also sensitive, having minimum mechanical loading effects. For signal processing, these sensors are usually incorporated either in ac deflection bridge circuits or oscillator circuits. In practice, capacitive sensors are not pure capacitances but have associated resistances representing losses in the dielectric. This can have an important influence in the design of circuits, particularly in oscillator circuits. Some of the signal processing circuits are discussed below.

Operational Amplifiers and Charge Amplifiers

One method of eliminating the nonlinearity of the relationship between the physical variable, (e.g., two-plate displacement sensors) and capacitance C is through the use of operational amplifiers, as illustrated in Figure 6.34. In this circuit, if the input impedance of the operational amplifier is high, the output is not saturated, and the input voltage is small, it is possible to write:

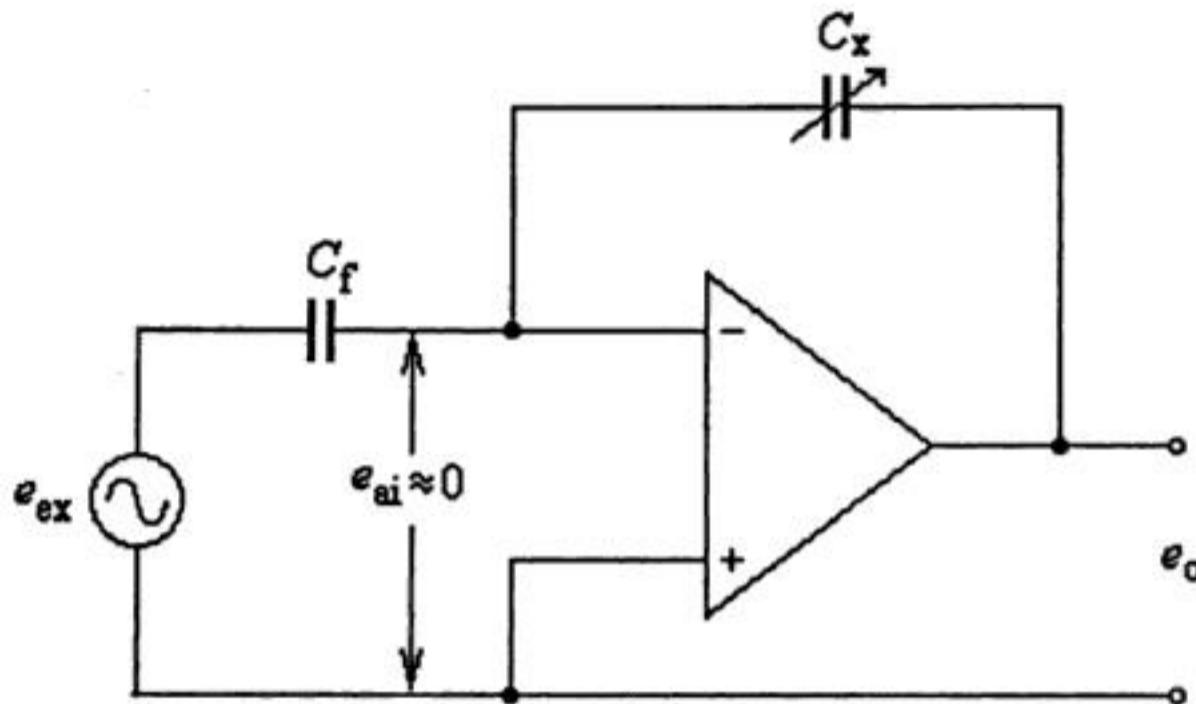


FIGURE 6.34 An operational amplifier signal processor. This method is useful to eliminate the nonlinearity in the signals generated by capacitive sensors. By this type of arrangement, the output voltage can be made directly proportional to variations in the signal representing the nonlinear operation of the device.

$$1/C_f = \int i_f dt = e_{ex} - e_{ai} = e_{ex} \tag{6.45}$$

$$1/C_x = \int i_x dt = e_o - e_{ai} = e_o \tag{6.46}$$

$$i_f + i_x - i_{ai} = 0 = i_f + i_x \tag{6.47}$$

Manipulation of these equations yields:

$$e_o = -C_f e_{ex} / C_x \tag{6.48}$$

Substituting the value of C_x yields:

$$e_o = -C_f x e_{ex} / \epsilon A \tag{6.49}$$

Equation 6.49 shows that the output voltage is directly proportional to the plate separation x , thus giving linearity for all variations in motion.

However, a practical circuit requires a resistance across C_f to limit output drift. The value of this resistance must be greater than the impedance of C_f at the lowest frequency of interest. Also, because the transducer impedance is assumed to be purely capacitive, the effective gain is independent of frequency.

A practical charge amplifier circuit is depicted in Figure 6.35. In this case, the effective feedback resistance R_{ef} is given by:

$$R_{ef} = R_3 (R_1 + R_2) / R_2 \tag{6.50}$$

It is possible to reduce the output drift substantially by selecting the resistors suitably. The accuracy of this circuit can be improved further by cascading two or more amplifiers. In this way, a substantial improvement in the signal-to-noise ratio can also be achieved. In the inverting input, the use of resistor R_4 is necessary because of bias currents.

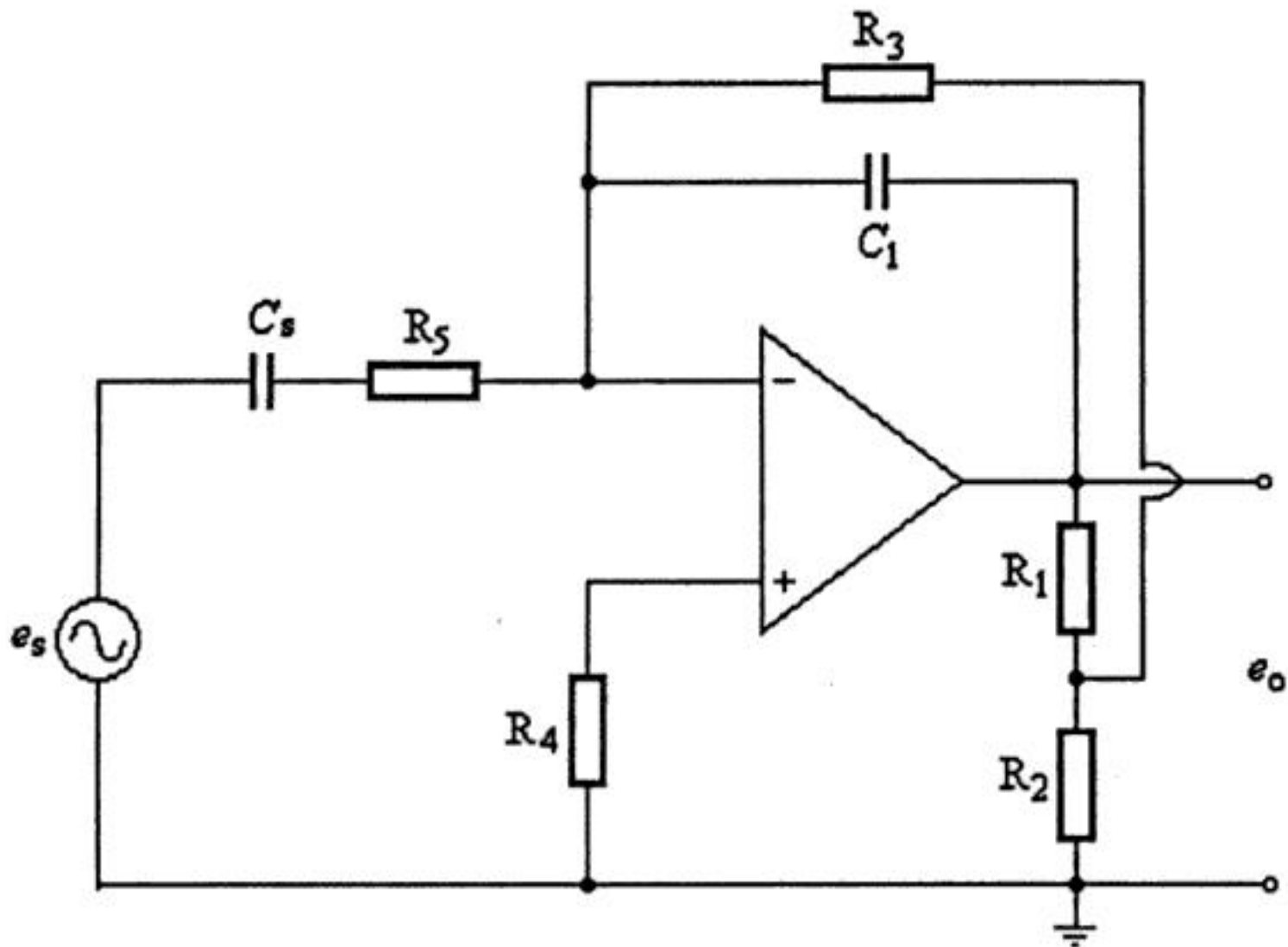


FIGURE 6.35 A practical charge amplifier. The effective feedback resistance is a function of other resistances. It is possible to reduce the output drift substantially by selecting the resistors suitably. The accuracy of this circuit can be improved further by cascading two or more amplifiers, thereby substantially improving the signal-to-noise ratio.

Pulse Width Modulation

As in the case of some capacitive vibrational displacement sensors, the output of the sensor may be an amplitude-modulated wave as shown in Figure 6.36. When rectified, the average value of this wave gives the mean separation of the plates. The vibration amplitude around this mean position may be extracted by a demodulator and a low-pass filter circuit. The output of the low-pass filter is a direct indication of vibrations, and the waveform can be viewed on an oscilloscope.

Square Wave Linearization

Another linearization technique applied in capacitive pressure transducers and accelerometers is pulse width modulation. The transducer consists of two differential capacitors as shown in Figure 6.37. The voltages of these capacitors, e_1 and e_2 , switch back and forth with a high excitation frequency (e.g., 400 kHz) between excitation voltage and ground. The system is arranged in such a way that the output voltage is the average voltage difference between e_1 and e_2 . At null position, $e_1 = e_2$, the output is a symmetrical square wave with zero average value. As the relative positions of the plates change, due to vibration, the average value of the output voltage shifts from the zero average value and becomes positive or negative depending on the direction of the displacement. Hence, the output voltage can be expressed by:

$$e_o = e_{ex} (C_1 - C_2) / (C_1 + C_2) \quad (6.51)$$

Substituting:

$$C_1 = C_0 x_0 / (x_0 - x_i) \quad \text{and} \quad C_2 = C_0 x_0 / (x_0 + x_i)$$

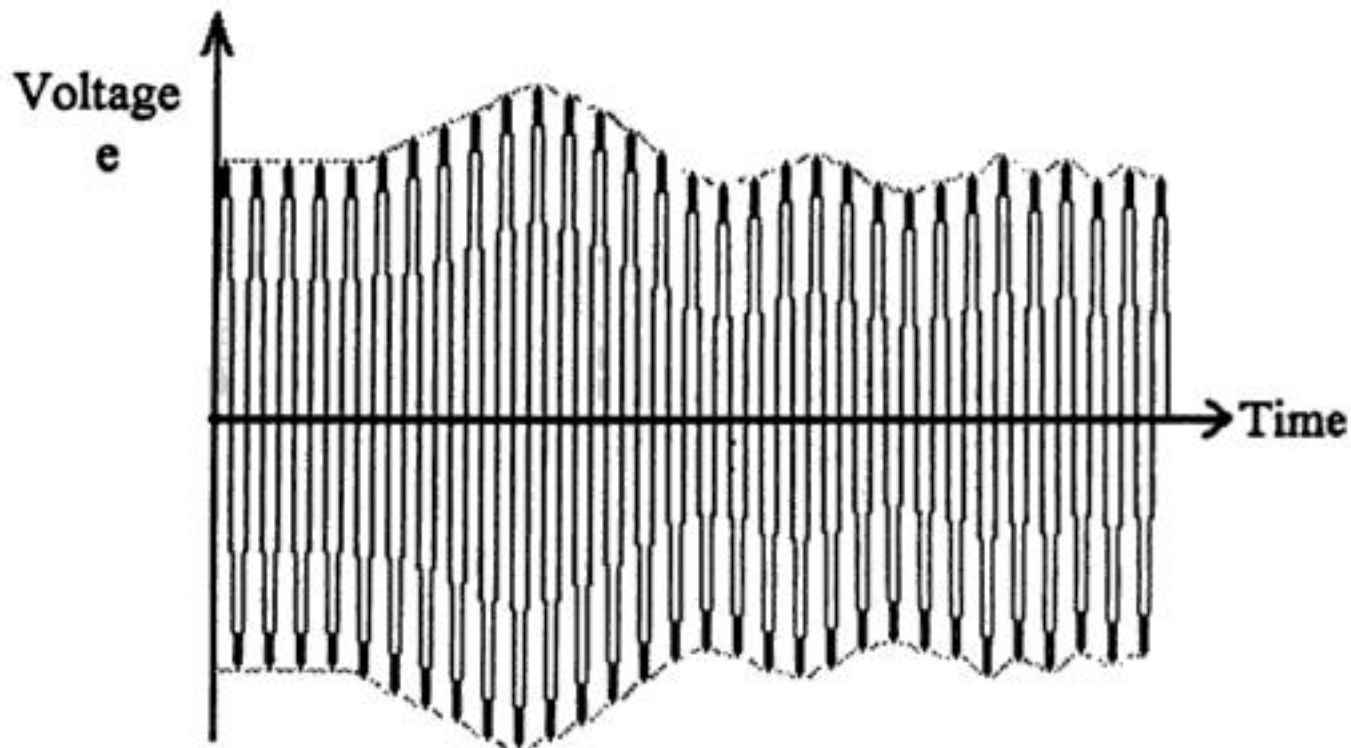


FIGURE 6.36 A amplitude modulated signal. It is possible to configure some sensors to give a amplitude-modulated signals, as in the case of capacitive vibrational displacement sensors. When rectified, the average value of this wave gives the mean separation of the plates. The vibration amplitude around this mean position can be extracted by a demodulator and low-pass filter circuit. The output of the low-pass filter is a direct indication of vibrations.

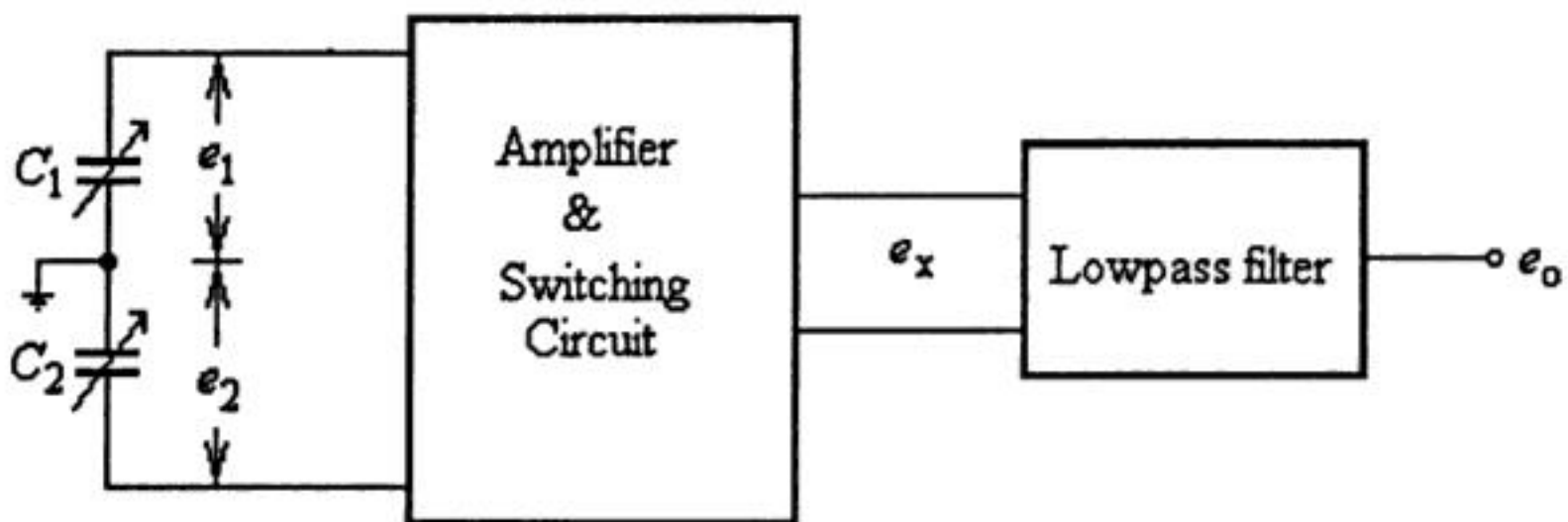


FIGURE 6.37 Block diagram of a square-wave linearization circuit. This is particularly useful for differential capacitance type sensors. The voltages of these two capacitors are made to switch back and forth with a high excitation frequency between excitation voltage and ground. As the relative positions of the plates change due to vibration, the average value of the output voltage becomes positive or negative, depending on the direction of the displacement.

yields

$$e_o = e_{ex} x_i / x_0 \tag{6.52}$$

Thus, the output is directly proportional to the variable x_i .

Feedback Linearization

Linearization of a capacitance transducer can also be obtained using a feedback system that adjusts capacitor current amplitude so that it stays constant at a reference value for all displacements. This is accomplished by obtaining a dc signal proportional to capacitor current from a demodulator, comparing this current with the reference current, and adjusting the voltage amplitude of the system excitation oscillator until the two currents agree. If the capacitor current is kept constant irrespective of capacitor motion, then the voltage amplitude is linearly related to x as:

$$e = K x_1 \quad (6.53)$$

where

$$K = \left| \frac{i_c}{\omega C_0} \right| x_0 \quad (6.54)$$

Oscillator Circuits

In many applications, the resultant changes in the capacitance of capacitive transducers can be measured with a suitable ac bridge such as Wein bridge or Schering bridge. However, in a majority of cases, improvised versions of bridges are used as oscillator circuits for capacitive signal processing. The transducer is configured as a part of the oscillatory circuit that causes changes in the frequency of the oscillations. This change in frequency is scaled to be a measure of the magnitude of the physical variable.

As part of the oscillator circuits, the capacitive transducer has excellent frequency response and can measure both static and dynamic phenomena. Its disadvantages include sensitivity to temperature variations and the possibility of erratic or distorted signals due to long lead lengths. Also, the receiving instrumentation may be large and complex, and it often includes a second fixed-frequency oscillator for heterodyning purposes. The difference frequency thus produced can be read by an appropriate output device such as an electronic counter.

References

1. J. P. Bentley, *Principles of Measurement Systems*, 2nd ed., United Kingdom: Longman Scientific and Technical, 1988.
2. E. O. Doebelin, *Measurement Systems: Application and Design*, 4th ed., New York: McGraw-Hill, 1990.
3. J. P. Holman, *Experimental Methods for Engineers*, 5th ed., New York: McGraw-Hill, 1989.
4. F. T. Noth and G. C. M. Meijer, A Low-Cost, Smart Capacitive Position Sensor, *IEEE Trans. Instrum. Meas.*, 41, 1041-1044, 1992.

Appendix to Section 6.3

List of Manufacturers

ANALITE Inc.

24-T Newtown Plaza
Plainview, NY 11803
Tel: (800) 229-3357

FSI/FORK Standards Inc.

668 Western Avenue
Lombard, IL 60148-2097
Tel: (708) 932-9380

Gordon Engineering Corp.

67 Del Mar Drive
Brookfield, CT 06804
Tel: (203) 775-4501

Hecon Corp.

15-T Meridian Rd.
Eatontown, NJ 07724
Tel: (800) 524-1669

Kistler Instrumentation Corp.

Amherst, NY 14228-2171
Tel: (716) 691-5100
Fax: (716) 691-5226

Locon Sensor Systems, Inc.

1750 S. Eber Road
P.O. Box 789
Holland, OH 43526
Tel: (419) 865-7651
Fax: (419) 865-7756

Rechner Electronic Industries Inc.

8651 Buffalo Avenue, Box 7
Niagara Falls, NY 14304
Tel: (800) 544-4106

RDP Electrosense, Inc.

2216-Dept. B
Pottstown, PA
Tel: (800) 334-5838

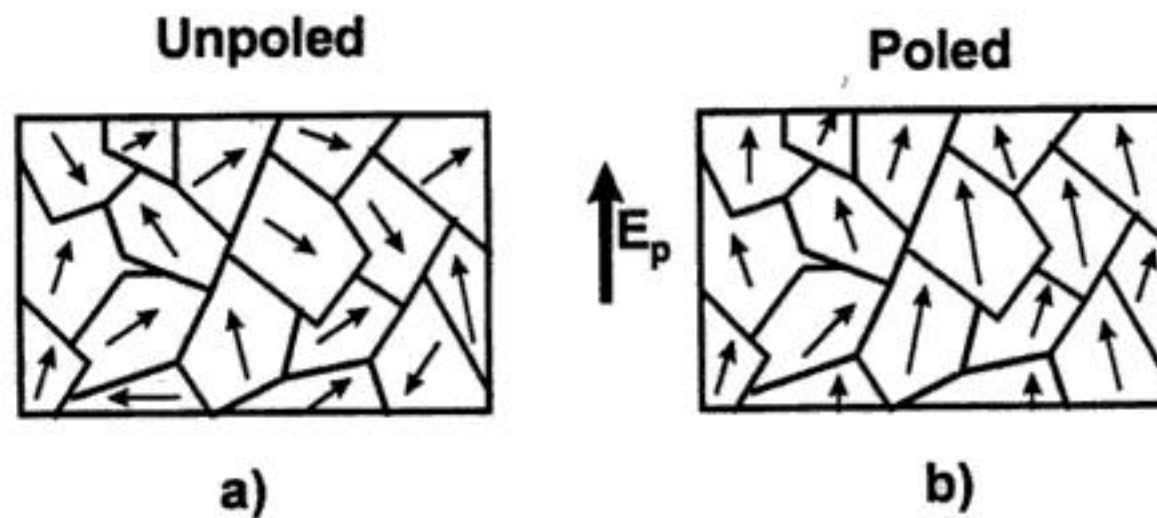


FIGURE 6.39 Schematic of the poling process in piezoelectric ceramics: (a) in the absence of an electric field, the domains have random orientation of polarization; (b) the polarization within the domains are aligned in the direction of the electric field.

ceramic. By applying a strong dc electric field at a temperature just below the Curie temperature, the spontaneous polarization in each grain gets oriented toward the direction of the applied field. This is schematically shown in Figure 6.39(b). Although all of the domains in a ceramic can never be fully aligned along the poling axis due to symmetry limitations, the ceramic ends up with a net polarization along the poling axis.

The largest class of piezoelectric ceramics is made up of mixed oxides containing corner-sharing octahedra of O^{2-} ions. The largest structure type, built with corner-shared oxygen octahedra, is the perovskite family, which is discussed in the following section.

Perovskites

Perovskite is the name given to a group of materials with general formula ABO_3 having the same structure as the mineral calcium titanate ($CaTiO_3$). Piezoelectric ceramics having this structure include barium titanate ($BaTiO_3$), lead titanate ($PbTiO_3$), lead zirconate titanate ($PbZr_xTi_{1-x}O_3$, or PZT), lead lanthanum zirconate titanate [$Pb_{1-x}La_x(Zr_yTi_{1-y})_{1-x/4}O_3$, or PLZT], and lead magnesium niobate [$PbMg_{1/3}Nb_{2/3}O_3$, or PMN]. Several of these ceramics are discussed below.

The piezoelectric effect in $BaTiO_3$ was discovered in the 1940s [4], and it became the first piezoelectric ceramic developed. It replaced Rochelle salt because it is more stable, has a wider temperature range of operation, and is easily manufacturable. The Curie point, T_0 , is about $130^\circ C$. Above $130^\circ C$, a nonpiezoelectric cubic phase is stable, where the center of positive charge (Ba^{2+} and Ti^{4+}) coincides with the center of the negative charge (O^{2-}) (Figure 6.40(a)). When cooled below the Curie point, a tetragonal structure (shown in Figure 6.40(b)) develops where the center of positive charge is displaced relative to the O^{2-} ions, leading to the formation of electric dipoles. Barium titanate has a relative dielectric constant ϵ_{33} of 1400 when unpoled, and 1900 when poled [2, 4]. The d_{15} and d_{33} coefficients of $BaTiO_3$ are 270 and $191 \times 10^{-12} C N^{-1}$, respectively. The k for $BaTiO_3$ is approximately 0.5. The large room temperature dielectric constant in barium titanate has led to its wide use in multilayer capacitor applications.

Lead titanate, $PbTiO_3$, first reported to be ferroelectric in 1950 [4], has a similar structure to $BaTiO_3$, but with a significantly higher Curie point ($T_0 = 490^\circ C$). Pure lead titanate is difficult to fabricate in bulk form. When cooled through the Curie point, the grains go through a cubic to tetragonal phase change, leading to large strain and ceramic fracturing. This spontaneous strain can be decreased by the addition of dopants such as Ca, Sr, Ba, Sn, and W. Calcium-doped $PbTiO_3$ [7] has a relative dielectric constant ϵ_{33} of 200, a d_{33} of $65 \times 10^{-12} C/N$, and a k of approximately 0.5. The addition of calcium results in a lowering of the Curie point to $225^\circ C$. The main applications of lead titanate are hydrophones and sonobuoys.

Lead zirconate titanate (PZT) is a binary solid solution of $PbZrO_3$ (an antiferroelectric orthorhombic structure) and $PbTiO_3$ (a ferroelectric tetragonal perovskite structure) [2–4]. It has a perovskite structure, with the Zr^{4+} and Ti^{4+} ions occupying the B site of the ABO_3 structure at random. At the morphotropic

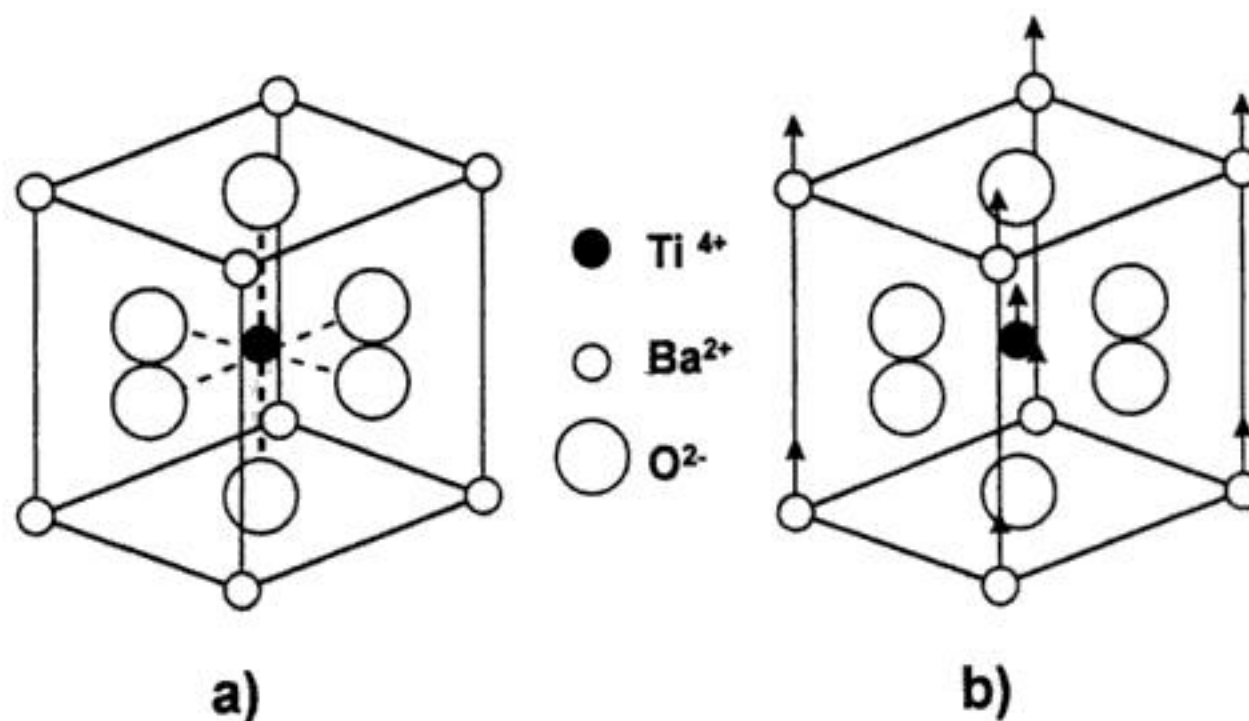


FIGURE 6.40 The crystal structure of BaTiO₃: (a) above the Curie point, the cell is cubic; (b) below the Curie point, the cell is tetragonal with Ba²⁺ and Ti⁴⁺ ions displaced relative to O²⁻ ions.

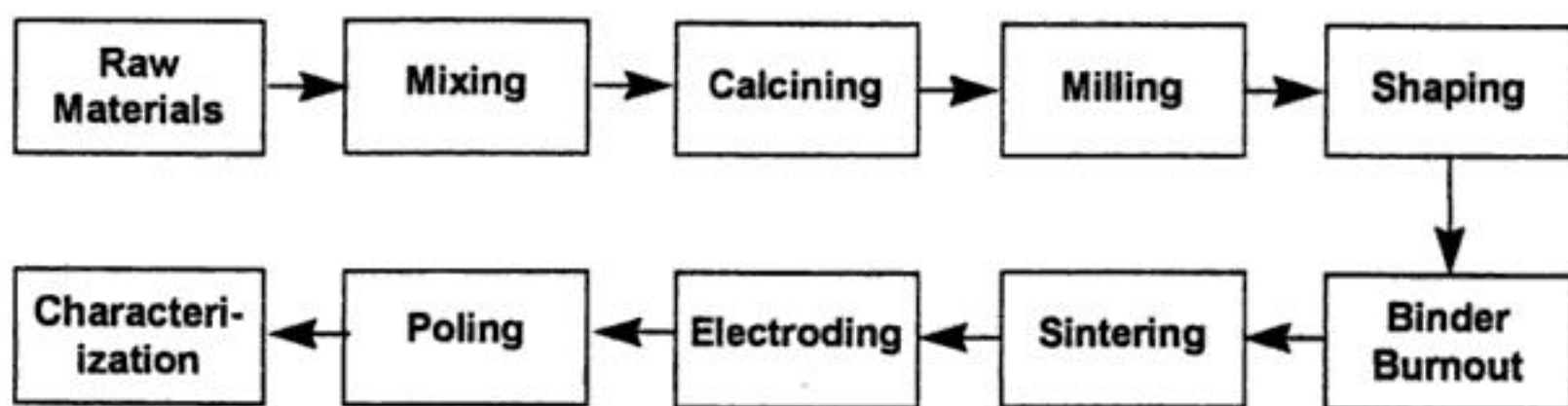


FIGURE 6.41 Flow chart for the processing of piezoelectric ceramics.

phase boundary (MPB) separating the tetragonal and orthorhombic phases, PZT shows excellent piezoelectric properties. At room temperature, the MPB is at a Zr/Ti ratio of 52/48, resulting in a piezoelectric ceramic which is extremely easy to pole. Piezoelectric PZT at the MPB is usually doped by a variety of ions to form what are known as "hard" and "soft" PZTs. Hard PZT is doped with acceptor ions, such as K⁺ or Na⁺ at the A site, or Fe³⁺, Al³⁺, or Mn³⁺ at the B site. This doping lowers the piezoelectric properties, and makes the PZT more difficult to pole or depole. Typical piezoelectric properties of hard PZT include [5, 7]: Curie point, T_0 , of 365°C, ϵ_{33} of 1700–1750 (poled), a piezoelectric charge coefficient d_{33} of 360 to $370 \times 10^{-12} \text{ C N}^{-1}$, and a coupling coefficient of about 0.7. Soft PZT is doped with donor ions such as La³⁺ at the A site, or Nb⁵⁺ or Sb⁵⁺ at the B site. It has very high piezoelectric properties, and is easy to pole or depole. Typical piezoelectric properties of soft PZT include [5, 7]: Curie point, T_0 , of 210°C, relative dielectric constant ϵ_{33} of 3200–3400 (poled), a d_{33} of 580 to $600 \times 10^{-12} \text{ C N}^{-1}$, and a coupling coefficient k_{33} of 0.7.

Processing of Piezoelectric Ceramics

The electromechanical properties of piezoelectric ceramics are largely influenced by their processing conditions. Each step of the process must be carefully controlled to yield the best product. Figure 6.41 is a flow chart of a typical oxide manufacturing process for piezoelectric ceramics. The high-purity raw materials are accurately weighed according to their desired ratio, and mechanically or chemically mixed. During the calcination step, the solid phases react to yield the piezoelectric phase. After calcining, the solid mixture is ground into fine particles by milling. Shaping is accomplished by a variety of ceramic

TABLE 6.7 Advantages (+) and Disadvantages (-) of Piezoelectric Ceramics, Polymers and Composites

Parameter	Ceramic	Polymer	Ceramic/Polymer Composite
Acoustic impedance	High (-)	Low (+)	Low (+)
Coupling factor	High (+)	Low (-)	High (+)
Spurious modes	Many (-)	Few (+)	Few (+)
Dielectric constant	High (+)	Low (-)	Medium (+)
Flexibility	Stiff (-)	Flexible (+)	Flexible (+)
Cost	Cheap (+)	Expensive (-)	Medium (+)

Adapted from T. R. Gururaja, *Amer. Ceram. Soc. Bull.*, 73, 50, 1994.

processing techniques, including powder compaction, tape casting, slip casting, or extrusion. During the shaping operation, organic materials are typically added to the ceramic powder to improve its flow and binding characteristics. These organics are removed in a low temperature (500 to 600°C) binder burn-off step.

After burnout, the ceramic structure is sintered to an optimum density at an elevated temperature. For the lead-containing piezoelectric ceramics (PbTiO₃, PZT, PLZT), sintering is performed in sealed crucibles with an optimized PbO atmosphere. This is because lead loss occurs in these ceramics above 800°C. As mentioned earlier (Figure 6.39), the randomness of the ceramic grains yields a nonpiezoelectric material. By electroding the ceramic and applying a strong dc electric field at high temperature, the ceramic is poled. At this point, the piezoelectric ceramic is ready for final finishing and characterization.

Piezoelectric Polymers

The piezoelectric behavior of polymers was first reported in 1969 [8]. The behavior results from the crystalline regions formed in these polymers during solidification from the melt. When the polymer is drawn, or stretched, the regions become polar, and can be poled by applying a high electric field. The most widely known piezoelectric polymers are polyvinylidene fluoride [9, 10], also known as PVDF, polyvinylidene fluoride — trifluoroethylene copolymer, or P(VDF-TrFE) [9, 10], and odd-number nylons, such as Nylon-11 [11].

The electromechanical properties of piezoelectric polymers are significantly lower than those of piezoelectric ceramics. The d_{33} values for PVDF and P(VDF-TrFE) are approximately $33 (\times 10^{-12} \text{ C N}^{-1})$, and the dielectric constant ϵ is in the range 6 to 12 [12, 13]. They both have a coupling coefficient (k) of 0.20, and a Curie point (T_0) of approximately 100°C. For Nylon-11, ϵ is around 2 [11], while k is approximately 0.11.

Piezoelectric Ceramic/Polymer Composites

As mentioned above, a number of single-crystal, ceramic, and polymer materials exhibit piezoelectric behavior. In addition to the monolithic materials, composites of piezoelectric ceramics with polymers have also been formed. Table 6.7 [14] summarizes the advantages and disadvantages of each type of material. Ceramics are less expensive and easier to fabricate than polymers or composites. They also have relatively high dielectric constants and good electromechanical coupling. However, they have high acoustic impedance, and are therefore a poor acoustic match to water, the media through which it is typically transmitting or receiving a signal. Also, since they are stiff and brittle, monolithic ceramics cannot be formed onto curved surfaces, limiting design flexibility in the transducer. Finally, they have a high degree of noise associated with their resonant modes. Piezoelectric polymers are acoustically well matched to water, are very flexible, and have few spurious modes. However, applications for these polymers are limited by their low electromechanical coupling, low dielectric constant, and high cost of fabrication. Piezoelectric ceramic/polymer composites have shown superior properties when compared to single-phase materials. As shown in Table 6.7, they combine high coupling, low impedance, few spurious modes, and an intermediate dielectric constant. In addition, they are flexible and moderately priced.

TABLE 6.8 Suppliers of Piezoelectric Materials and Sensors

Name	Address	Ceramic	Polymer	Composite
AMP Sensors	950 Forge Ave. Morristown, PA 19403 Phone: (610) 650-1500 Fax: (610) 650-1509		X	
Krautkramer Branson	50 Industrial Park Rd. Lewistown, PA 17044 Phone: (717) 242-0327 Fax: (717) 242-2606			X
Materials Systems, Inc.	531 Great Road Littleton, MA 01460 Phone: (508) 486-0404 Fax: (508) 486-0706			X
Morgan Matroc, Inc.	232 Forbes Rd. Bedford, OH 44146 Phone: (216) 232-8600 Fax: (216) 232-8731	X		
Sensor Technology Ltd.	20 Stewart Rd. P.O. Box 97 Collingwood, Ontario, Canada Phone: +1 (705) 444-1440 Fax: +1 (705) 444-6787	X		
Staveley Sensors, Inc.	91 Prestige Park Circle East Hartford, CT 06108 Phone: (860) 289-5428 Fax: (860) 289-3189			X
Valpey-Fisher Corporation	75 South Street Hopkinton, MA 01748 Phone: (508) 435-6831 Fax: (508) 435-5289		X	
Vermont U.S.A.	6288 SR 103 North Bldg. 37 Lewistown, PA 17044 Phone: (717) 248-6838 Fax: (717) 248-7066	X		
TRS Ceramics, Inc.	2820 E. College Ave. State College, PA 16801 Phone: (814) 238-7485 Fax: (814) 238-7539	X		

Suppliers of Piezoelectric Materials

Table 6.8 lists a number of the suppliers of piezoelectric materials, their addresses, and whether they supply piezoelectric ceramic, polymers, or composites. Most of them tailor the material to specific applications.

Measurements of Piezoelectric Effect

Different means have been proposed to characterize the piezoelectric properties of materials. The resonance technique involves the measurement of the characteristic frequencies when the suitably shaped specimen (usually ceramic) is driven by a sinusoidally varying electric field. To a first approximation, the behavior of a piezoelectric sample close to its fundamental resonance frequency can be represented by an equivalent circuit as shown in Figure 6.42(a). The schematic behavior of the reactance of the sample as a function of frequency is represented in Figure 6.42(b). By measuring the characteristic frequencies

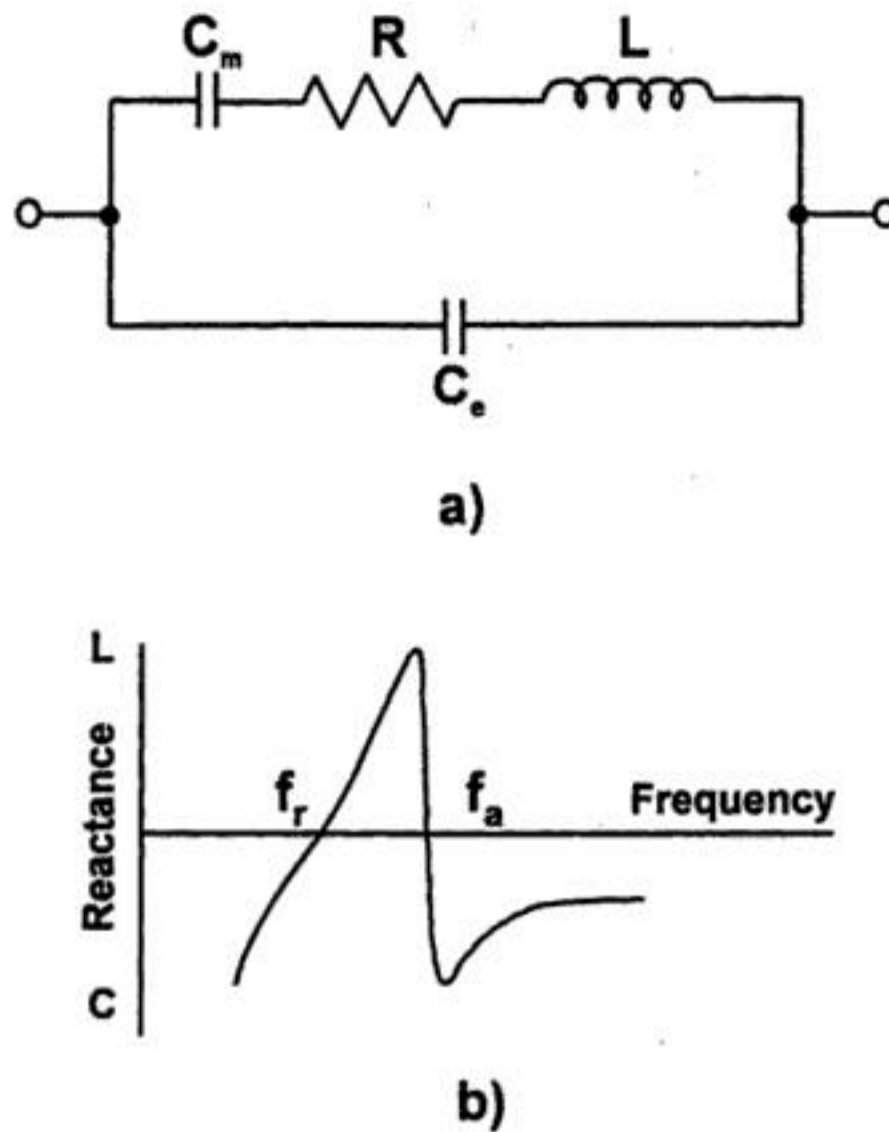


FIGURE 6.42 (a) Equivalent circuit of the piezoelectric sample near its fundamental electromechanical resonance (top branch represents the mechanical part and bottom branch represents the electrical part of the circuit); (b) electrical reactance of the sample as a function of frequency.

of the sample, the material constants including piezoelectric coefficients can be calculated. The equations used for the calculations of the electromechanical properties are described in the IEEE Standard on piezoelectricity [15]. The simplest example of piezoelectric measurements by resonance technique relates to a piezoelectric ceramic rod (typically 6 mm in diameter and 15 mm long) poled along its length. It can be shown that the coupling coefficient k_{33} is expressed as a function of the series and parallel resonance frequencies, f_s and f_p , respectively:

$$k_{33}^2 = \frac{\pi}{2} \frac{f_s}{f_p} \tan \left(\frac{\pi}{2} \frac{f_p - f_s}{f_p} \right) \quad (6.63)$$

The longitudinal piezoelectric coefficient d_{33} is calculated using k_{33} , elastic compliance s_{33}^E and low-frequency dielectric constant ϵ_{33}^X :

$$d_{33} = k_{33} \left(\epsilon_{33}^X s_{33}^E \right)^{1/2} \quad (6.64)$$

Similarly, other electromechanical coupling coefficients and piezoelectric moduli can be derived using different vibration modes of the sample. The disadvantage of the resonance technique is that measurements are limited to the specific frequencies determined by the electromechanical resonance. It is used mostly for the rapid evaluation of the piezoelectric properties of ceramic samples whose dimensions can be easily adjusted for specific resonance conditions.

Subresonance techniques are frequently used to evaluate piezoelectric properties of materials at frequencies much lower than the fundamental resonance frequency of the sample. They include both the measurement of piezoelectric charge under the action of external mechanical force (direct effect) and the

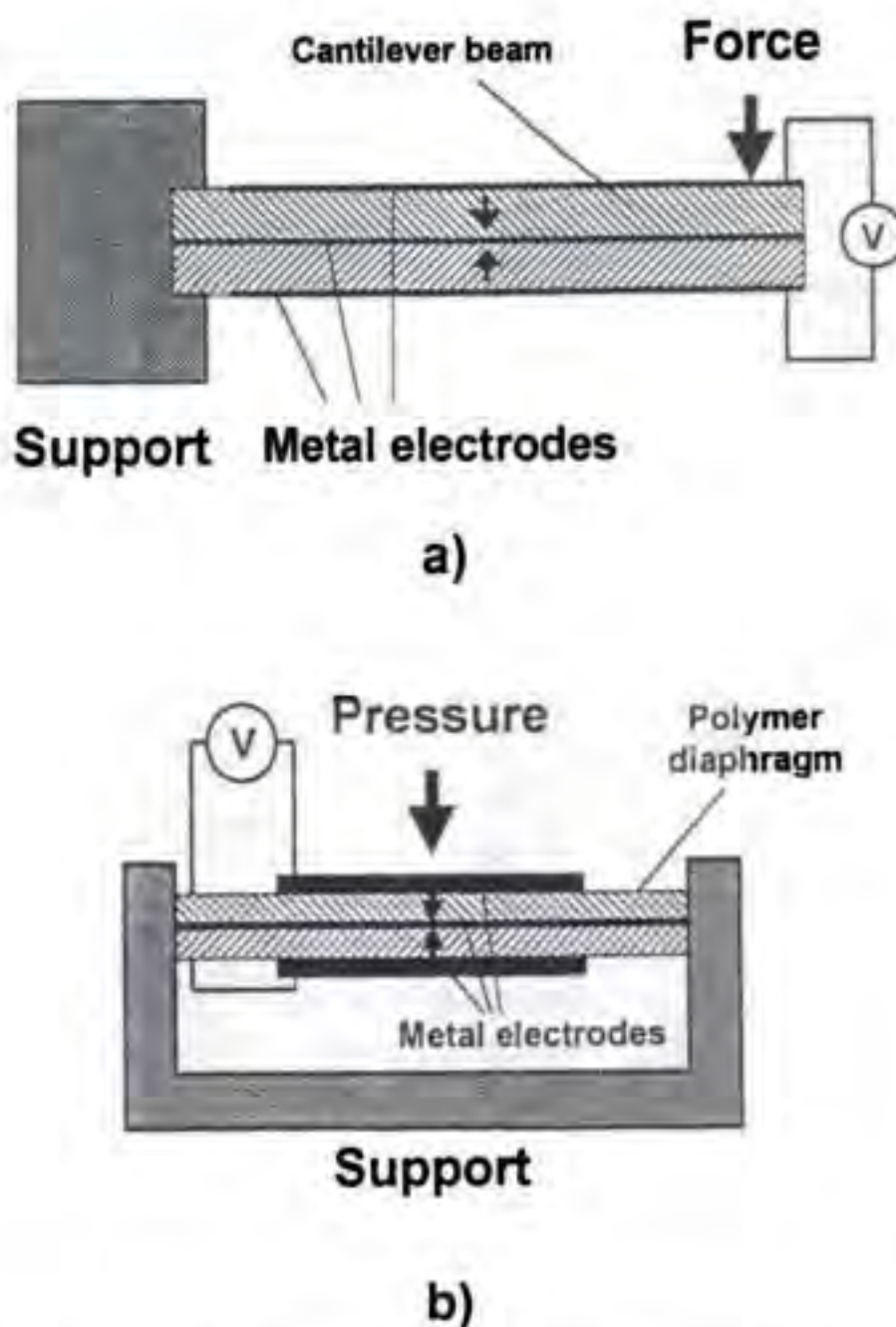


FIGURE 6.45 Schematic designs of the displacement sensor based on piezoelectric ceramic (a) and of the pressure sensor based on piezoelectric polymer film (b). Arrows indicate the directions of ferroelectric polarization in the piezoelectric material.

transducer, using the converse piezoelectric properties of the material, changes its dimensions and sends an acoustic signal into a medium. Active mode applications include nondestructive evaluation, fish/depth finders, ink jet printers, micropositioners, micropumps, and medical ultrasonic imaging. Often, the same transducer is used for both sensor and actuator functions.

Two examples of piezoelectric sensors are given below. The first example is the ceramic transducer, which relates the deformation of the piezoelectric sensor to the output voltage via direct piezoelectric effect. Piezoceramics have high Young's moduli; therefore, large forces are required to generate strains in the transducer to produce measurable electric response. Compliance of the piezoelectric sensor can be greatly enhanced by making long strips or thin plates of the material and mounting them as cantilevers or diaphragms. Displacement of the cantilever end will result in a beam bending, leading to the mechanical stress in the piezoelectric material and the electric charge on the electrodes. A common configuration of the piezoelectric bender is shown in Figure 6.45(a). Two beams poled in opposite directions are cemented together with one common electrode in the middle and two electrodes on the outer surfaces. Bending of such a bimorph will cause the upper beam to stretch and the lower beam to compress, resulting in a piezoelectric charge of the same polarity for two beams connected in series. To the first approximation, the charge Q appearing on the electrodes is proportional to the displacement Δl of the end of the bimorph via Equation 6.68 [18]:

$$Q = \frac{3}{8} \frac{Hw}{L} e_{31} \Delta l, \tag{6.68}$$

where H , w , and L are the thickness, the width, and the length of the bimorph, respectively, and e_{31} is the transverse piezoelectric coefficient relating electric polarization and strain in a deformed piezoelectric material. The charge can be measured either by the voltage amplifier (Figure 6.43) or by the charge amplifier (Figure 6.44).

In certain applications, the parameters of piezoelectric sensors can be improved by using ferroelectric polymers instead of single crystals and piezoceramics. Although the electromechanical properties of polymers are inferior to those of piezoelectric ceramics, their low dielectric constant offers the higher voltage response since they possess higher g piezoelectric coefficients. Also, the polymers are more mechanically robust and can be made in the form of thin layers (down to several micrometers). An example using the polymer bimorph as a pressure sensor is shown in Figure 6.45(b). A circular diaphragm composed of two oppositely poled polymer films is clamped along its edges to a rigid surround, forming a microphone. The voltage appearing on the electrodes is proportional to the applied pressure p by Equation 6.69 [19]:

$$V = \frac{3}{16} \frac{d_{31}}{\epsilon_{33}} \frac{D^2}{h} (1-\nu)p \quad (6.69)$$

where D and h are the diameter and thickness of the diaphragm, respectively, and ν is the Poisson ratio. The high d_{31}/ϵ_{33} value for polymer sensors is advantageous to obtain higher voltage response. According to Equation 6.66, this advantage can be realized only if the high input impedance amplifier is used in close proximity to the transducer to reduce the influence of the connecting cables.

Defining Terms

Piezoelectric transducer: Device that converts the input electrical energy into mechanical energy and vice versa via piezoelectric effect.

Coupling coefficients: Materials constants that describe an ability of piezoelectric materials to convert electrical energy into mechanical energy and vice versa.

Piezoelectric coefficients: Materials constants that are used to describe the linear coupling between electrical and mechanical parameters of the piezoelectric.

Ferroelectrics: Subgroup of piezoelectric materials possessing a net dipole moment (ferroelectric polarization) that can be reversed by the application of sufficiently high electric field.

Poling: Process of aligning the ferroelectric polarization along a unique (poling) direction.

Piezoelectric composites: Materials containing two or more components with different piezoelectric properties.

Charge amplifier: An operational amplifier used to convert the input charge into output voltage by means of the capacitor in the feedback loop.

References

1. J. F. Nye, *Physical Properties of Crystals*, Oxford: Oxford University Press, 1985.
2. Y. Xu, *Ferroelectric Materials and Their Applications*, Amsterdam: North-Holland, 1991.
3. L. E. Cross, Ferroelectric ceramics: tailoring properties for specific applications, In N. Setter and E. L. Colla (ed.), *Ferroelectric Ceramics: Tutorial Reviews, Theory, Processing, and Applications*, Basel: Birkhauser, 1993.
4. B. Jaffe, W. R. Cook, Jr., and H. Jaffe, *Piezoelectric Ceramics*, Marietta, OH: R. A. N., 1971.
5. *The User's Guide to Ultrasound & Optical Products*, Hopkinton, MA: Valpey-Fisher Corporation, 1996.
6. S.-E. Park and T. R. Shrout, Relaxor based ferroelectric single crystals with high piezoelectric performance, *Proc. of the 8th US-Japan Seminar on Dielectric and Piezoelectric Ceramics*: October 15-18, Plymouth, MA, 1997, 235.

7. *Piezoelectric Products*, Sensor Technology Limited, Collingwood, Ontario, Canada, 1991.
8. H. Kawai, The piezoelectricity of poly(vinylidene fluoride), *Japan. J. Appl. Phys.*, 8, 975, 1969.
9. L. F. Brown, Ferroelectric polymers: Current and future ultrasonic applications, *Proc. 1992 IEEE Ultrasonics Symposium*: IEEE, New York, 1992, 539.
10. T. Furukawa, Recent advances in ferroelectric polymers, *Ferroelectrics*, 104, 229, 1990.
11. L. F. Brown, J. I. Scheinbeim, and B. A. Newman, High frequency dielectric and electromechanical properties of ferroelectric nylons, *Proc. 1994 IEEE Ultrasonics Symposium*: IEEE, New York, 1995, 337.
12. *Properties of Raytheon Polyvinylidene Fluoride (PVDF)*, Raytheon Research Division, Lexington, MA, 1990.
13. *Standard and Custom Piezo Film Components*, Atochem Sensors Inc., Valley Forge, PA, 1991.
14. T. R. Gururaja, Piezoelectric transducers for medical ultrasonic imaging, *Amer. Ceram. Soc. Bull.*, 73, 50, 1994.
15. IEEE Standards on Piezoelectricity, *IEEE Std. 176*, 1978.
16. W. Y. Pan and L. E. Cross, A sensitive double beam laser interferometer for studying high-frequency piezoelectric and electrostrictive strains, *Rev. Sci. Instrum.*, 60, 2701, 1989.
17. A. L. Kholkin, Ch. Wuethrich, D. V. Taylor, and N. Setter, Interferometric measurements of electric field-induced displacements in piezoelectric thin films, *Rev. Sci. Instrum.*, 67, 1935, 1996.
18. A. J. Moulson and J. M. Herbert, *Electroceramics: Materials, Properties, Applications*, London: Chapman and Hall, 1990.
19. J. M. Herbert, *Ferroelectric Transducers and Sensors*, New York: Gordon and Breach, 1982.

6.5 Laser Interferometer Displacement Sensors

Bernhard Günther Zagar

In the past few years, very high precision, numerically controlled machine tools have been developed. To achieve the potential precision of these tools, length and displacement measurements whose resolution exceeds the least significant digit of the tool must be made. The measurement equipment typically would not rely on mechanical scales.

Laser interferometers compare the changes in optical path length to the wavelength of light, which can be chosen from atomic constants that can be determined with very little uncertainty.

In 1983, there was a redefinition of the meter [1] that was previously defined in 1960. The old definition was based on the wavelength of a certain radiation (the krypton-86 standard) that could not be realized to better than 4 parts in 10^9 . The new definition, being based on frequency but not related to any particular radiation, opened the way to significant improvements in the precision with which the meter can be realized. As recommended in resolution 2 for the practical realization of the meter, the wavelength in vacuum λ_v of a plane electromagnetic wave of frequency f is $\lambda_v = c/f$, where c is the speed of light in vacuum, $c = 299,792,458 \text{ m s}^{-1}$ exactly. This way, the wavelength is related to *frequency* and *time*, which can be measured with the highest precision of all units within the *Système International (SI)*.

In order to be independent of any environmental parameters, the meter is defined using the speed of light in a vacuum. However, interferometers usually must operate in ambient air. Thus, environmental parameters that influence the speed of light in a particular medium (air) will affect and degrade the precision of the measurement.

Three major factors limit the absolute accuracy attainable with laser interferometers operating in ambient air: (1) the uncertainties of the vacuum wavelength, λ_v , of the laser source; (2) the uncertainty of the refractive index of the ambient air; and (3) the least count resolution of the interferometer.

This chapter section is organized as follows. First, some basic laser principles are detailed, including ways to stabilize the vacuum wavelength of the laser. The effect most often used to stabilize lasers in commercial interferometers is the Zeeman effect, which yields relative uncertainties of 10^{-8} .

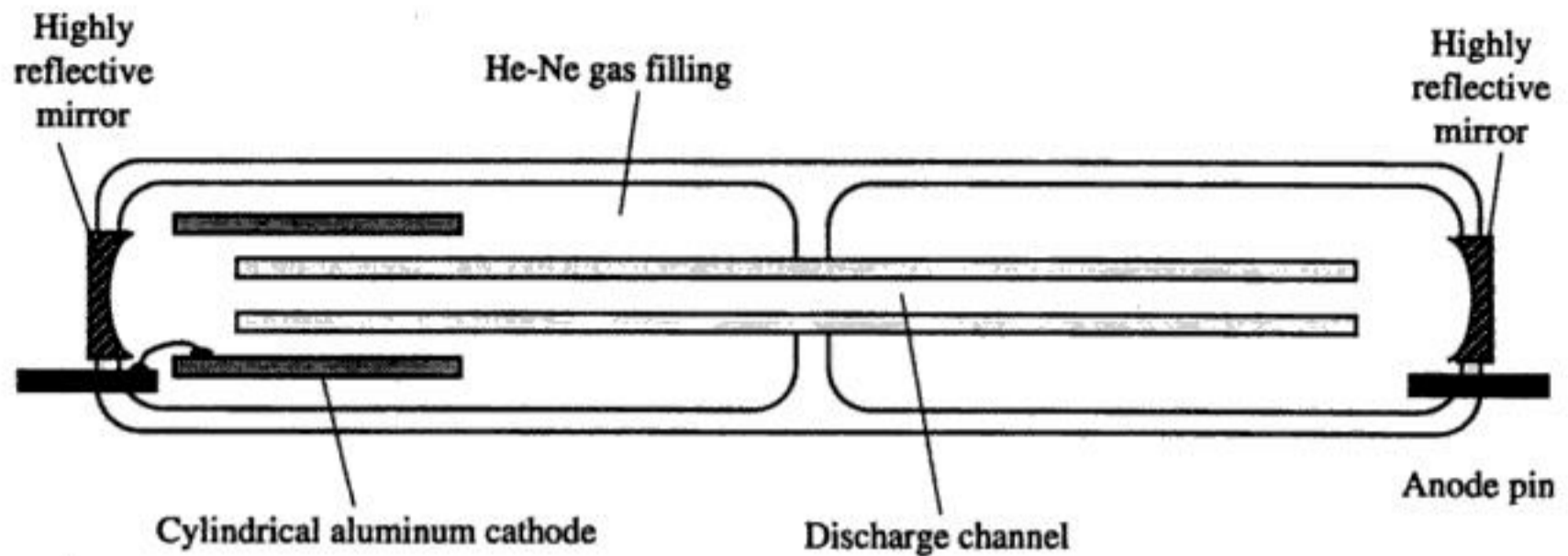


FIGURE 6.46 Schematics of the helium–neon laser (Reprinted with permission of University Science Books [3]).

Second, the refractive index of air as another major factor limiting the attainable accuracy of laser interferometers operated in air is addressed. It is shown that it cannot be determined currently with uncertainty better than 5×10^{-8} .

And finally, the chapter section describes the most widely used Michelson interferometer and two of its variants for long-travel length measurement and gives their resolution.

Helium–Neon Laser

In order to attain the best possible accuracy, great care must be taken to ensure the highest wavelength stability of the light source. Almost all interferometric dimensional gages utilize a helium–neon laser because it has proven reliable, its emitted wavelength is in the visible range at about 633 nm, and it can be stabilized sufficiently well utilizing the Zeeman effect and to an even higher degree with the use of a very well-defined iodine absorption line also at ≈ 633 nm [1, 2].

The helium–neon laser consists of a discharge tube as shown in Figure 6.46 [3] filled with the single-isotope gases helium (He^3) at a partial pressure of ≈ 105 Pa and neon (Ne^{20}) with a partial pressure of ≈ 13 Pa. It is pumped electrically using a voltage on the order of kilovolts with a current of a few milliamperes to excite both helium and neon atoms. Since the helium gas is the majority component, it dominates the discharge properties of the laser tube. Neutral helium atoms collide with free electrons that are accelerated by the axial voltage and become excited and remain in two rather long-lived metastable states. These are close enough to the energy levels of certain excited states of neon atoms so that collisional energy transfer can take place between these two groups of atoms. Excited helium atoms may drop down to the ground state, while simultaneously neon atoms take up almost exactly the same amount of energy. Therefore helium only serves to excite neon atoms, they do not contribute to the emission of light. The excited neon atoms remain in the excited state for a rather long period of time (on the order of 10^{-3} s). They return to lower energetic levels by stimulated emission of highly coherent light. This stimulated emission comes into effect when light emitted by some neon atoms also prompts other atoms to emit. The mirrors of the laser cavity, by reflecting most of the incident light cause the light to traverse multiple paths through the active laser volume, thereby greatly amplifying the light if the cavity length L is an integer multiple m of half the wavelength λ .

$$L = m \frac{\lambda}{2} \quad (6.70)$$

The emitted light is fairly monochromatic, but still has some finite spectral linewidth determined by the random emissions of Ne^{20} .

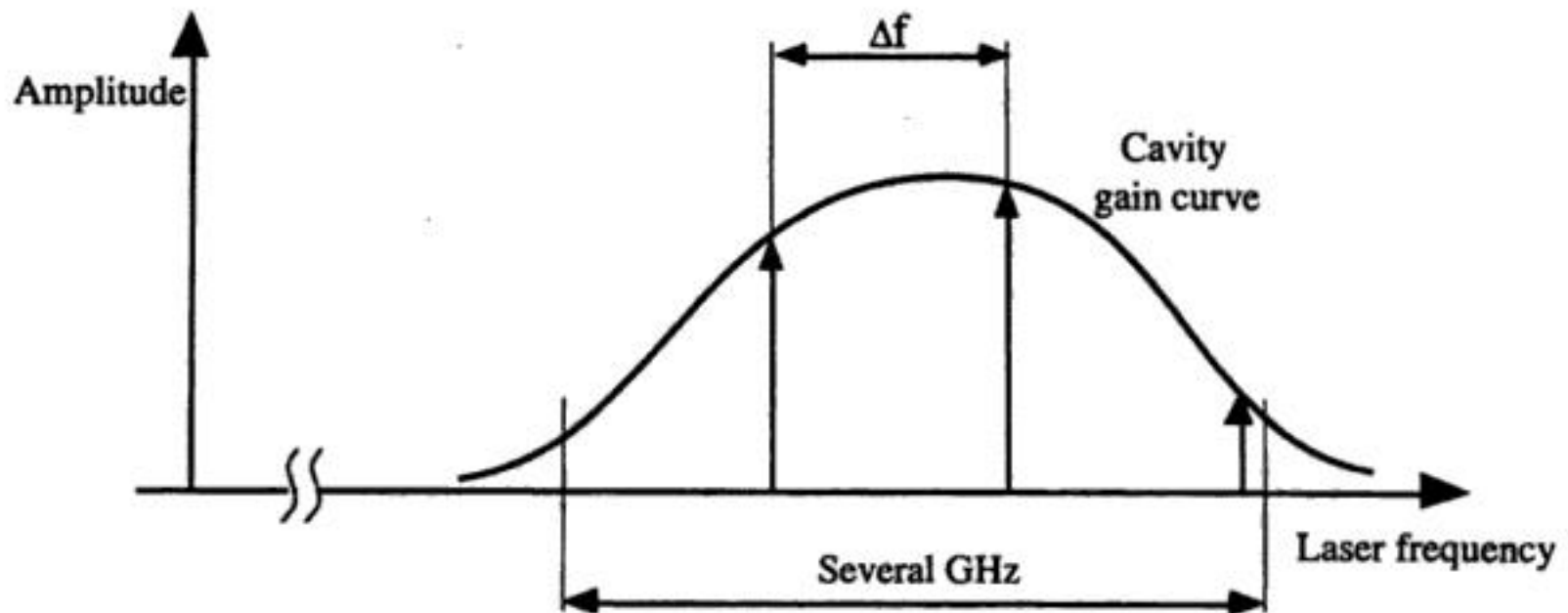


FIGURE 6.47 A He–Ne laser can have multiple resonating modes (shown for ≈ 20 cm cavity length) (Reprinted with permission of University Science Books [3]).

Brewster angle [4] end windows of the discharge tube transmit light of the proper linear polarization if desired. The end mirrors have to be carefully polished. Their curvature radii have to satisfy the condition for stability. They have wavelength–selective dielectric coatings of very high reflectivity sometimes exceeding 99%.

Unless special precautions are taken, a He–Ne laser will emit several axial modes as shown schematically in Figure 6.47, resulting in a beat frequency that limits the temporal coherence and renders the laser unsuitable for interferometric purposes. Also, due to thermal expansion of the laser tube, the end mirrors will change their relative distance, thereby effectively tuning the wavelength within the linewidth of the gain curve.

Another important property of a high-quality He–Ne laser is its Gaussian cross-sectional profile, which is maintained along a propagating wave, i.e., fundamental lateral mode. It is also a necessary condition for the wavefronts to remain quasiplanar.

Frequency Stabilization of He–Ne Lasers

The resonant frequency (and the wavelength λ) of the laser is determined in part by the distance between the two end mirrors and also by the refractive index n_M of the active medium (the He–Ne mixture). Since the linewidth of the gain profile of the active medium is usually in the gigahertz range, multiple axial modes can resonate in the cavity as is shown in Figure 6.47. The frequency difference Δf between two adjacent longitudinal modes is the *free spectral range* (FSR), which is given by Equation 6.71 and depends only on the cavity length, L , the refractive index of the active medium, n_M , and the speed of light in vacuum, c .

$$\Delta f = \frac{c}{2n_M L} \quad (6.71)$$

Due to thermal expansion of the laser cavity and/or thermally induced change in refractive index of the active medium, all resonating laser modes will move within the envelope of the gain profile. The effort undertaken in stabilizing the laser wavelength or equivalently stabilizing its frequency is aimed at locking the modes with respect to the gain profile and reducing the number of modes resonating simultaneously.

Longitudinal Zeeman Effect

One of the most often used effects in stabilizing the frequency of a He–Ne laser for distance measurements is the Zeeman effect [3, 5].

$$A = 78.603 \left[1 + 0.540(D - 0.0003) \right] \frac{P}{TZ} \times 10^{-8}$$

$$B = (0.00042066 f_E e_s H) \times 10^{-8} \quad (6.73)$$

In Equation 6.73, P is the atmospheric pressure in pascals, T is the absolute temperature in kelvin, H is the relative humidity in %, and D is the concentration of CO_2 in percent by volume. There are three additional factors in Jones' formulation that take the nonideal behavior of moist air as compared to an ideal gas into account. They are Z , a compressibility factor that reflects the nonideality of the air–water vapor mixture and which, for air containing reasonable amounts of CO_2 at a temperature between 15°C and 28°C and pressure of between 7×10^4 Pa and 11×10^4 Pa, lies in the range between 0.99949 and 0.99979. f_E is an enhancement factor that expresses the fact that the effective saturation vapor pressure of water in air is greater than the saturation vapor pressure e_s . For the pressure and temperature ranges given above, f_E is bounded between 1.0030 and 1.0046 [16]. e_s is the saturation vapor pressure over a plane surface of pure liquid water and according to Jones is about 1705 Pa at a temperature of 15.0°C and about 3779 Pa for 28.0°C . Tables of Z , f_E and e_s are included in the Appendix of Jones' paper [16].

Table 6.10 gives an overview of the changes in environmental parameters that would cause a relative index change of 10^{-7} .

Edlén [8] and Jones [16] estimate that their empirical expressions for the dependency of the refractive index of air on the listed parameters has an absolute uncertainty of 5×10^{-8} .

Besides this fundamental limitation, there are some practical considerations that must be taken into account regarding the precision with which the environmental parameters can be measured. Estler [17] states that atmospheric pressure P can currently be determined with an uncertainty of ≈ 2.7 Pa, which can be assumed to be constant for the entire optical path of the interferometer if it is oriented horizontally. Please note that at sea level, the pressure gradient is ≈ -13 Pa m^{-1} , resulting in a pressure-induced change in n_A of $3.4 \times 10^{-8} \text{ m}^{-1}$ if the measuring equipment is not kept level.

In an exceptionally well-controlled laboratory environment where special care is devoted to keep temperature gradients from affecting the refractive index along the optical path as much as possible, uncertainties of the temperature measurement can be as low as 0.01°C according to [17]. Humidity measured with high accuracy dew-point hygrometers can have uncertainties down to 0.5%. Changes in carbon dioxide concentrations have to be very significant (20% of the natural concentration) to cause a $\Delta n/n$ of 10^{-8} .

Michelson Interferometer

The basis for most interferometers used in interferometric dimensional gages is the classical Michelson interferometer [4] which is shown in Figure 6.49. The coherent monochromatic light of a wavelength-stabilized He–Ne laser is incident onto a beam splitter which splits the light into two equally intense beams (1) and (2).

TABLE 6.10 Parameters of Standard Air and Their Deviation to Cause a $\Delta n/n$ of 10^{-7}

Parameter	Standard value	Variation for $\Delta n/n = +1 \times 10^{-7}$
Pressure P	101.3 kPa	+37.3 Pa
Temperature T	20.0°C	-0.1°C
Humidity H	40%	-10.0%
CO_2 concentration	350 ppm	+670 ppm

Reprinted with permission of *J. Applied Optics* [17].

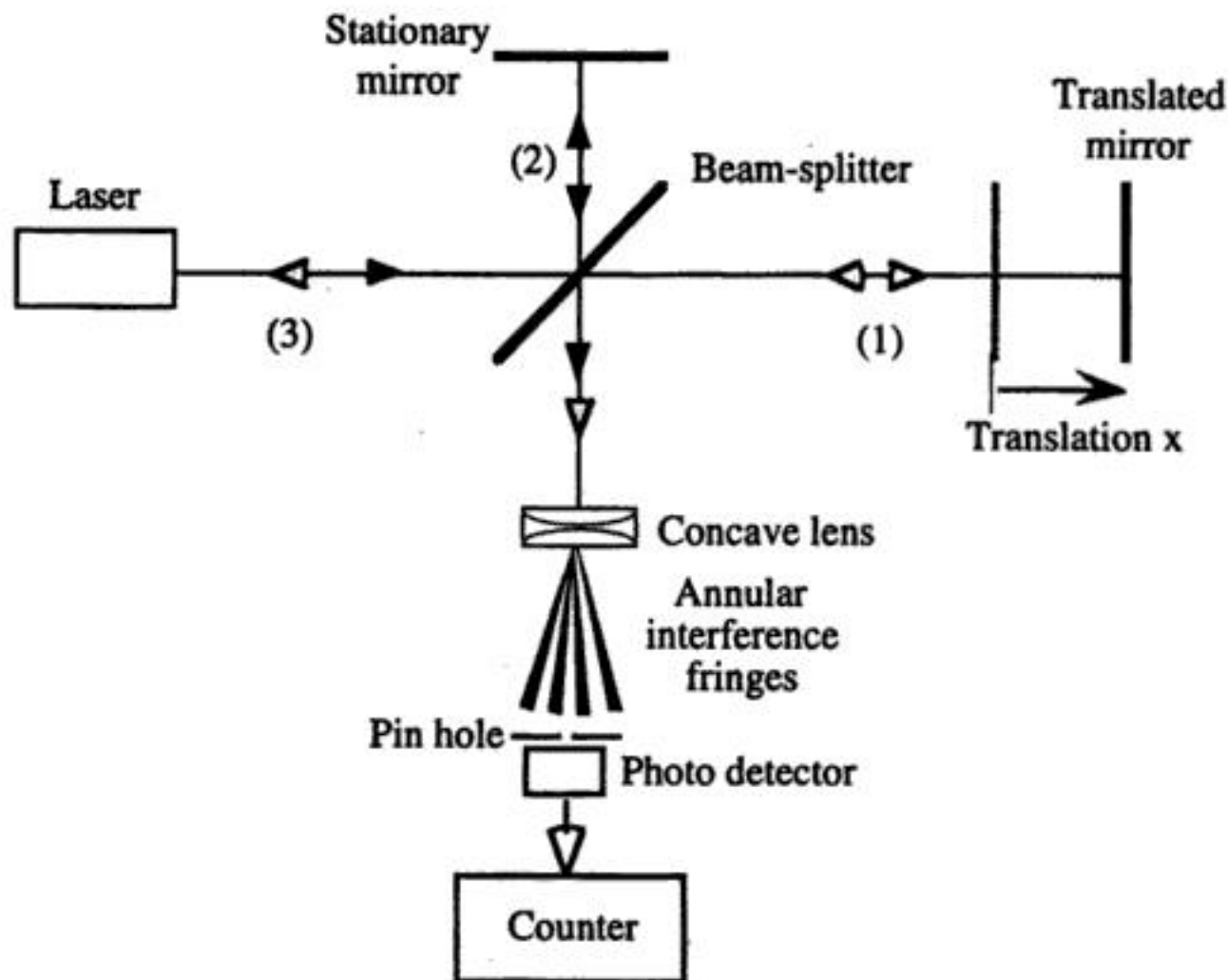


FIGURE 6.49 Schematics of the basic Michelson interferometer.

They are reflected off of both the stationary and the translatable mirror whose displacement x is to be measured, and recombined at the splitter, where they are redirected toward a concave lens. Due to the coherence of the laser light, the wavefronts have a well-defined phase relation with respect to each other. This phase is determined by the difference between the optical path lengths of the two beams in arms 1 and 2. If this path difference is continuously changed by translating one of the mirrors, a sinusoidal intensity variation can be observed at a fixed location in space behind the lens used to introduce a beam divergence, resulting in an annular fringe pattern. The pinhole is used to define an exact location of observation and the photodetector picks up the varying intensity for further processing. In the most basic signal processing setup, the number of bright and dark cycles are fed into a counter, which then counts changes in optical path length in integer multiples of $\lambda_A/2$. More sophisticated signal processing not only counts cycles but also determines relative phase changes in the sinusoidal varying intensity so that resolutions of $\lambda_A/512$ can ultimately be achieved.

When moving the mirror, one must guarantee a smooth motion without backward jitter of the mirror carriage to avoid double counts of interference fringes. Very high-quality linear bearings (such as air bearings) are necessary to accomplish just that.

As can be seen in Figure 6.49, the light reflected off both mirrors essentially retraces its own path and is at least partially incident onto the active volume of the laser source (3), thereby forming an external laser resonator which is able to detune the laser, effectively modulating its output power as well as its wavelength. To avoid this effect, commercial versions of Michelson interferometers employ corner-cube reflectors instead of plane mirrors as well as optical isolators, as shown in Figures 6.50 and 6.51. Some authors, however, report using optical arrangements that utilize this effect in conjunction with laser diodes to realize low-cost, short-travel displacement sensors [18, 19]. These setups will not be discussed here, however.

Two-Frequency Heterodyne Interferometer

Figure 6.50 shows the commercially available two-frequency Michelson interferometer operating with a Zeeman-stabilized He-Ne laser source. This laser emits two longitudinal modes with frequencies f_1 and f_2 that are both circularly polarized in opposite directions. By passing the modes through a quarter-wave

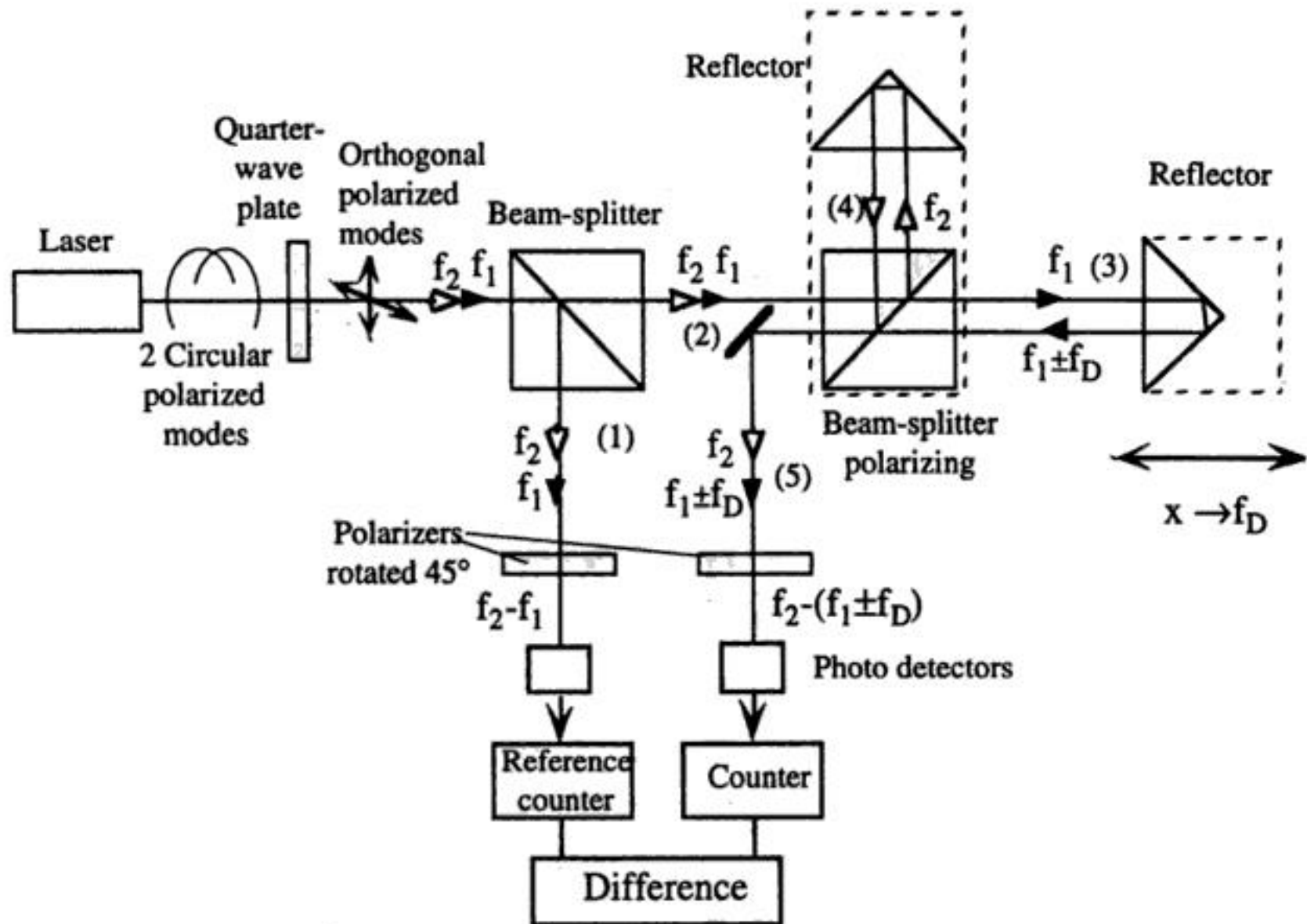


FIGURE 6.50 Two-frequency heterodyne interferometer (Courtesy of Spindler & Hoyer Inc.).

plate, two orthogonal linearly polarized waves are generated. Both are split by a nonpolarizing beam splitter. There is a polarizer located in arm 1 of that splitter, which is rotated by 45° with respect to both polarized waves impinging on it, thus effectively allowing them to interfere behind it yielding a difference frequency of $f_2 - f_1$ that is picked up by a photodetector and counted by a reference counter (frequency difference typically 1.5 MHz).

The orthogonal polarized waves in 2 are further split by a polarizing splitter. Spectral component $f_1 < f_2$ is transmitted into measuring arm 3 and frequency component f_2 is reflected into reference arm 4 of the interferometer. Due to the velocity v of the reflector in arm 3 resulting in a displacement x , the frequency f_1 is Doppler-shifted by f_D (Equation 6.74). Movement of the reflector toward the interferometer results in a positive Doppler frequency $f_D > 0$. After recombining both waves from 3 and 4 in the beam splitter again, they are sent through a polarizer in arm 5 that also is rotated by 45° with respect to the direction of polarization of both recombined waves, thereby allowing them to interfere, yielding a difference frequency of $f_2 - f_1 - f_D$ at the location of the photodetector, which is counted by a second counter. By continuously forming the difference of both counts, the measurand (the displacement of x in multiples of $\lambda_A/2$) is calculated.

With this type of interferometer, the direction of motion is given by the sign of the resulting count. One disadvantage of the two-mode heterodyne interferometer is its limited dynamic range for the velocity v of the reflector moving toward the interferometer, since the Doppler frequency f_D , given by:

$$f_D = \frac{2}{\lambda_A} v \quad (6.74)$$

is bound to be less than the initial frequency difference between f_2 and f_1 for stationary reflectors. Given a typical Zeeman effect-induced frequency difference f_2 of 1.5 MHz, the velocity v is therefore bound to be less than:

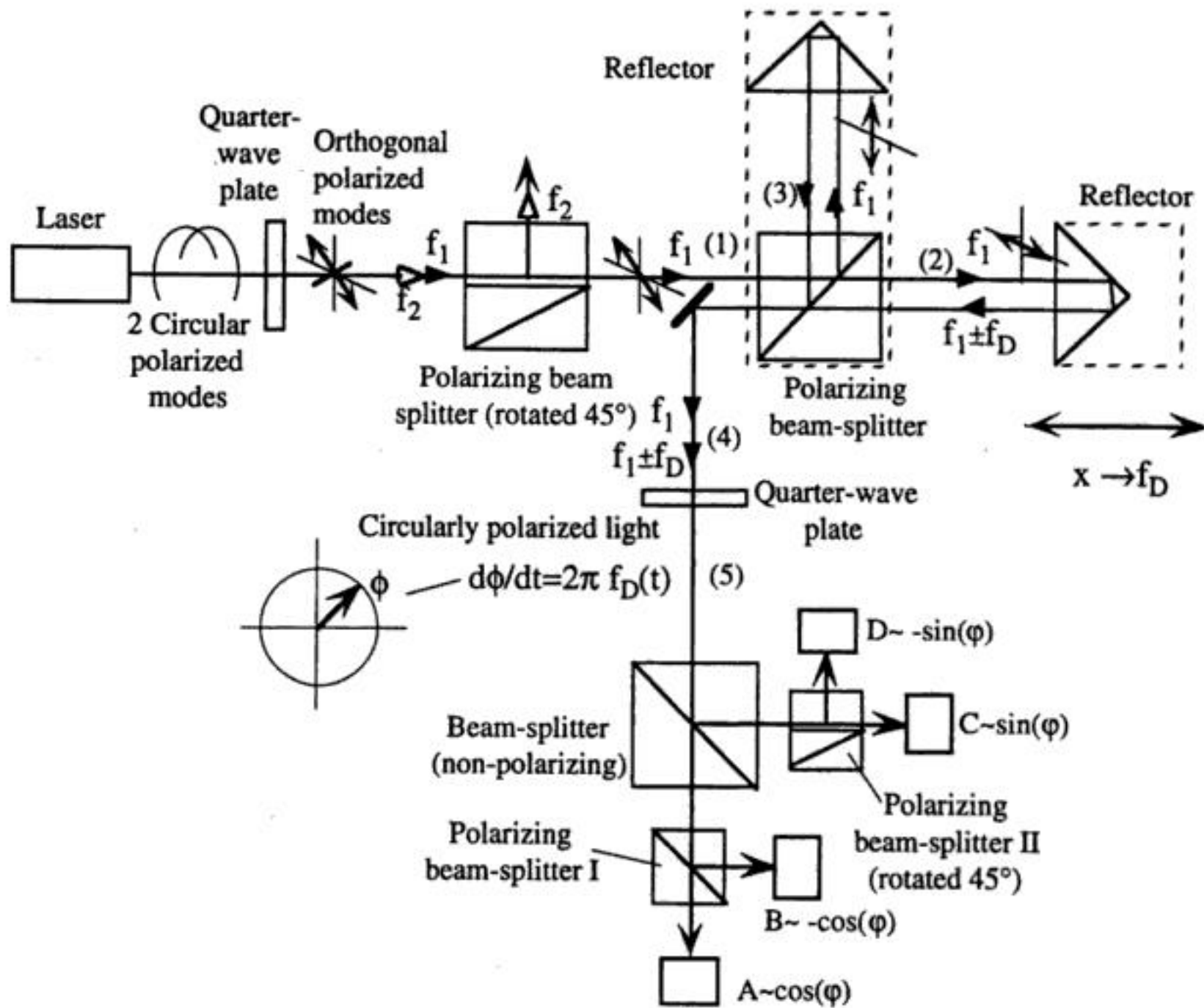


FIGURE 6.51 Single-mode homodyne interferometer (Courtesy of Spindler & Hoyer Inc.).

$$v < \frac{f_2}{2} \lambda_A = 0.474 \text{ m s}^{-1} \tag{6.75}$$

There is no such bound if the reflector is traveling away from the interferometer. By electronically interpolating the output signal of the photodetector, subwavelength resolution can be obtained [22].

Single-Mode Homodyne Interferometer

An interferometer setup that has no limitation on the maximum velocity in the above sense is the single-mode homodyne interferometer shown in Figure 6.51.

As in the two-frequency heterodyne interferometer, a Zeeman effect-stabilized laser source that emits two frequency-displaced circularly polarized axial modes is usually used. After passing through a quarter-wave plate, two orthogonal polarized waves are generated, only one of which (f_1) is further used. The other (f_2) is reflected out of the optical path by an appropriately oriented polarizing beam splitter. The plane of polarization in arm 1 of the interferometer is tilted by 45° with respect to the plane defined by the two arms 2 and 3. The second polarizing beam splitter will transmit one horizontally oriented component into the measuring arm 2 and reflect a vertically oriented component into the reference arm 3 of the interferometer. The two arms of the interferometer maintain their orthogonal polarizations. The frequency of the wavefront in arm 2 is shifted by the Doppler effect (Equation 6.74) due to the motion of the reflector. The light reflected by the two triple mirrors is recombined in the polarizing beam splitter and redirected by a mirror. Since in this particular optical setup, the polarization states of the two beams in the two arms of the interferometer are orthogonal, there is no interference after the redirecting mirror

in arm 4 as was the case with the basic Michelson interferometer setup (Figure 6.49). After passing a quarter-wave plate at 45° , two opposite circular polarized waves (one with frequency f_1 , the other with frequency $f_1 \pm f_D$) are generated and can be described by a rotating phasor (characterized by $\Phi(t)$) with constant amplitude whose rate of rotation is dependent on the Doppler frequency (in arm 5). Amplitude fluctuations can be observed at the photodetectors A–D after this phasor has passed polarizers, which it does after being split by a nonpolarizing beam splitter.

The output of an interferometer has the general form:

$$I(t) = I_0(t) \frac{1}{2} \left[1 + \cos(\Phi(t)) \right] \quad (6.76)$$

It is desired to infer $\Phi(t)$ from observation of $I(t)$. Note that $I_0(t)$, which is the intensity of the laser, can also fluctuate with time. The problems encountered with this are (1) the ambiguity in the sign of $\Phi(t)$ and (2) the dependence of the calculated phase on the intensity fluctuations due to the aging of the laser and optical components. The first problem stems from the fact that $\arccos(\dots)$ yields two solutions to Equation 6.76:

$$\Phi(t) = \pm \arccos \left[\frac{2I(t)}{I_0(t)} - 1 \right] \quad (6.77)$$

The sign ambiguity can be resolved by also generating a $\sin(\Phi(t))$ yielding quadrature signals. In order to do so, a second output of the interferometer of the form:

$$I_2(t) = I_0(t) \frac{1}{2} \left[1 + \sin(\Phi(t)) \right] \quad (6.78)$$

is sought. Equations 6.76 and 6.78 will determine $\Phi(t)$ unambiguously only in the region $[0, 2\pi)$, but there is still an ambiguity modulo 2π . The second problem can be dealt with by adding two more outputs of the form:

$$I_3(t) = I_0(t) \frac{1}{2} \left[1 - \cos(\Phi(t)) \right] \quad (6.79)$$

$$I_4(t) = I_0(t) \frac{1}{2} \left[1 - \sin(\Phi(t)) \right] \quad (6.80)$$

Taking Equations 6.76 to 6.79 and 6.78 to 6.80, it is possible to obtain a zero crossing at the linear most sensitive point of inflection of the fringe where the effect of intensity fluctuations on the phase measurement is minimal. The setup of Figure 6.51 attempts to obtain these four outputs. Signal A represents the intensity variations as given in Equation 6.76 and signal B due to the nature of the splitting action is shifted with respect to A by 180° (Equation 6.79). There is another arm to the right of the nonpolarizing beam splitter incorporating the polarizing beam splitter II, which is rotated by 45° with respect to beam splitter I so that the attached detectors C and D are generating the signals defined by Equations 6.78 and 6.80.

Since the Doppler frequency is time dependent according to the velocity v of the measuring reflector, the distance, x , traveled by the reflector up to time T is given by Equation 6.81.

$$x = \int_0^T v(t) dt = \int_0^T f_D(t) \frac{\lambda_\Lambda}{2} dt = \int_0^T \frac{\partial \Phi(t)}{\partial t} \frac{\lambda_\Lambda}{4\pi} dt \quad (6.81)$$

11. R. Revelle, Carbon dioxide and world climate, *Sci. Amer.*, 247(2), 35, 1982.
12. K. P. Birch, F. Reinboth, R. W. Ward, and G. Wilkening, Evaluation of the effect of variations in the refractive index of air upon the uncertainty of industrial length measurement, *Metrologia*, 30(1), 7-14, 1993.
13. K. P. Birch and M. J. Downs, An updated Edlén equation for the refractive index of air, *Metrologia*, 30, 155-162, 1993.
14. K. P. Birch and M. J. Downs, Corrections to the updated Edlén equation for the refractive index of air, *Metrologia*, 31, 315-316, 1994.
15. P. E. Ciddor, Refractive index of air: new equations for the visible and near infrared, *Appl. Optics*, 35(9), 1566-1573, 1996.
16. F. E. Jones, The refractivity of air, *J. National Bureau of Standards*, 86(1), 27-32, 1981.
17. W. T. Estler, High-accuracy displacement interferometry in air, *J. Appl. Optics*, 24(6), 808-815, 1985.
18. J. A. Smith, U. W. Rathe, and C. P. Burger, Lasers with optical feedback as displacement sensors, *Opt. Eng.*, 34(9), 2802-2810, 1995.
19. N. Takahashi, S. Kakuma, and R. Ohaba, Active heterodyne interferometric displacement measurement using optical feedback effects of laser diodes, *Opt. Eng.*, 35, 802-907, 1996.
20. K. Oka, M. Tsukada, and Y. Ohtsuka, Real-time phase demodulator for optical heterodyne detection processes, *Meas. Sci. Technol.*, 2, 106-110, 1991.
21. J. Waller, X. H. Shi, N. C. Altoveros, J. Howard, B. D. Blackwell, and G. B. Warr, Digital interface for quadrature demodulation of interferometer signals, *Rev. Sci. Instrum.*, 66, 1171-1174, 1995.
22. J. A. Smith and C. P. Burger, Digital phase demodulation in heterodyne sensors, *Opt. Eng.*, 34, 2793-2801, 1995.

Appendix to Section 6.5

In the appendix below some companies are listed that manufacture either complete interferometer systems or major components thereof, such as beam splitters, retroreflectors, refractometers, wavemeters, etc. This list is by no means exhaustive. Furthermore, no price information is included because the system cost is too much dependent on the particular choice of system components.

Companies that produce interferometers or significant components.

Manufacturer	Sub-systems	Complete systems	Manufacturer	Sub-systems	Complete systems
Aerotech Inc. 101 Zeta Drive Pittsburgh, PA 15238 Tel: (412) 963-7470	*	*	Oriel Instruments Inc. 250 Long Beach Blvd. Stratford, CT 06497-0872 Tel: (203) 380-4364	—	*
Burleigh Inc. Burleigh Park Fishers, NY 14453-0755 Tel: (716) 924-9355	—	*	Polytec PI Inc. Auburn, MA 01501 Tel: (508) 832-3456	*	*
Hewlett-Packard Inc. Test & Measurement Customer Business Center, P.O. Box 4026 Englewood, CO 80155-4026 Tel: (800) 829-4444	*	—	Spindler & Hoyer Inc. 459 Fortune Blvd. Milford, MA 01757-1745 Tel: (508) 478-6200	*	*
Melles Griot Inc. 4665 Nautilus Court South Boulder, CO 80301 Tel: (303) 581-0337	*	*	Zygo Corporation Middlefield, CT 06455-0448 Tel: (860) 347-8506	*	*

6.6 Bore Gaging Displacement Sensors

Viktor P. Astakhov

Dimensions are a part of the total specification assigned to parts designed by engineering. However, the engineer in industry is constantly faced with the fact that no two objects in the material world can ever be made exactly the same. The small variations that occur in repetitive production must be considered in the design. To inform the workman how much variation from exact size is permissible, the designer uses a tolerance or limit dimension technique. A *tolerance* is defined as the total permissible variation of size, or the difference between the limits of size. *Limit dimensions* are the maximum and minimum permissible dimensions. Proper tolerancing practice ensures that the finished product functions in its intended manner and operates for its expected life.

Bore Tolerancing

All bore dimensions applied to the drawing, except those specifically labeled as basic, gage, reference, maximum, or minimum, will have an exact tolerance, either applied directly to the dimension or indicated by means of general tolerance notes. For any directly tolerated decimal dimension, the tolerance has the same number of decimal places as the decimal portion of the dimension.

Engineering tolerances may broadly be divided into three groups: (1) *size tolerances* assigned to dimensions such as length, diameter, and angle; (2) *geometric tolerances* used to control a hole shape in the longitudinal and transverse directions; and (3) *positional tolerances* used to control the relative position of mating features. Interested readers may refer to [1, 2].

The ISO system of limits and fits (ISO Recommendation R 286) covers standard tolerances and deviations for sizes up to 3150 mm. The system is based on a series of tolerances graded to suit all classes of work from the finest to the most coarse, along with different types of fits that range from coarse clearance to heavy interference. Here, *fit* is the general term used to signify the range of tightness that may result from the application of a specific combination of tolerances in the design of mating parts.

There are 18 tolerance grades intended to meet the requirements of different classes of parts. These tolerance grades are referred to as ITs and range from IT 01, IT 02 (reserved for the future), and IT 1, to IT 16 (for today's use). In each grade, the tolerance values increase with size according to a formula that relates the value of a given constant to the mean diameter of a particular size range. The system provides 27 different fundamental deviations for sizes up to and including 500 mm, and 14 for larger sizes to give different type of fits ranging from coarse clearance to heavy interference. Interested readers may refer to [3].

Bore Gage Classification and Specification

To measure the above-listed tolerances, modern manufacturing requires the use of gages. A *gage* is defined as a device for investigating the dimensional fitness of a part for specific function. *Gaging* is defined by ANSI as a process of measuring manufactured materials to assure the specified uniformity of size and contour required by industries. Gaging thereby assures the proper functioning and interchangeability of parts; that is, one part will fit in the same place as any similar part and perform the same function, whether the part is for the original assembly or replacement in service.

Bore gages may be classified as follows:

1. Master gages
2. Inspection gages
3. Manufacturer's gages
4. Gages that control dimensions
5. Gages that control various parameters of bore geometry
6. Fixed limit working gages

7. Variable indicating gages
8. Post-process gages
9. In-process gages

Master gages are made to their basic dimensions as accurately as possible and are used for reference, such as for checking or setting inspection of manufacturer's gages. *Inspection gages* are used by inspectors to check the manufactured products. *Manufacturer's gages* are used for inspection of parts during production.

Post-process gages are used for inspecting parts after being manufactured. Basically, this kind of gage accomplishes two things: (1) it controls the dimensions of a product within the prescribed limitations, and (2) it segregates or rejects products that are outside these limits. Post-process gaging with feedback is a technique to improve part accuracy by using the results of part inspection to compensate for repeatable errors in the machine tool path. The process is normally applied to CNC (computer numerically controlled) machines using inspection data to modify the part program, and on tracer machines using the same data to modify the part template.

In-process gages are used for inspecting parts during the machining cycle. In today's manufacturing strategy, in-process gages and data-collection software provide faster feedback on quality. Indeed, the data-collection and distribution aspect of 100% inspection has become as important as the gaging technology itself. Software specifically designed to capture information from multiple gages, measure dozens of products types and sizes, and make it available to both roving inspectors and supervising quality personnel as needed, is quickly becoming part of quality control strategies. In conjunction with computer numerically controlled (CNC) units, in-process gaging can automatically compensate for workpiece misalignment, tool length variations, and errors due to tool wear.

Gages That Control Dimensions

Gages that control dimensions are used to control bore diameter. These gages can be either post-process or in-process gages. Further, these gages can be either *fixed limit gages* or *variable indicating gages*.

A *plug gage* is a fixed limit working bore gage. These inexpensive gages do not actually measure dimensions or geometry. They simply tell the operator whether the bore is oversized or undersized. The actual design of most plug gages is standard, being covered by American Gage Design (AGD) standards. However, there are many cases where a special plug gage must be designed.

A plug gage is usually made up of two members. One member is called the go end, and the other the no-go or not-go end. The gage commonly has two parts: the gaging member, and a handle with the sign, go or no-go, and the gagemaker's tolerance marked on it. There are generally three types of AGD standard plug gages. First is the single-end plug gage (Figure 6.53(a)); the second is the double-end (Figure 6.53(b)); and the third is the progressive gage (Figure 6.53(c)). Interested readers may refer to [4].

Fixed-limit gage tolerance is generally determined from the amount of workpiece tolerance. A 10% rule is generally used for determining the amount of gage tolerance for fixed, limit-type gages. Four classes of gagemakers' tolerances have been established by the American Gage Design Committee and are in general use [4]. These four classes establish maximum variation for any designed gage size. The degree of accuracy needed determines the class of gage to be used. Table 6.11 shows these four classes of gagemakers' tolerances. Referring to Table 6.11, class XX gages are used primarily as master gages and for final close tolerance inspection. Class X gages are used for some types of master gage work and as close tolerance inspection and working gages. Class Y gages are used as inspection and working gages. Class Z are used as working gages where part tolerances are large. Table 6.12 shows the diameter ranges and prices of the plug gages manufactured by the Flexbar Machine Corp.

Variable indicating gages allow the user to inspect some bore parameters and get numbers for charting and statistical process control (commonly abbreviated as SPC). These gages have one primary advantage over fixed gages: they show how much a hole is oversized or undersized. When using a variable indicating gage, a master ring gage to the nominal dimension to be checked must be used to preset the gage to zero. Then, in applying the gage, the variation from zero is read from the dial scale. Figure 6.54 shows industry's

TABLE 6.12 Premium Quality Hardened Steel GO/NO GO Plug Gages by the Flexbar Machine Corp.

Size range	Class	Price (\$)			Handle	
		1	2-4	5-10	No	Price (\$)
0.01 in. to 0.030 in. 0.25 mm to 0.762 mm	XX	35.45	28.65	22.26	1W	8.00
	X	31.15	23.35	17.00		
	Y	28.35	22.15	15.05		
	Z	25.65	20.70	12.55		
0.03 in. to 0.075 in. 0.762 mm to 1.91 mm	XX	19.25	15.30	13.63	1W	8.00
	X	15.05	11.70	10.60		
	Y	13.90	10.80	9.40		
	Z	11.40	9.35	8.15		
0.075 in. to 1.80 in. 1.91 mm to 4.57 mm	XX	21.15	17.00	14.80	2W	8.30
	X	18.15	14.45	12.85		
	Y	17.00	13.90	11.25		
	Z	14.00	11.10	9.65		
0.180 in. to 0.281 in. 4.57 mm to 7.14 mm	XX	22.00	17.80	15.30	3W	9.00
	X	18.90	15.30	13.10		
	Y	17.85	13.55	11.95		
	Z	14.30	11.40	9.90		
0.281 in. to 0.406 in. 7.17 mm to 10.31 mm	XX	24.20	19.50	17.80	4W	9.35
	X	21.15	17.00	14.80		
	Y	19.30	15.65	13.80		
	Z	14.90	12.00	10.50		
0.406 in. to 0.510 in. 10.31 mm to 12.95 mm	XX	25.35	20.35	17.80	5W	9.70
	X	22.25	17.80	15.80		
	Y	20.45	16.50	14.55		
	Z	16.30	13.15	11.30		
0.510 in. to 0.635 in. 12.95 mm to 16.13 mm	XX	26.70	21.45	18.50	6W	11.40
	X	23.70	19.25	16.30		
	Y	21.85	17.60	15.50		
	Z	17.50	14.30	12.55		
0.635 in. to 0.760 in. 16.13 mm to 19.30 mm	XX	28.15	22.60	19.85	7W	13.50
	X	25.35	20.05	17.50		
	Y	23.25	18.45	16.40		
	Z	18.65	15.15	13.45		
0.780 in. to 1.010 in. 19.30 mm to 25.65 mm	XX	44.25	34.25	32.80	8W	20.00
	X	39.56	29.80	27.30		
	Y	36.30	27.25	24.90		
	Z	32.40	25.65	24.20		

Gages with Manual Probe Head Systems

Geometry gages with manual probe head systems are rapidly becoming common in many high-precision metalworking applications. The simplest form of manual probe head systems in common use is the air plug gage (spindle) (Figure 6.57). Compressed air from the air gage indicating unit is pressed to the plug gage and allowed to escape from two or more jets in the periphery. When the air plug gage is inserted into a hole, the air escaping from the jets is limited by the clearance between the jet faces and the hole. The small changes in clearance, arising when the air plug gage is inserted in successive holes of different sizes, produce changes in the flow rate or back pressure in the circuit. The magnification and datum setting of systems with variable control orifices and zero bleeds is carried out with master holes. Some errors of form that can be detected with air plug gages are: (1) taper, (2) bell mouthing, (3) barreling, (4) ovality, and (5) lobing. Dearborn (Dearborne Gage Company, MI) open-orifice air spindles are available as standard to use in measuring thru, blind, and counterbored holes ranging in diameter from 0.070 in. (2 mm) to 6.000 in. (154 mm).



FIGURE 6.54 Industry's most popular dial and electronic bore gages (Courtesy of The L.S. Starrett Co.).

Another type of geometry gage with manual probe head system is the electronic bore gage. These gages measure bores at various depths to determine conditions such as bellmouth, taper, convexity, or concavity. They are also able to determine out-of-roundness conditions when equipped with a 3-point measuring system. Figure 6.58 shows a TRIOMATIC® electronic bore gage (Brown & Sharpe), and



FIGURE 6.54 (continued)

runout, vertical straightness, and vertical parallelism. The machines are supplied with MeasurLink® data acquisition software for Windows™, which allows immediate measurement data analysis and feedback for variable, attribute, and short inspection runs. The software gives the quality control/production manager the ability to create traceability lists of unlimited size. Information such as machine center, operator, materials used, assignable causes, and other relevant data can be stored and attached to measurement values. See Table 6.16 for a list of companies that make bore gages.

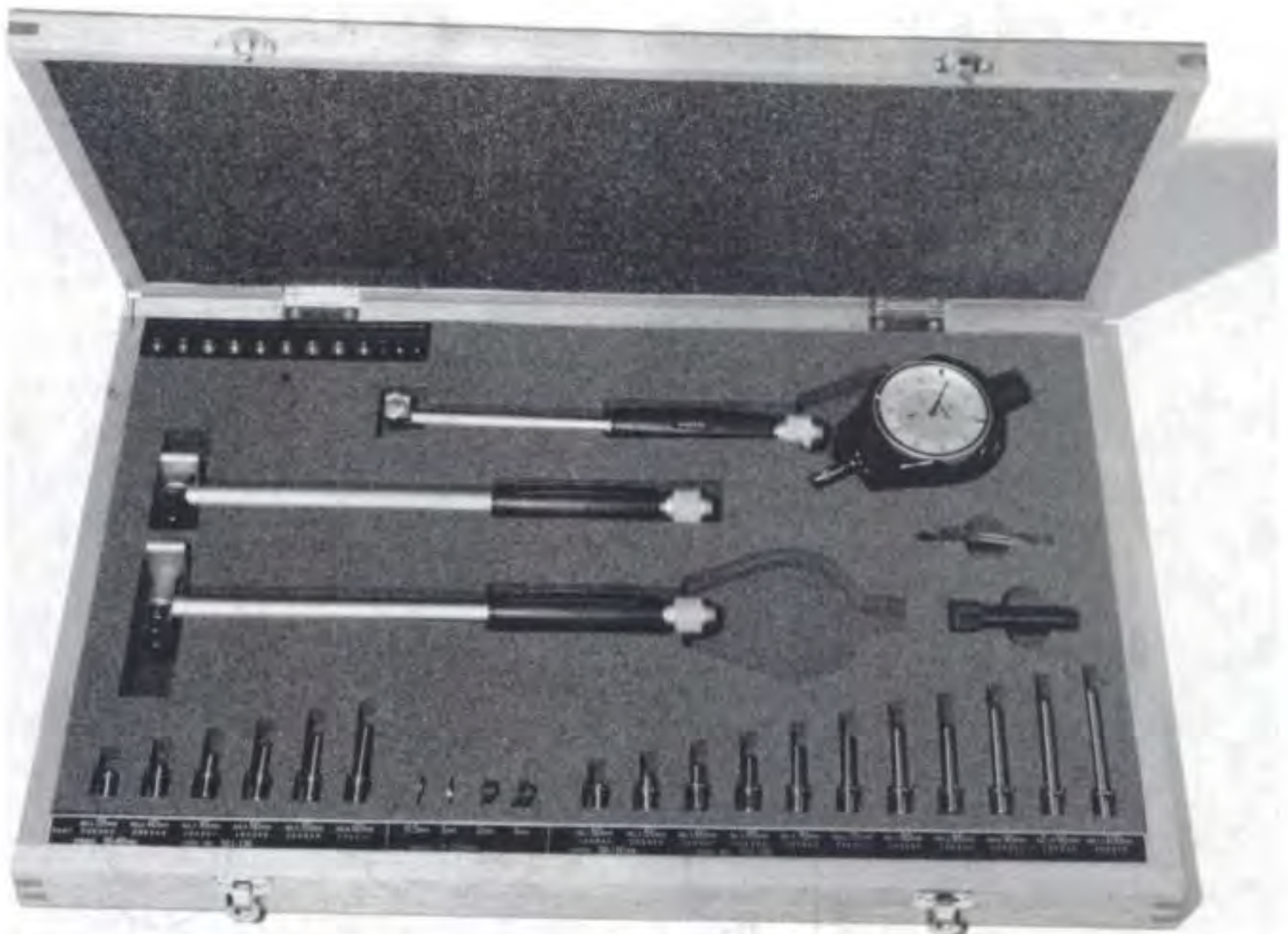


FIGURE 6.55 Set of dial bore gages (Courtesy of MITUTOYO/MTI Corporation).



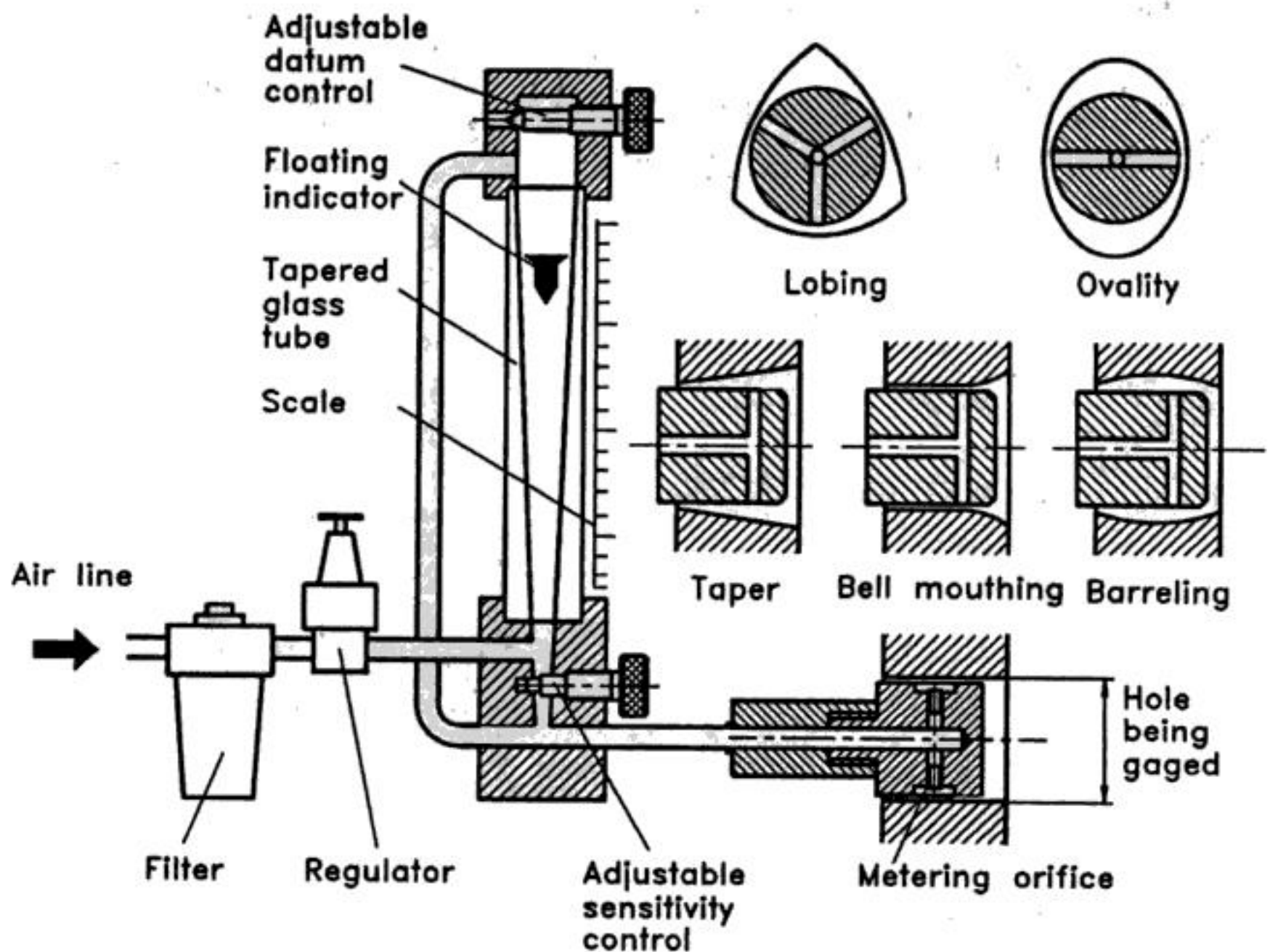
FIGURE 6.56 Intrimic® plus internal micrometer (Courtesy of Brown & Sharpe Manufacturing Company).

GAGE R AND R Standards

Gage Repeatability and Reproducibility (GAGE R AND R) capability standards have direct implications for parts makers and for gage manufacturers. Repeatability is the ability of an operator using a single gage to obtain the same measurements during a series of tests. Reproducibility is the ability of different

TABLE 6.13 Intrimik® Plus Internal Micrometers by Brown & Sharpe

Range	B&S Tool No.	Price (\$)
0.275 in. to 0.350 in. (6–8 mm)	599-290-35	745.40
0.350 in. to 0.425 in. (8–10 mm)	599-290-42	745.40
0.425 in. to 0.500 in. (10–12 mm)	599-290-50	745.40
0.500 in. to 0.600 in. (12–14 mm)	599-290-60	826.70
0.600 in. to 0.700 in. (14–17 mm)	599-290-70	826.70
0.700 in. to 0.800 in. (17–20 mm)	599-290-80	826.70
0.800 in. to 1.0 in. (20–25 mm)	599-290-100	843.60
1.0 in. to 1.2 in. (25–30 mm)	599-290-120	843.60
1.2 in. to 1.4 in. (30–35 mm)	599-290-140	854.30
1.4 in. to 1.6 in. (35–40 mm)	599-290-160	854.30
1.6 in. to 2.0 in. (40–50 mm)	599-290-200	933.10
2.0 in. to 2.4 in. (50–60 mm)	599-290-240	933.10
2.4 in. to 2.8 in. (60–70 mm)	599-290-280	933.10
2.8 in. to 3.2 in. (70–80 mm)	599-290-320	950.20
3.2 in. to 3.6 in. (80–90 mm)	599-290-360	950.20
3.6 in. to 4.0 in. (90–100 mm)	599-290-400	950.20
Intrimik Plus Complete Set #5	599-290-5	3374.60

**FIGURE 6.57** Air plug gage.

operators to obtain similar results with the same gage. GAGE R AND R blends these two factors together to determine a measuring system's reliability and its suitability for a particular measuring application. For example, a gage design that meets the GAGE R AND R standards for 50 mm (2 in.) bores may be unsatisfactory on 250 mm (10 in.) bores. A gage that meets a tolerance of 2 μm , may not be satisfactory

TABLE 6.15 Rountest Machines by Mitutoyo/MIT Corp.

Model	RA-112	RA-334	RA-434	RA-661
Measuring range	11 in. (280 mm)	11.8 in. (300 mm)	11.8 in. (300 mm)	
Max. measuring dia.	8.6 in. (220 mm)	27.6 in. (700 mm)	13.8 in. (350 mm)	
Max. measuring height	—	21.7 in. (550 mm)	20.4 in. (520 mm)	
Max. loading dia.	22 lb (10 kg)	66.1 lb (30 kg)	132 lb (60 kg)	
Max. loading capacity	± 0.01 in. ($\pm 250 \mu\text{m}$)	± 0.012 in. ($\pm 300 \mu\text{m}$)	± 0.012 in. ($\pm 300 \mu\text{m}$)	
Detector	Range	7-10 gf	7-10 gf	
Turntable	Measuring force	(1.6 + 0.6H) μinch	(1.6 + 0.6H) μinch	
	Rotating accuracy	(0.04 + 0.3H) μm	(0.04 + 0.6H) μm	
	Centering adj. range	± 0.08 in. (± 2 mm)	± 0.2 in. (± 5 mm)	
	Leveling adj. range	$\pm 1^\circ$	$\pm 1^\circ$	
Z-axis column	Rotating speed	6 rpm	2, 4, 6 rpm	8 $\mu\text{inch}/8$ in.
	Straightness	—	40 $\mu\text{inch}/7.9$ in.	0.2 $\mu\text{m}/200$ mm
	Parallelism	—	1 $\mu\text{m}/200$ mm	
			120 $\mu\text{inch}/7.9$ in.	
			80 $\mu\text{inch}/13.8$ in.	
Measuring magnifications	Stroke	10 in. (25 mm)	3 $\mu\text{m}/200$ mm	2 $\mu\text{m}/200$ mm
Dimensions	Measuring unit	100—20,000 \times	18.9 in. (480 mm)	13.8 in.
$W \times D \times H$	Electric unit	9.9 \times 15.9 \times 21.5 in.	100—50,000 \times	100—100,000 \times
	Measuring unit	251 \times 404 \times 576 mm	24.4 \times 19.7 \times 36.2 in.	28.7 \times 23.2 \times 62.2 in.
Mass	Measuring unit	11.4 \times 11.8 \times 3.6 in.	620 \times 500 \times 920 mm	730 \times 590 \times 1580 mm
	Measuring unit	290 \times 300 \times 92 mm	9.8 \times 16.1 \times 13 in.	30.7 \times 23.3 \times 28.8 in.
	Electric (analyzer) unit	61.7 lb (28 kg)	250 \times 410 \times 330 mm	780 \times 592 \times 732 mm
Base price (\$)	Base price (\$)	276 lb (125 kg)	24.2 lb (11 kg)	298 lb (135 kg)
		45,000	60,000	770 lb (350 kg)
				110 lb (50 kg)
				60,000

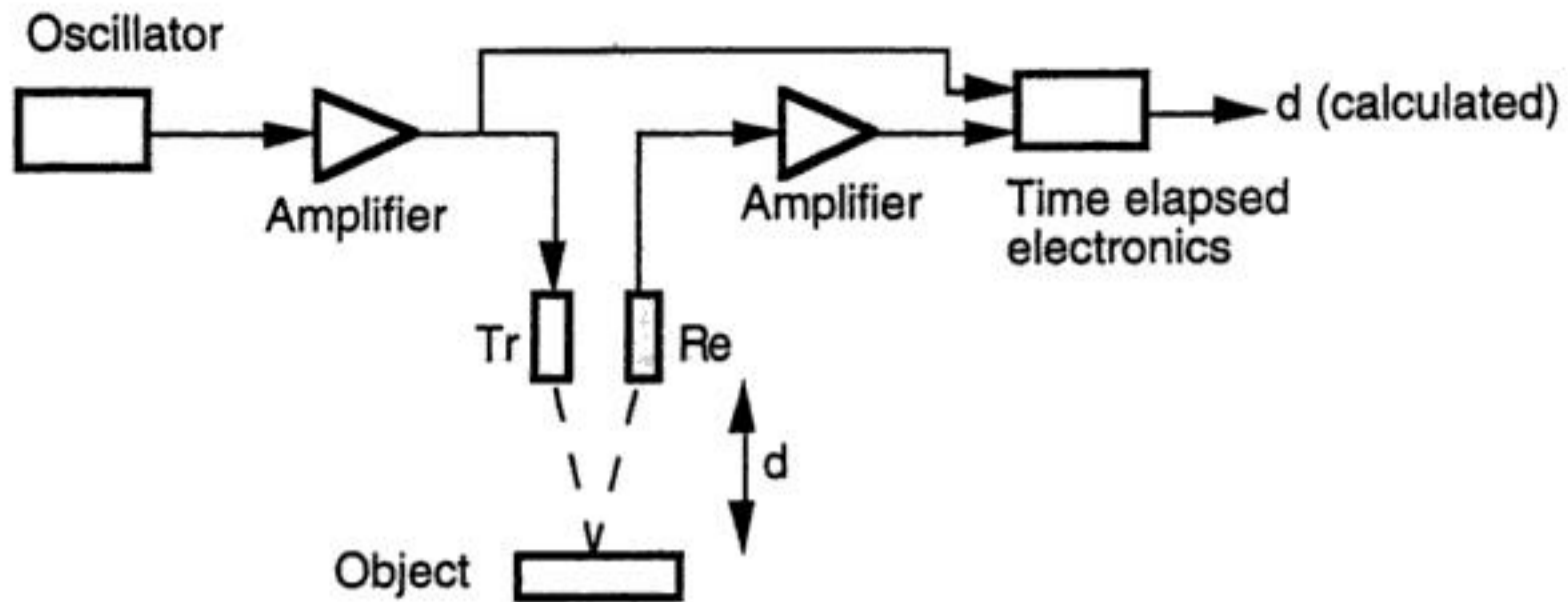


FIGURE 6.61 Principle of a pulse-echo ultrasound system for distance measurements (Tr = transmitter, Re = receiver).

$$c = \sqrt{\frac{1}{K\rho}} \quad (6.85)$$

In gases, the velocity of sound is described by Equation 6.86. Here g represents the ratio of the specific heat at constant pressure (c_p) to the specific heat at constant volume (c_v), p is pressure, R is the universal gas constant, T is the absolute temperature, and M is the molecular weight.

$$c = \sqrt{\frac{gRT}{M}} = \sqrt{\frac{c_p}{c_v} \frac{p}{\rho}} \quad (6.86)$$

An important quantity is the *specific acoustic impedance*. It is, in general, a complex quantity but in the far field (Figure 6.63), the imaginary component disappears, leaving a real quantity. This real quantity is the product of the density ρ and the sound speed c in the medium. This product is called the characteristic impedance R_s (Equation 6.87).

$$R_s = \rho c \quad (6.87)$$

The characteristic impedance is thus independent of the sound frequency.

An acoustic wave has an intensity I (rate of flow of energy per unit area), which can be expressed in watts per square meter (W m^{-2}). A usually unwanted phenomenon arises when the sound wave has to pass from one medium with characteristic impedance R_1 to another medium with characteristic impedance R_2 . If R_1 and R_2 have different values, a part of the wave intensity will reflect at the boundary between the two media (see Figure 6.61 and 6.62). The two media are said to be mismatched, or poorly coupled, if a major part of the wave intensity is reflected and a minor part is transmitted. The relative amounts of reflected and transmitted wave intensities can be defined by:

$$\text{Reflection coefficient} = \frac{I_{\text{refl}}}{I_{\text{incident}}} \quad (6.88a)$$

$$\text{Transmission coefficient} = \frac{I_{\text{trans}}}{I_{\text{incident}}} \quad (6.88b)$$

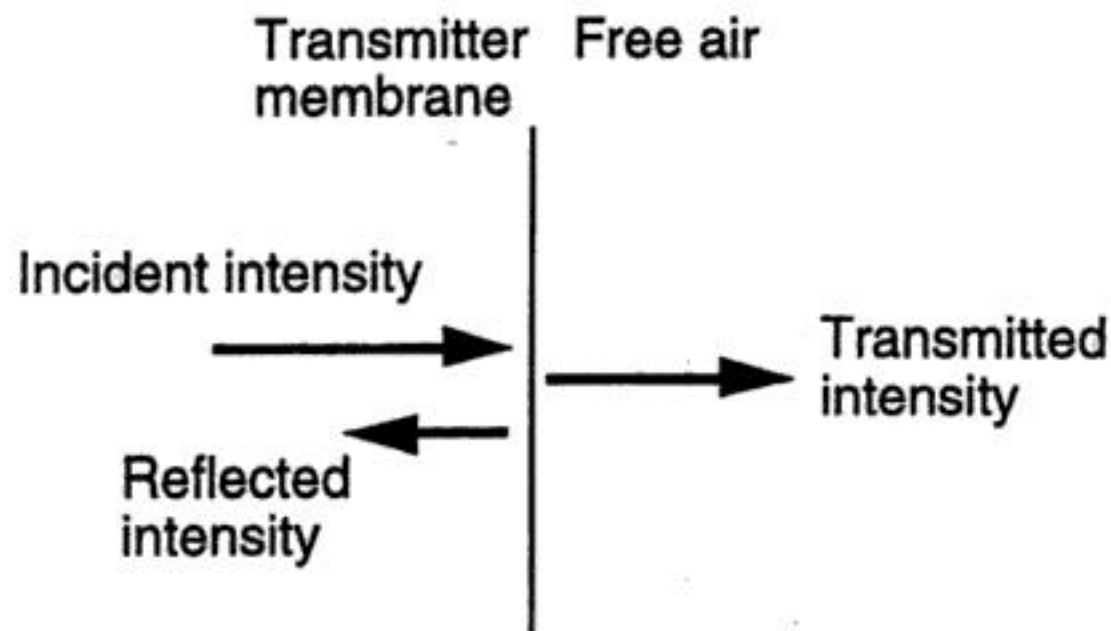


FIGURE 6.62 Reflection and transmission of a sound wave at the interface between media of different characteristic impedances.

It can be shown [1] that these coefficients have simple relations to the previously mentioned characteristic impedances.

$$\text{Reflection coefficient} = \frac{(R_1 - R_2)^2}{(R_1 + R_2)^2} \quad (6.89a)$$

$$\text{Transmission coefficient} = \frac{4R_1R_2}{(R_1 + R_2)^2} \quad (6.89b)$$

The practical importance of the acoustic impedance is realized when the ultrasonic pulse-echo system shown in Figure 6.61 is considered. First, the electric energy is converted into mechanical vibrations of a membrane in the transmitter. Second, the vibrations (the sound wave) have to pass through the boundary between the membrane (usually a solid material) and free air. Because the transmitter membrane and the free air have different characteristic impedances, much of the acoustic intensity is reflected (Figure 6.62).

The transmitted ultrasound in free air will first propagate in a parallel beam (near field of the transducer); but after a distance L , the beam diverges (the far field of the transducer). See Equation 6.90 and Figure 6.63.

$$L = \frac{D^2}{4\lambda} \quad (6.90)$$

D is the diameter of the circular transmitter and λ is the wavelength of the ultrasound.

The sound intensity in the near field is complicated due to interference effects of sound originating from different parts of the transducer membrane. In the far field, the intensity is approximately uniform and the beam spread follows:

$$\sin\beta = 1.22 \frac{\lambda}{D} \quad (6.91)$$

where β is the half lobe angle.

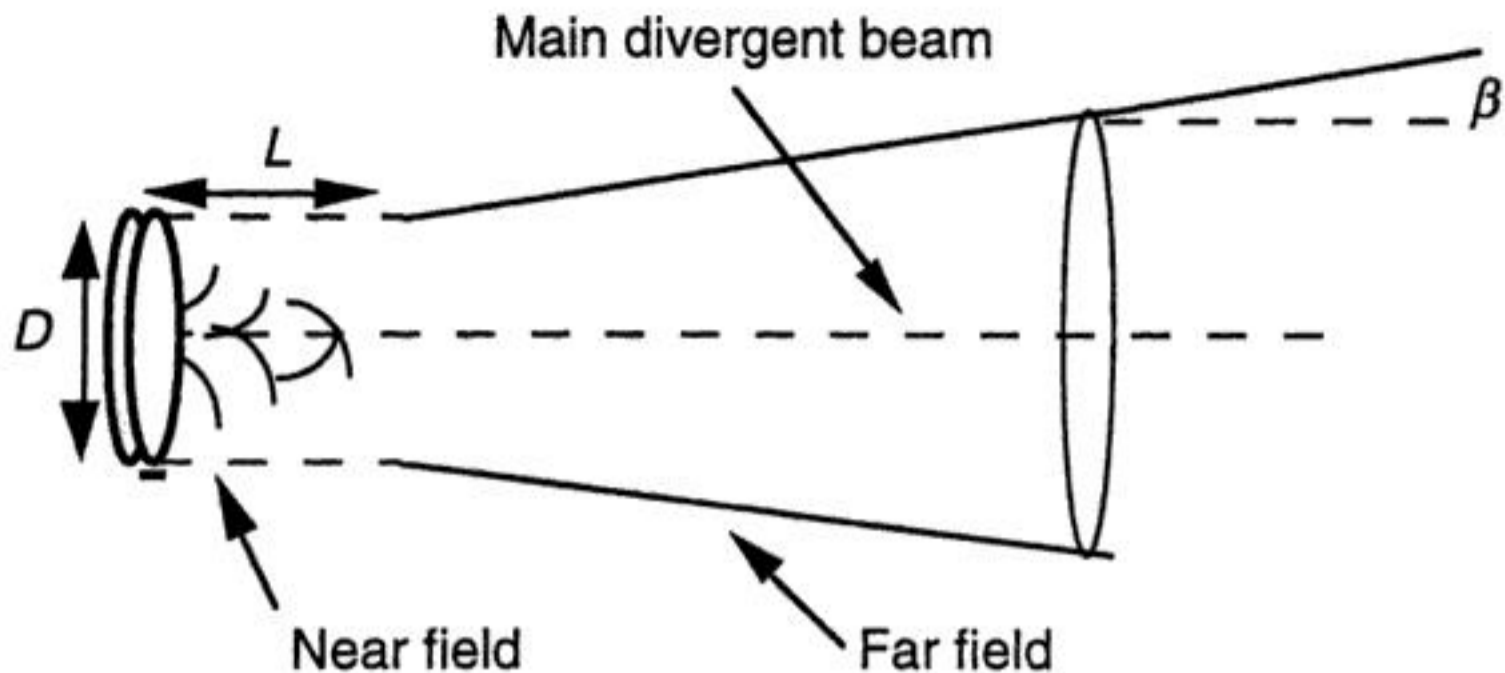


FIGURE 6.63 Illustration of the ultrasound beam in the near field and the far field of the transducer.

To get a narrow beam, the transmitter membrane diameter must be large with respect to the wavelength. High-frequency ultrasound cannot be the general solution as ultrasound of a high frequency is absorbed faster than ultrasound of a low frequency [1-4].

Ultrasound Transducers

Most ultrasound transducers convert electric energy to mechanical energy and vice versa. The most common types of in-air transducers are [5-7]:

1. Mechanical
2. Electromagnetic
3. Piezoelectric
4. Electrostatic
5. Magnetostrictive

The simplest type, mechanical transducers such as whistles and sirens, are used up to approximately 50 kHz. This type works only as a transmitter.

Electromagnetic transducers such as loudspeakers and microphones can be used for ultrasonic wave generation, but they are mainly suited for lower frequencies.

The piezoelectric transducer (Figure 6.64) is more suitable for use in ultrasonics and is quite common. It uses a property of piezoelectric crystals: they change dimensions when they are exposed to an electric field. When an alternating voltage is applied over the piezoelectric material, it changes its dimensions with the frequency of the voltage. The transducer is mainly suited for use at frequencies near the mechanical resonance frequency of the crystal. The piezoelectric transducer can be both a transmitter and a receiver: when a piezoelectric material is forced to vibrate by a sound pulse, it generates a voltage. Some natural crystals, such as quartz, are piezoelectric. Ceramics can be polarized to become piezoelectric; so can some polymers like PVDF (polyvinylidene fluoride). Polymers are suitable as transducers in air since their acoustic impedance is low [8-10] compared with other standard piezoelectric materials.

The electrostatic transducer (Figure 6.64) is a plate capacitor with one plate fixed and the other free to vibrate as a membrane. When a voltage is applied between the plates, the electrostatic forces tend to attract or repel the plates relative to each other depending on the polarity of the voltage. This transducer can be used both as a transmitter and a receiver [11].

The magnetostrictive transducer is based on the phenomenon of magnetostriction, which means that the dimensions of a ferromagnetic rod change due to the changes of an externally applied magnetic field. This transducer can also act as both a receiver and a transmitter.



FIGURE 6.64 Ultrasonic transducers — piezoelectric (left) and electrostatic (right).

Principles of Time-of-Flight Systems

There are several techniques for ultrasonic range measurements [12–15].

The previously described *pulse echo method* is the simplest one. Usually, this method has a low signal-to-noise ratio (SNR) because of the low transmitted energy due to the short duration of the pulse. Multireflections are detectable.

In the *phase angle method*, the phase angle is measured between the continuous transmitted signal and the continuous received signal and is used as a measure of the distance. The method is relatively insensitive to disturbances. Multireflections are not detectable in a meaningful way. When the distance is longer than one wavelength, another method must be used to monitor the distance.

The *frequency modulation method* uses transmitted signals that are linearly frequency modulated. Thus, detected signals are a delayed replica of the transmitted signal at an earlier frequency. The frequency shift is proportional to the time-of-flight. The method is robust against disturbing signals, and multireflections are detectable.

The *correlation method* (Figure 6.65) determines the cross-correlation function between transmitted and received signals. When the transmitted signal is a random sequence, i.e., white Gaussian noise, the cross-correlation function estimates the impulse response of the system, which, in turn, is a good indicator of all possible time delays. The method is robust against disturbances, and multireflections are detectable.

Industrial acoustic noise can affect the received signals in an ultrasound time-of-flight system. The noise can be generated from leaking compressed air pipes, noisy machines, or other ultrasonic systems. This external noise is not correlated with the relevant echo signals of the transmitted noise and can therefore be eliminated by the use of correlation methods. Disturbances correlated with the relevant echo signal (e.g., unwanted reflections) will not be eliminated by the use of correlation methods.

The impulse response $h(t, t_0)$ is used as a sensitive indicator of time delay between transmitted signal at time t_0 and received signal at time t . The impulse response is given by [14]:

$$h(t, t_0) = F^{-1} \begin{bmatrix} S_{xy} \\ S_{xx} \end{bmatrix} \quad (6.92)$$

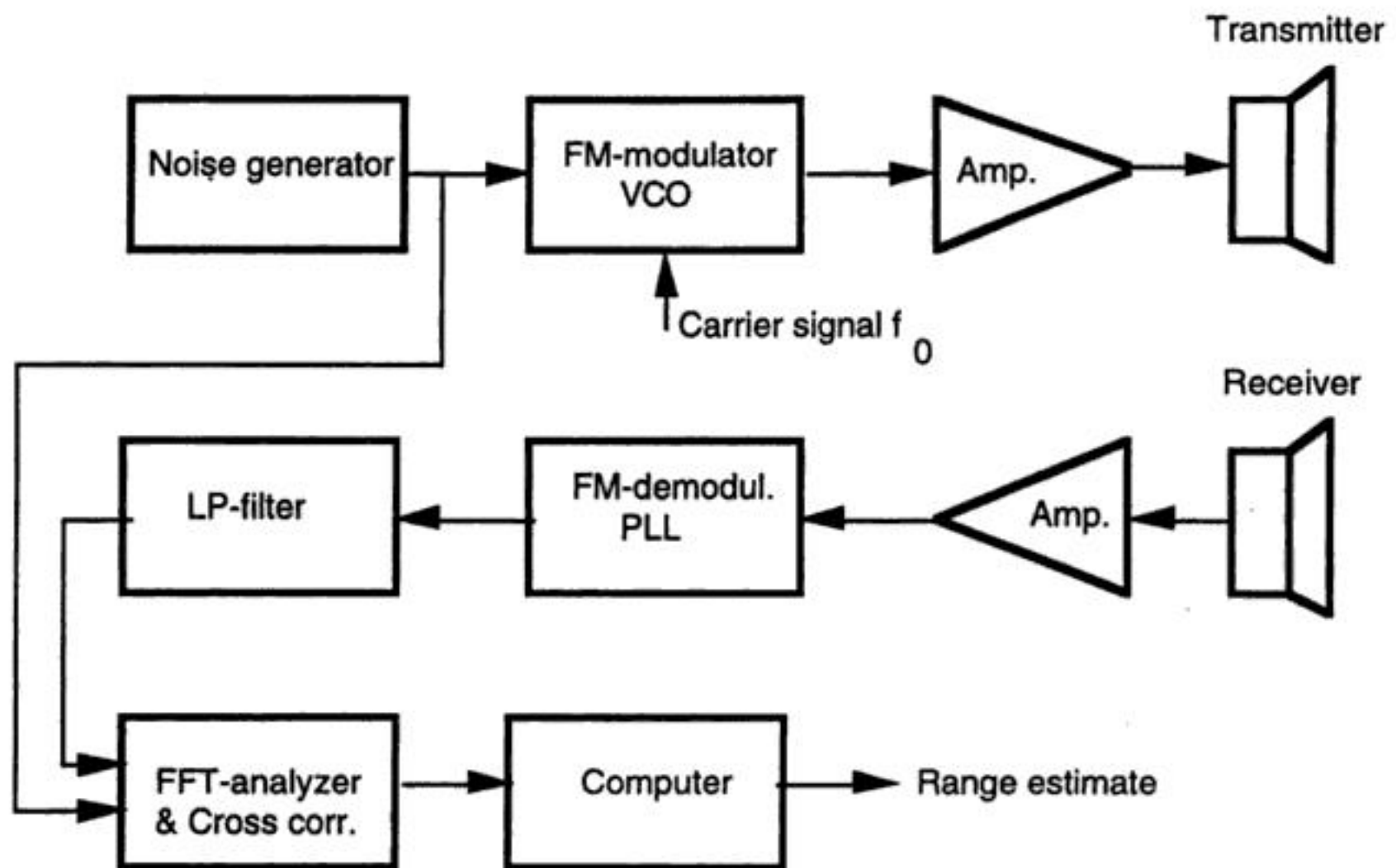


FIGURE 6.65 Diagram of a correlation-based time-of-flight system.

where F^{-1} is the inverse Fourier transform, x is the transmitted signal, y is the received signal, $S_{xy}(f)$ is the cross-spectral density function [the Fourier transform of the cross-correlation function of the transmitted signal $x(t)$ and the received signal $y(t)$] and S_{xx} is the power density function (the Fourier transform of the auto-correlation function of the transmitted signal).

To analyze the transfer channel by data acquisition requires a high sampling rate, theoretically at least two times the highest frequency component in the received signal (and in practice as high as 10 times the highest frequency). One way to reduce the sampling rate is to first convert the signal from its bandpass characteristics around a center frequency f_0 (approx. 50 kHz) to lowpass characteristics from dc to $B/2$, where B is the appropriate bandwidth. The accuracy of the range estimate, and hence the time interval $t - t_0$, can be improved by processing the estimate of the impulse response $h(t - t_0)$ with a curve-fitting (least square) method and digital filtering in a computer. A block diagram of a correlation-based time-of-flight system is shown in Figure 6.65. Further details and complete design examples can be found in the literature [12–15].

Table 6.17 lists some advantages and drawbacks of the described time-of-flight methods.

TABLE 6.17 Advantages and Disadvantages of Time-of-Flight Methods

Method	Main advantage	Main disadvantage
Pulse echo method	Simple	Low signal-to-noise ratio
Phase angle method	Rather insensitive to disturbances	Cannot be used directly at distances longer than the wavelength of the ultrasound
Frequency modulation method	Robust against disturbances; multireflections detectable	Can give ambiguous results measurements on long and short distances can give the same result (compare with phase angle method)
Correlation method	Very robust against disturbances	Make relatively high demands on hardware and/or computations

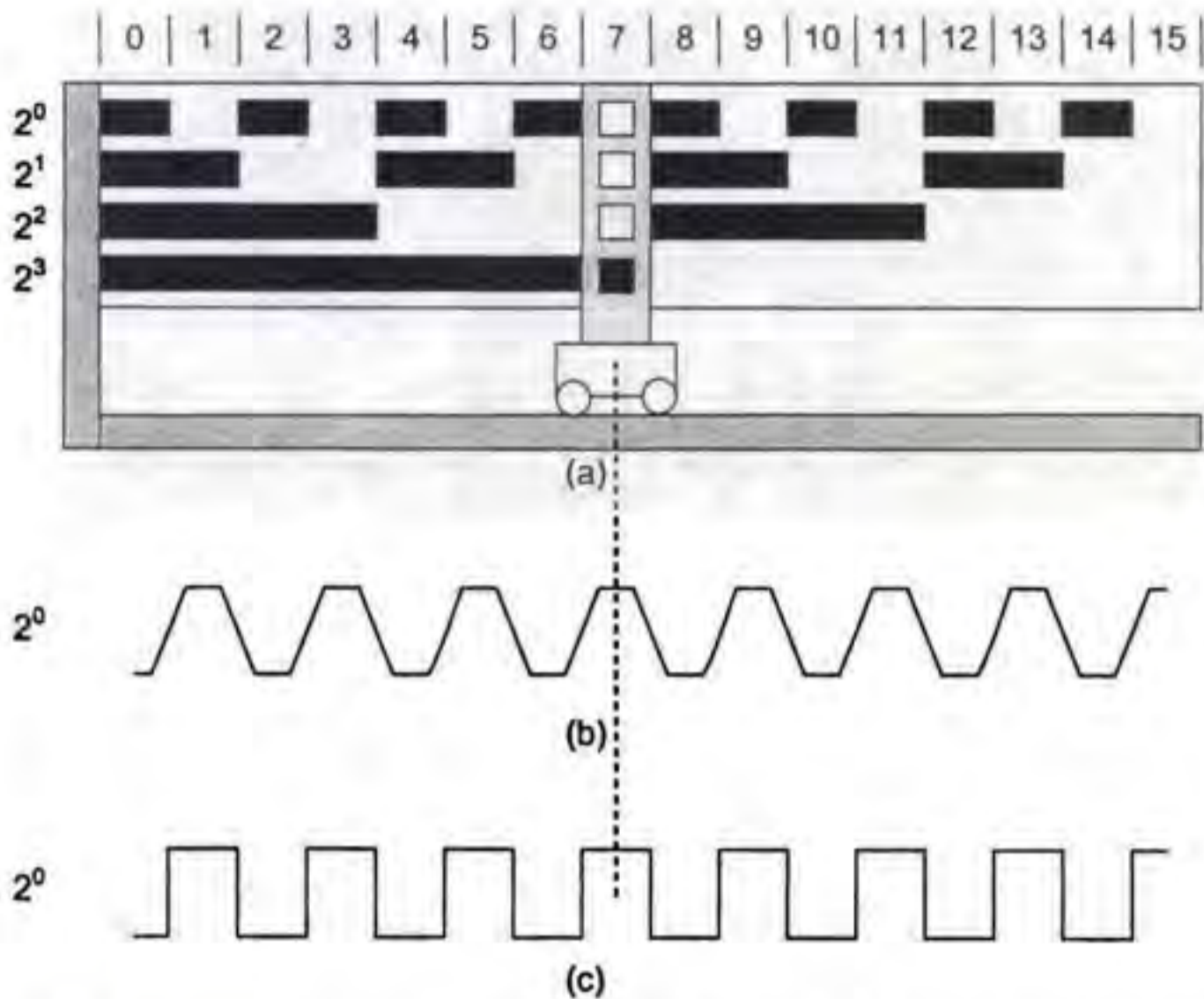


FIGURE 6.67 (a) Absolute encoders using a natural binary code of four digits. Four tracks are required. The moving read head has four apertures and is shown in position 7 along the scale. (b) The output of the read head aperture corresponding to the least significant track. It represents the proportion of light area covering the aperture. (c) The binary digit obtained after squaring the raw output signal.

Absolute encoders are classified according to the type of code used. The main four codes are Gray, binary (usually read by vee-scan detection), optical resolving, and pseudorandom. All absolute encoders use geometric masking to generate the code.

Encoder Signals and Processing Circuitry

Absolute Encoders

Direct Binary

Figure 6.67(a) illustrates the concept of an absolute linear optical encoder using a direct binary encoded scale. The fixed scale has n tracks (here n is 4), each providing one bit of a direct binary number. The lowest track (first track from the center of the disk for a rotary encoder) is the most significant digit and has a weight, 2^{n-1} (here 2^3), while the upper track is the least significant digit with a weight 2^0 . The track providing the least significant digit has 2^{n-1} cycles of light and dark records, while the most significant track has 2^0 or 1 such cycle. For each track, the moving read head has a readout unit consisting of a light source, a mask, and a photodetector. Figure 6.67(b) shows the output from the photodetector, which represents the total intensity of light reaching its surface. As the mask passes over a clear region of the grating, the photodetector output increases, and then decreases. In theory, a truncated triangular wave is obtained, which can easily be converted to a square wave (Figure 6.67(c)) by a suitably chosen thresholding level. The result is a high or 1 for a light record and a low or 0 for a dark one. The position, in base 10, corresponding to the reading head position in Figure 6.67 is

$$1 \times 2^0 + 1 \times 2^1 + 1 \times 2^2 + 0 \times 2^3 = 7 \quad (6.93)$$

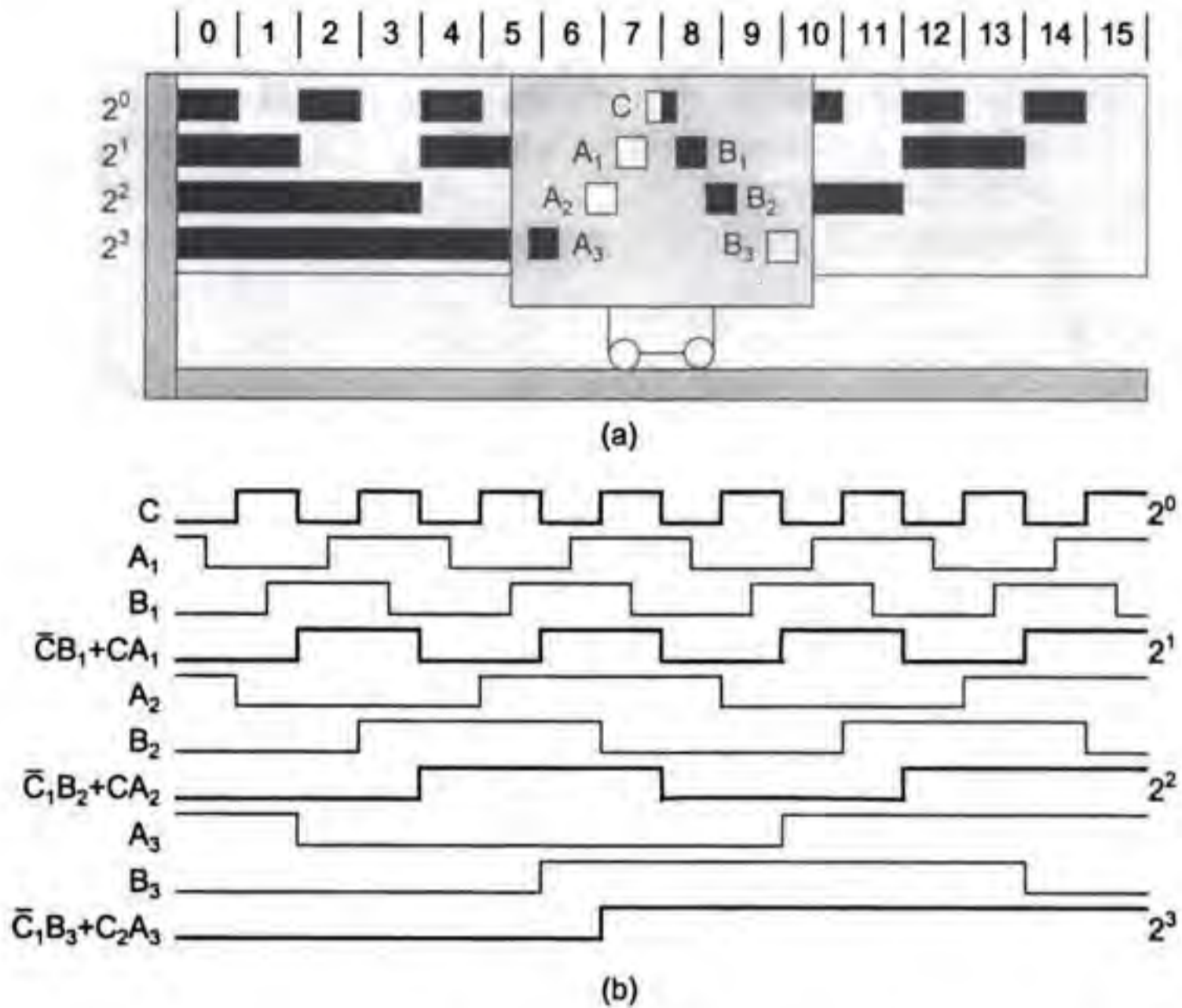


FIGURE 6.68 The vee-scan configuration of reading units in (a) removes the ambiguity associated with a natural binary scale. Simple combinational logic is then used to generate the natural binary readout in (b).

The code configuration just described is not suitable for practical use because some transitions require that two or more bit values change simultaneously. For example, from position 7 to position 8, all bits change values. Unless the change is simultaneous, an incorrect position readout results at some position. This would require that the scale is geometrically perfect, that the read head be perfectly aligned with the scale, and that the electronics are perfectly adjusted and stable over time. This problem is solved either by the use of a vee-scan detection method or the use of a unit-distance code such as the Gray code.

Vee-scan

The vee-scan method uses a V-shape pattern of readout units that removes the potential reading ambiguity of direct binary scales. Stephens et al. [1] indicate the read points at which transitions are detected in Figure 6.68(a). They also describe the conversion of the thresholded output signals to binary code using combinational logic. The primary advantage is that the location tolerance of the transition point of each reading unit need only be $\pm 1/8$ of the cycle length for that particular track. For example, $\pm 45^\circ$ for the most significant track of a rotary encoder disk. Figure 6.68(b) shows a direct binary word obtained through logic combinations of the vee-scan readings.

Gray Code

The use of vee-scan requires additional reading heads as well as processing electronics. The Gray code is a unit-distance code and so only one bit of data changes between representations of two consecutive numbers or successive positions. This removes the possibility of ambiguous readout. It has the following advantages: (1) it is easily converted to direct binary code, and (2) the finest tracks are twice the width of equivalent direct binary code tracks. Figure 6.69(a) shows a Gray code linear scale. Figure 6.69(b) shows a scheme for the conversion from Gray code to binary code and proceeds as follows: (1) the most significant bit (msb) of the binary code equals the msb of the Gray-coded number; (2) add (modulo-2)

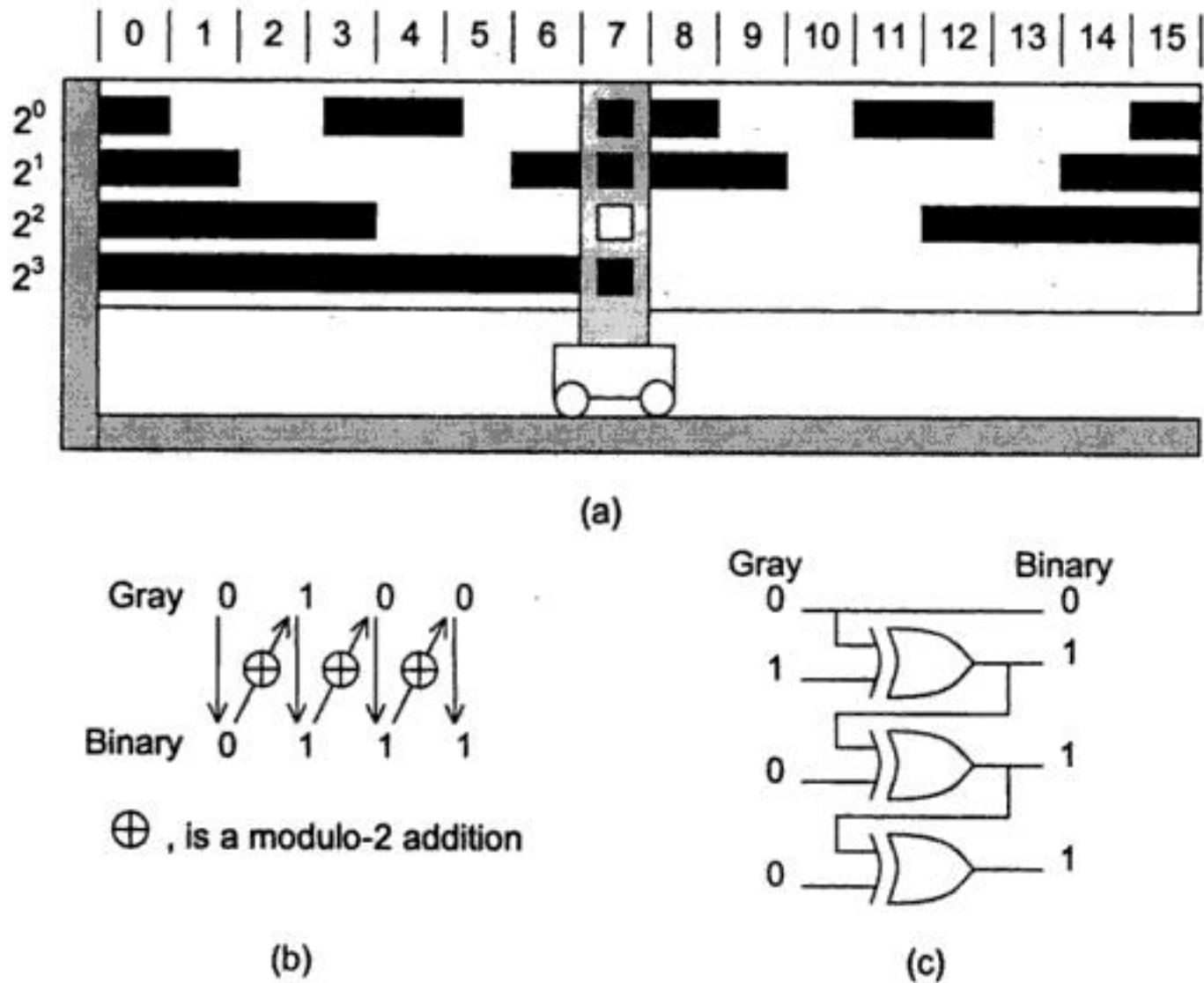


FIGURE 6.69 (a) Gray code allows a transition on only one track between each successive position so that no ambiguity arises. A scheme based on modulo-2 additions converts the Gray code to natural binary code in (b). Exclusive-ORs implement the conversion in (c).

the msb of the binary number to the next significant bit of the Gray-coded number to obtain the next binary bit; and (3) repeat step (2) until all bits of the Gray-coded number have been added modulo-2. The resultant number is the binary equivalent of the Gray-coded number. The modulo-2 addition is equivalent to the action of an exclusive-OR. Figure 6.69(c) shows a simple circuit using combinational logic to perform Gray to binary conversion. Sente et al. [2] suggest the use of an external ROM to convert the code where the input coded word is the address and the data is the output coded word. Using a 16-bit ROM for a 12-bit code, Sente et al. [2] suggest using the remaining 4 bits to implement a direction signal. Stephens et al. [1] describe the use of vee-scan with a gray code scale for even better robustness.

Pseudorandom Code

Pseudorandom encoding allows the use of only two tracks to produce an absolute encoder. One track contains the pattern used to identify the current position, while the other is used to synchronize the reading of the encoded track and remove ambiguity problems. A pseudorandom binary sequence (PRBS) is a series of binary records or numbers, generated in such a way that any consecutive series of n digits is unique. Such code is called chain code, and it has the property that the first $n-1$ digits of an n -bit word are identical to the last $n-1$ digits of the previous code word. This allows their partial overlapping on a single track. A PRBS of length $2^n - 1$ is defined by:

$$XN(j) \mid j = 0, 1, \dots, 2^n - 1 \quad (6.94)$$

The code can be generated by reading the n th stage of a feedback shift register after j shifts. The register must be initialized so that at least one of the registers is nonzero and the feedback connection implements the formula

$$X(0) = X(n) \oplus c(n-1) X(n-1) \oplus \dots \oplus c(1) X(1) \quad (6.95)$$

TABLE 6.18 Shift-Register Feedback Connections for Generating Pseudorandom Binary Sequences

n	Length	Direct sequence	Reverse sequence
4	15	1, 4	3, 4
5	31	2, 5	3, 5
6	63	1, 6	5, 6
7	127	3, 7	4, 7
8	255	2, 3, 4, 8	4, 5, 6, 8
9	511	4, 9	5, 9
10	1023	3, 10	7, 10
11	2047	2, 11	9, 11
12	4095	1, 4, 6, 12	6, 8, 11, 12
13	8191	1, 3, 4, 13	9, 10, 12, 13
14	16,383	1, 6, 10, 14	4, 8, 13, 14

where the c coefficients are 0 or 1. The feedback registers for which c is 1 are listed in Table 6.18 for values of n from 4 to 14. The following is a PRBS for $n = 4$ with the pseudocode obtained using all registers set to 1 initially, 111101011001000. For a rotary encoder disk, the 15 sectors would have a 24° width.

Petriu [3, 4] describes a possible configuration of a PRBS disk that uses a PRBS track and a synchronization track (Figure 6.70), together with the processing method to reconstitute the position in natural binary. Table 6.18 gives the reverse feedback configuration. The shift register is initially loaded with the current n -tuple. Then the reverse logic is applied recurrently until the initial sequence of the PRBS is reached. At this point, the n -bit counter represents the value of j . For a rotary encoder $(j * 360)/(2^n - 1)$ is the current angular position; whereas for a linear encoder, the position is $j * P$ where P is the scale record length or pitch. Petriu [4] suggests that in order to allow nonambiguous bidirectional reading, $n + 1$ heads are used on the PRBS track. The synchronization track has a series of 0s and 1s in records of the same width as the PRBS track and in phase. There is a $P/2$ shift, where P is the record's length between the A head on the synchronization track and the $n + 1$ read heads on the PRBS track. The $n + 1$ records are updated on a trigger from the A signal. This ensures that the $n + 1$ heads are closely aligned with the PRBS record mid-position. A second read head called B on the synchronization track is shifted by $P/2$ relative to A. A and B are in quadrature, which allows their simultaneous use to generate a motion direction signal. The correct n -tuple, i.e., the lower or upper subset, is selected on the basis of the moving direction and is then converted to natural binary by reverse feedback. Petriu [4] also suggests a simple means of increasing the resolution by a factor of 2 using some additional electronics. He also proposes a scheme to use an arbitrary (not $2^n - 1$) number of sectors, but this requires a third track and some additional correction electronics to handle the last $n - 1$ records of a disk, since these are no longer PRBS patterns. Tomlinson [5] proposes another method for truncation of the PRBS sequence that does not require a third track. Instead, particular codes were removed by applying additional logic in the direct and reverse feedback logic.

Ross and Taylor [6] and Arsic and Denic [7] suggest ways of reducing the number of reading heads by accumulating readings into a shift register so that a minimum of two heads are sufficient to read the PRBS track. However, on start-up, the correct position is not known until the encoder has moved so that all registers have been updated. This type of encoder is therefore not completely absolute because it does not indicate its correct position on start-up. Finally, Arazi [8] mentions the use of a ROM that stores the translation table, as an alternative to a logic circuit.

Optical Resolving

This method has similarities with its electromagnetic counterpart and depends on the generation of a sine and a cosine signal pair per encoder shaft revolution. The resolution and accuracy of this encoder depend on its ability to generate signals that conform to their ideal waveforms and the resolving power of the electronic circuit responsible for performing the rectangular to polar conversion to produce angular

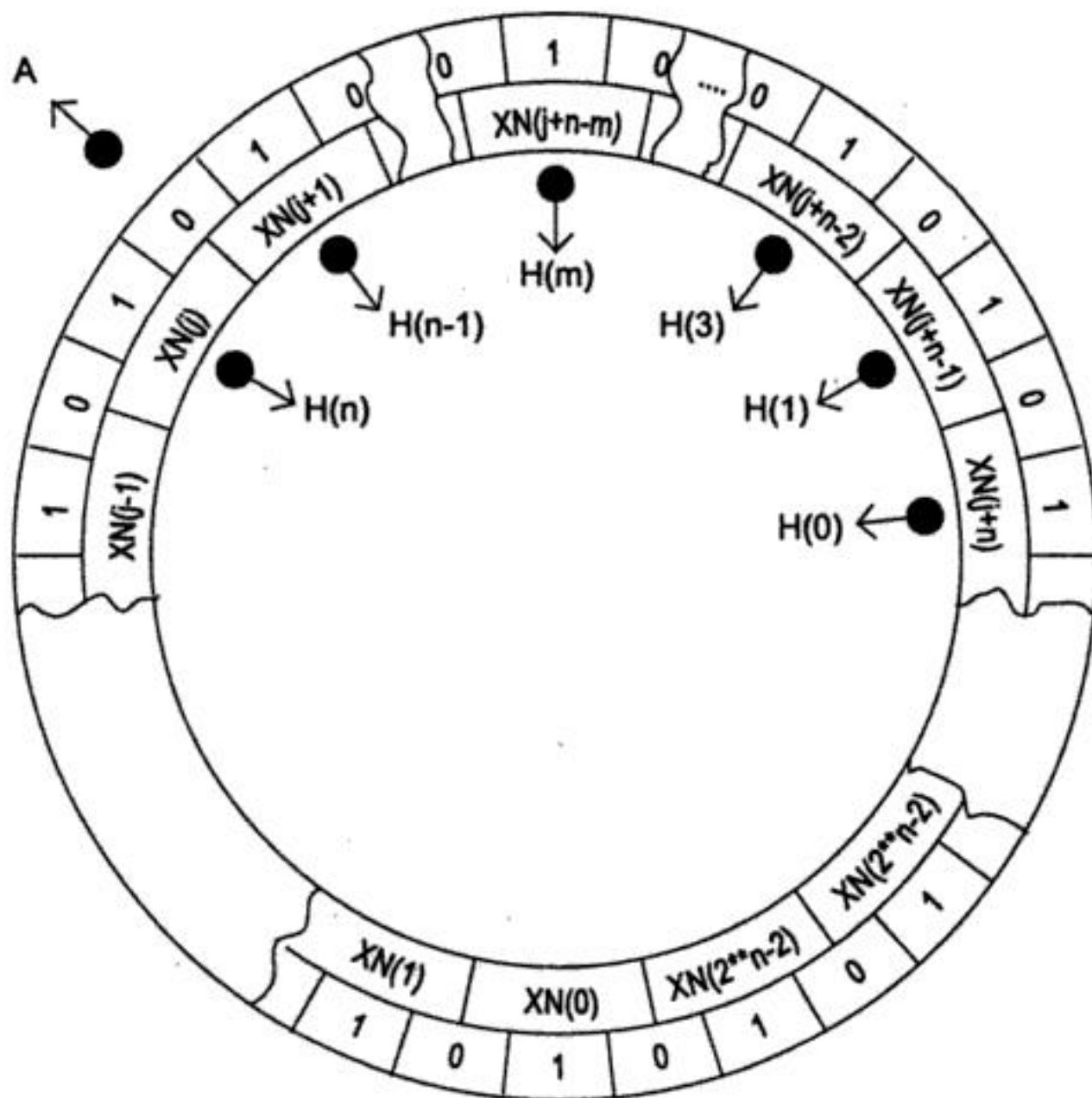
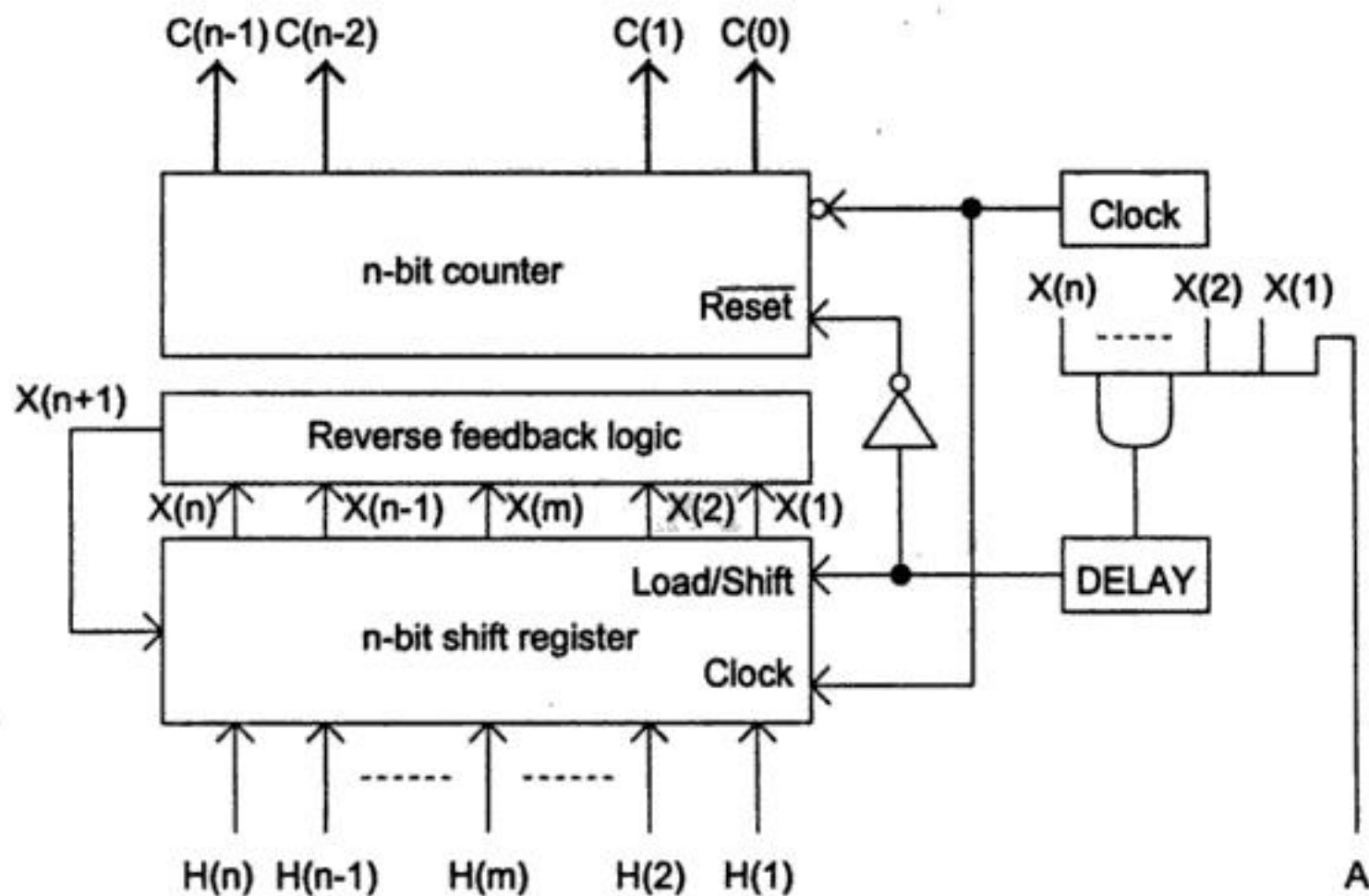


FIGURE 6.70 Pseudorandom shaft encoder with a simple synchronization track to validate the readout from the n -tuple read by the reading heads $H(i)$, $i = n, \dots, 1$. The circuit is a simplified code conversion to binary based on the reverse feedback logic. The counter counts the number of steps required to return to the initial sequence.

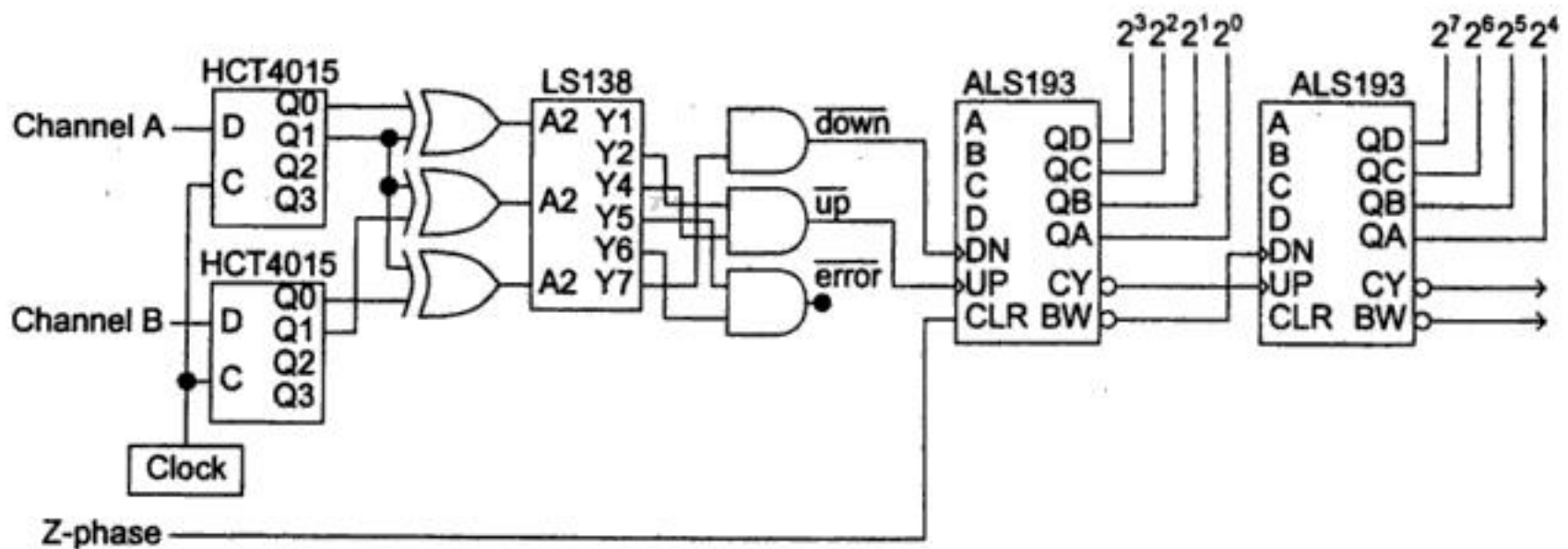


FIGURE 6.73 Divide-by-four circuit producing four counts per cycle of the quadrature signals. Based on Butler's design [14], except for the LS138 demultiplexer that replaces the suggested 4051 multiplexer because of the latter's unequal fall time and rise time. The clocked shift registers generate a 4-bit code to the demultiplexer for each clock cycle. The clock frequency should be 8 times the maximum frequency of the quadrature signals.

circuit drives a 74HC193 counter with a down clock and an up clock. Both flip-flops are normally in the high state. Whenever one of them switches to a low state, it is soon reset by the use of that output as a Set signal. That low output is only possible when phase B is low and there is a transition on A. Depending on the direction of that transition, one of the flip-flops produces the brief low state, causing an appropriate up or down count. The problem associated with small oscillations around a position is reduced since the up and down counts occur at the same encoder position. Venugopal [11] proposes a circuit, in Figure 6.72(c), that produces a similar effect. A count can only occur when B is low and there is a transition on A. Depending on the transition direction, one of two monostables triggers and effectively allows the count at the appropriate input of the 74LS193 counter.

In cases where a noisy environment is present, Holle [12] describes a digital filter to clean the A and B signals further, a filter that uses a small number of logic gates to form a 4-bit delay filter. Wigmore [13] proposes other circuits for count and direction detection.

Tables 6.23 and 6.24 list commercial chips and their suppliers' details. Some chips produce count and direction signals or up and down clocking signals. Others also include counters.

The above circuits do not fully exploit the information contained in the quadrature signals since only one of the four edges (or states) within one quadrature cycle is used to count. Figure 6.73 shows a slightly modified version of a divide-by-four counter circuit proposed by Butler [14]. Two 4-bit shift registers, three exclusive-OR gates, and an eight channel demultiplexer derive the up and down count signals. The clocked shift registers generate a 4-bit code to the demultiplexer for each clock cycle. To ensure that no encoder transitions are missed, the clock frequency should be at least $8NS$, where N is the number of cycles produced by the encoder for each shaft revolutions and S is the maximum speed in revolutions per second. The up and down signals can be fed to cascaded 74ALS193 counters. This circuit analyzes the current and previous states of the A and B channels to either count up, count down, or issue an error flag when an improper sequence occurs.

Kuzdrall [10] proposes to view a quadrature signal cycle as a Gray code so that the two least significant bits of the count are obtained by a Gray code to binary code conversion as in Figure 6.72(b). Phase A generates bit 1 of the natural binary position, and phase B is exclusive-ORed with phase A to produce the least significant bit. The counter provides the remaining bits of the natural binary position. Marty [15] proposes a state machine that stores both the actual A and B values and the previous values using D-type flip-flops (in a similar way to Butler) to form a hexadecimal number. In all, eight different hex-digits are generated, four for each direction of motion. A 4-line to 16-line decoder then feeds into two 4-input NAND gates to produce an up or down count signal. As for Butler, this last circuit is not dependent on propagation delays as with Kuzdrall's circuit.

Analog Quadrature Signals

Interpolation by resistor network

Some optical encoders deliver analog quadrature signals in the form of a $\sin(\theta)$ signal, A, and a $\cos(\theta)$ signal, B, where θ is the phase (electrical degree) within one cycle of the quadrature signal. θ does not equal shaft angle but is related to it: 360 electrical degrees of θ corresponds to $360/N$ mechanical degrees, where N is the number of analog quadrature signal cycles per shaft revolution. Indeed, θ still exists for linear encoders of this type. Although these signals can be squared and fed to a divide-by-four counter, they can also be processed directly to generate a finer resolution. The main techniques are (1) multiple phase shifted signals, (2) lookup table, and (3) arctangent processor.

The multiple phase-shifted signals method relies solely on electronics to increase the frequency of the final digital quadrature signals by an integer amount. Benzaid et al. [16] propose the circuit of Figure 6.74(a), which has been designed in this particular case for a fourfold frequency increase. This can then be followed by a digital divide-by-four circuit (not shown). The A and B signals are combined to produce an additional six phase-shifted signals, three of which are shifted by $\pi/8$, $\pi/4$, and $3\pi/8$, respectively, from the A signal, and three others that are shifted similarly from the B signal. The result is a total of 8 available sinusoidal signals, phase-shifted by $\alpha = i\pi/8$ with $i = 0, 1, \dots, 7$. This can be generalized, saying that for an m -fold increase in resolution, $2m$ signals that are phase-shifted by $\alpha i\pi/2m$ with $i = 0, 1, \dots, 2m - 1$ are required. The phase-shifted sinusoidal signals are then squared using TTL converters. Figure 6.74(b) shows the resulting square waves. The addition, modulo-2, of the signals for even values of i gives A' , and similarly for the signals with odd values of i gives B' . The modulo-2 sum is performed via exclusive-ORs.

The vector additions of the initial A and B signals result in phase-shifted signals. The weights of A and B are calculated from the trigonometric relations:

$$\sin(\theta + \alpha) = \cos(\alpha) \sin(\theta) + \sin(\alpha) \cos(\theta) \quad (6.96)$$

and

$$\cos(\theta + \alpha) = -\sin(\alpha) \sin(\theta) + \cos(\alpha) \cos(\theta) \quad (6.97)$$

where $\sin(\theta) = A$ and $\cos(\theta) = B$.

Thus,

$$\sin(\theta + \alpha) = \cos(\alpha) A + \sin(\alpha) B \quad (6.98)$$

and

$$\cos(\theta + \alpha) = -\sin(\alpha) A + \cos(\alpha) B \quad (6.99)$$

Note that the amplitudes of the phase-shifted signals are not critical, since it is their zero-crossing points that produce the square-wave transition. Also, in the circuit of Figure 6.74(a), the weights were slightly modified to simplify the implementation, which results in a small variation of the duty cycle of A' and B' .

The phase-shifted signals can alternatively be produced using voltage dividers and Schmitt triggers, where the divider resistors are in the ratio $\tan(\alpha)$ [17]. As m increases, the precision of the weights become more stringent and the speed of the electronics processing the high-frequency square signals might limit the upper value of m possible with this method.

Interpolation by Sampling and Numerical Processing

A number of methods digitize the analog quadrature signals in order to perform a digital interpolation within one cycle of the quadrature signals. These techniques permit an even higher interpolation. The signals are periodically sampled in sample-and-hold circuitry and digitized by an analog to digital converter (ADC). One technique uses an interpolation table, while the other performs arctangent calculations.

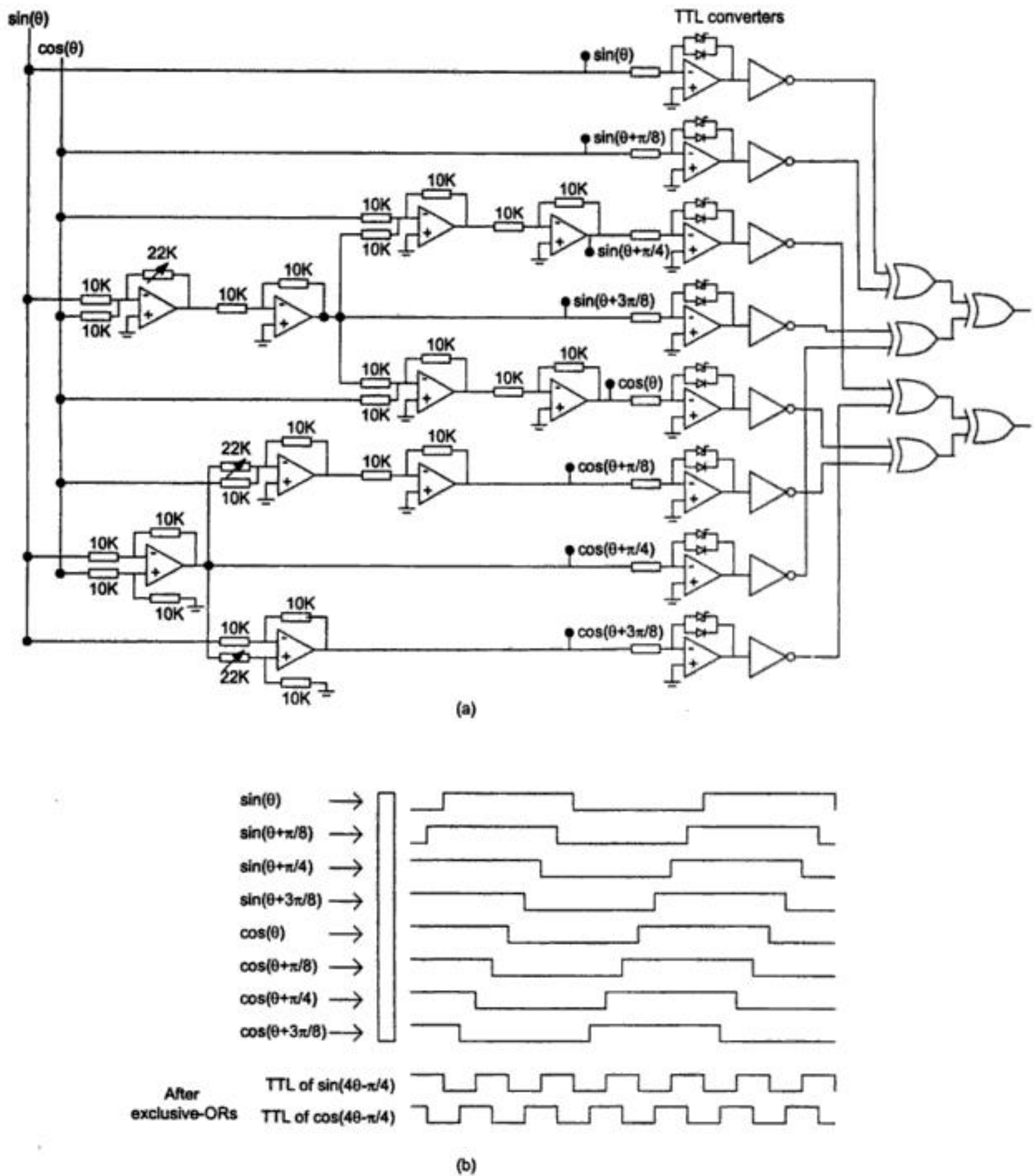


FIGURE 6.74 Interpolation circuit by Benzaid et al. [16]. The sine and cosine signals from the encoder are vector-combined in (a) to produce various phase-shifted signals. Following TTL conversion, they are exclusive-ORed. The squared phase-shifted signals and the final result are in (b).

Hagiwara [18] uses an m -bit ADC to digitize the analog A and B signals into the binary numbers D_a and D_b . Figure 6.75(a) shows how these binary numbers are used as the addresses of a grid, here built for $m = 3$ as a 2^m by 2^m grid. Hagiwara then associates a phase angle with the center of each cell of the matrix using simple arctangent calculations. The angle is calculated with respect to the center of the grid. This phase angle is then associated with one of $2^n - 1$ phase codes using:

$$\text{Phase code} = \text{integer part of } \left(2^n \theta / 2\pi \right) \tag{6.100}$$

as shown in Figure 6.75(b).

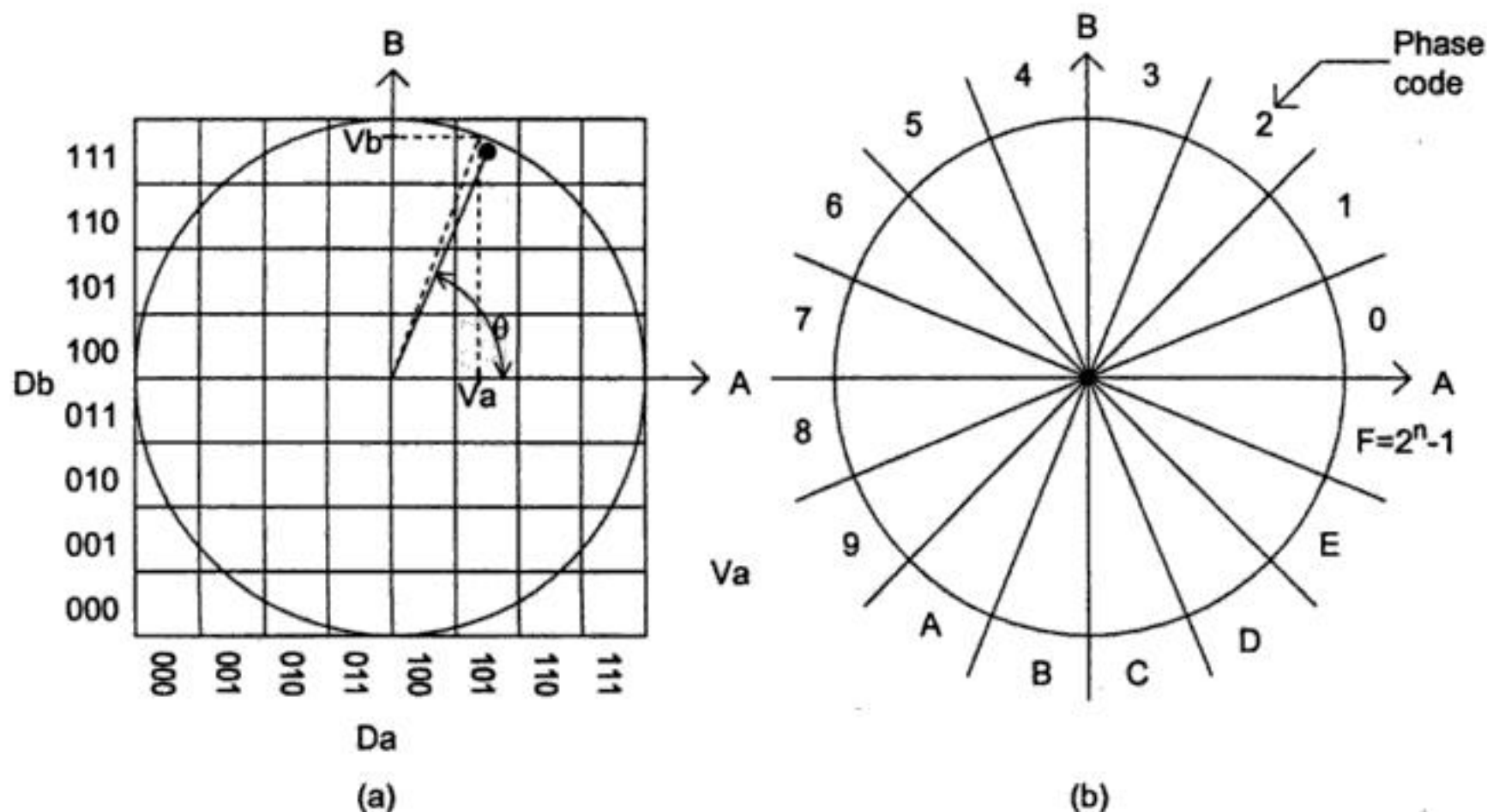


FIGURE 6.75 Hagiwara [18] proposes a two-dimensional look-up table that associates a quantized phase value to a set of digitized values of the quadrature signals. The phase code associated with a grid address may be adjusted to compensate for known errors in the quadrature signals.

The proposed circuit includes the means of performing the accumulation of interpolated values to deliver the encoder absolute position. A modification to this circuit is also proposed that can compensate the phase code for known inaccuracies in the encoder signals with respect to the actual encoder position.

For the highest possible level of interpolation, Mayer [19] describes an arctangent calculator method. It uses a microprocessor to perform high-level trigonometric calculations on the digitized analog signals. As this process is performed, the analog signals are continuously squared and fed to a divide-by-four counter. The two information sources are combined to produce the encoder's position with a very high level of resolution. Note, however, that although the resulting theoretical resolution is limited only by the ADC, in practice it is the quality of the analog signals in relation to the physical position being measured that will limit the precision obtained.

Encoding Principles

Optical encoders use one of three techniques to generate the electrical signals from the relative movement of the grating and the reading heads. They are (1) geometric masking, (2) Moiré effects, and (3) laser interference. Note that absolute encoders mainly use geometric masking. Geometric masking relies on geometric optics theory and considers light as traveling in a straight line. However, as the grating period reduces and the resolution of the encoder increases, Moiré fringe effects are observed and used to produce the signals. Although diffraction effects then become nonnegligible, their influence can be controlled by careful design. Finally, for very high resolution, diffraction effects are directly exploited to perform the measurements.

Geometric Masking

Geometric masking is applied to absolute and incremental encoders. The electromagnetic field associated with the propagation of visible light is characterized by very rapid oscillations (frequencies of the order of 10^{14} s^{-1}). It may therefore be expected that a good first-order approximation to the propagation laws of light is obtained by neglecting the wavelength of light. In this case, diffraction phenomena may be ignored and light may be thought to propagate in a straight line. Geometry can then be used to analyze

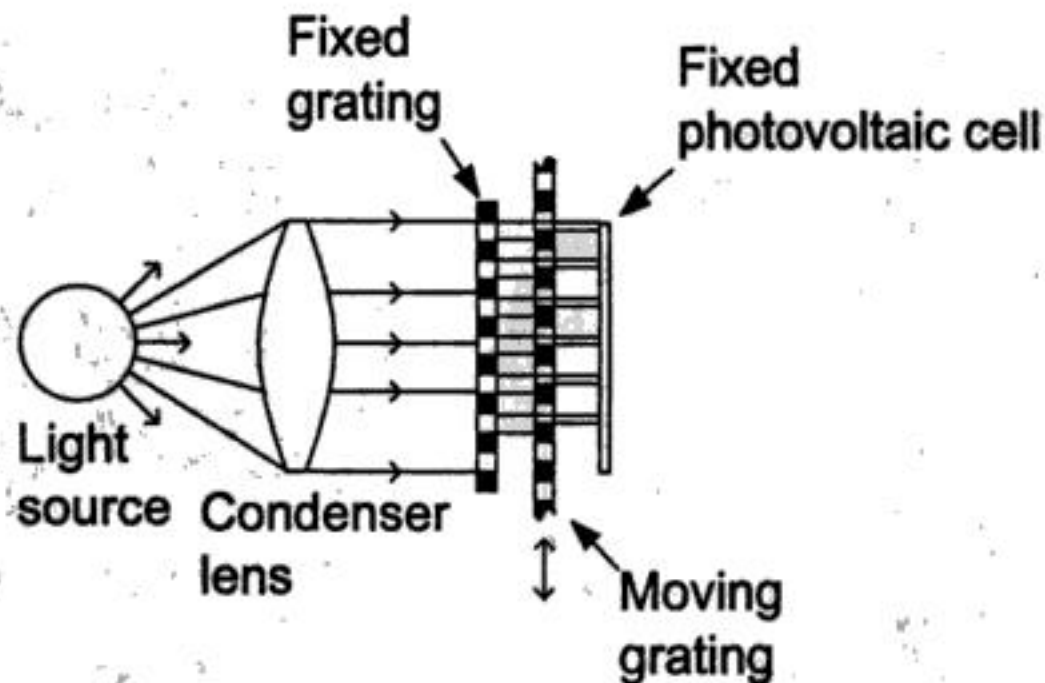


FIGURE 6.76 The reading head contains a light source that can be an incandescent bulb or a light emitting diode. The light may require a condenser in order to collimate it. This redirects the light rays from the source so that they travel perpendicular to the gratings. The stationary index grating structures the light beam into bands that then cross the moving grating. Depending on the relative position of the fixed and the moving gratings, more or less light reaches the photodetector. The detector produces a signal proportional to the total amount of light on its sensitive surface.

the behavior of light. The approximation is valid whenever light rays propagate through an encoder grating with fairly coarse pitch (say, more than $10\ \mu\text{m}$). Figure 6.76 shows a portion of an encoder scale being scanned by a reading head with multiple slits in order to send more light onto the photosensitive area. The light source is collimated (rays are made parallel to each other after having been emitted in multiple directions from a small source), passes through the stationary grating or index grating, and propagates through the moving grating. In the case of an incremental encoder, a second head would read from the same moving grating but be displaced by $n + 1/4$ pitch in order to produce a second signal with a phase of 90° . When the slits of the two gratings are aligned, a maximum amount of light reaches the detector. Similarly, a minimum amount of light is transmitted when the gratings are out of phase by 180° . Depending on the design of the grating, such as the duty cycle of the transmitting and opaque sectors, various cyclic signal patterns are obtained. Generally, it is noteworthy that optical encoders can also be used as a reflective as opposed to transmissive design. Also, in some absolute optical encoders, the encoding principle does not rely on a primary grating, per se, or the detection system does not necessarily incorporate any masking element.

Moiré Fringes

Moiré fringe methods are primarily associated with incremental encoders. Moiré fringes are observed when light passes through two similar periodic patterns brought close to each other and with their line pattern nearly parallel. Figure 6.77 shows linear and radial gratings producing Moiré fringes. Take a linear grating with sinusoidal amplitude transmittance as proposed by Gasvik [20]:

$$f(x, y) = a + a \cos(2\pi x/P) \quad (6.101)$$

where P is the grating period, a is the amplitude, and x is measured perpendicularly to the grating lines. It is also possible to represent a square wave type grating using a Fourier series in the case of a radial grating and a Fourier integral for a linear grating. When two linear gratings, a and b , are laid in contact, the resulting transmittance, f_c , is the product of their individual transmittances f_a and f_b :

$$f_c = f_a \times f_b \quad (6.102)$$

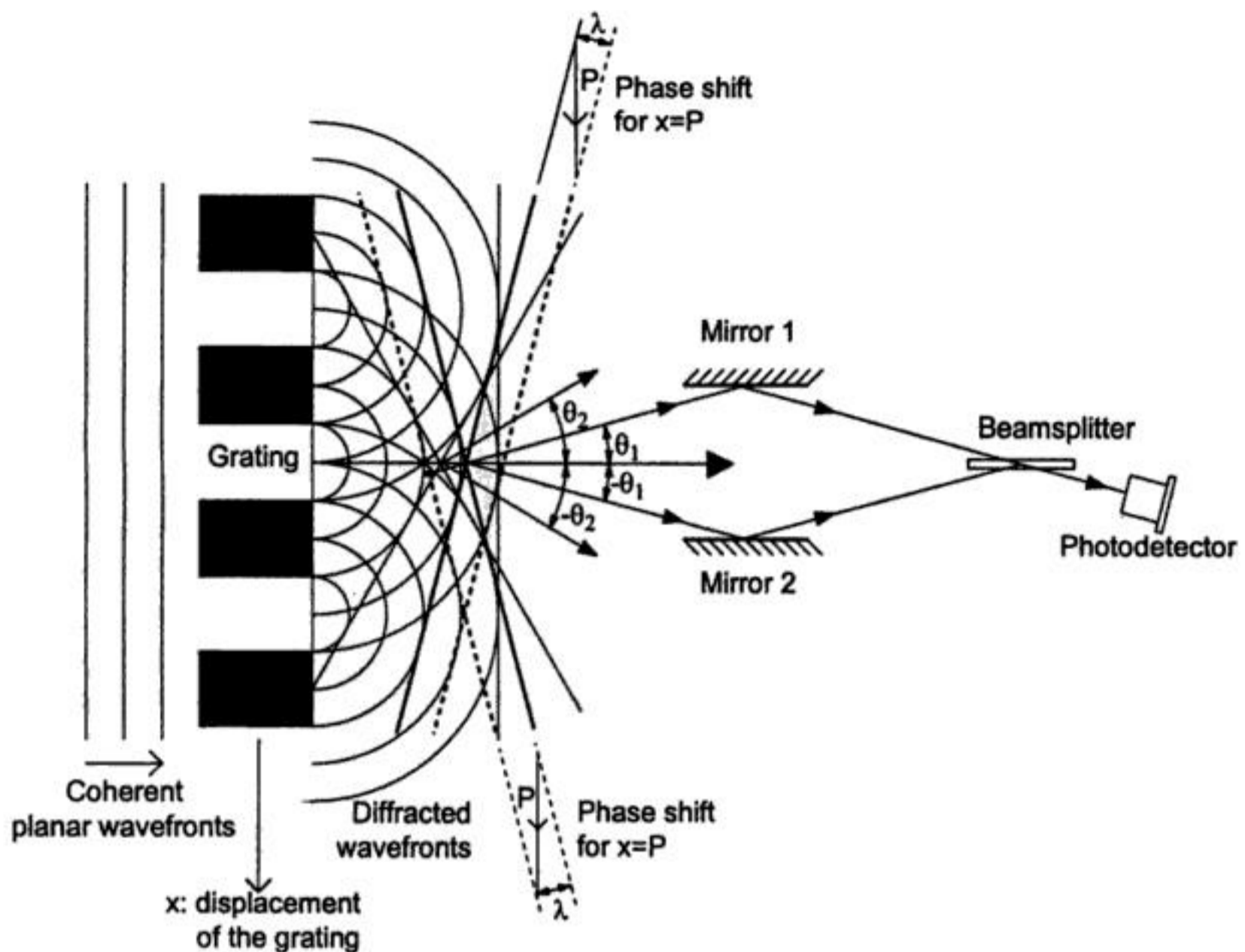


FIGURE 6.78 The grating diffracts the coherent planar wavefront coming from the left into a series of circular wavefronts. The coherent circular wavefronts are in phase along certain directions resulting in diffracted planar wavefronts. These planar wavefronts correspond to the common tangents of the circles. For example, the third innermost circle of the uppermost slit is in phase with the fourth innermost circle of the middle slit and also with the fifth innermost circle of the lowest slit. This diffracted wavefront corresponds to the first order diffraction ($m = 1$) and makes an angle θ_1 with the initial planar wavefront direction. A similar order wavefront ($m = -1$) has a direction $-\theta_1$. These two wavefronts are initially in phase. A displacement of the grating downward by x causes the $m = 1$ wavefront to move by a distance $x\lambda/P$ and the $m = -1$ wavefront by $-x\lambda/P$, thus causing a relative phase shift between the two wavefronts of $2x\lambda/P$. This relative phase shift produces a movement of interference fringes at the photodetector.

where m is a positive or negative integer and is the order of diffraction, and P is the diffraction grating period. Suppose now that the slits are moving by a distance x as shown in Figure 6.78. Then the phase of the wavefront at a stationary location will increase for m positive and decrease for m negative by $2\pi mx/P$. When the slit pattern has moved by one pattern cycle, then the two wavefronts will have developed a relative phase change of $2m$ cycles. The two phase-shifted wavefronts may be recombined to produce interference fringes.

Rotary Encoders

The Canon laser rotary encoder uses an optical configuration that generates four cycles of interference fringes per cycle of the diffraction grating. This is achieved by splitting the original wavefront in two and interrogating the grating at two diametrically opposed locations. Reflecting both diffracted beams back through the grating results in a further doubling of the resolution. The two beams are finally recombined for interference fringe counting. Furthermore, as explained by Nishimura and Ishizuka [23], using diametrically opposed portions of the disk attenuates the effect of eccentricity. An encoder with an external diameter of 36 mm produces 81,000 analog quadrature cycles per revolution. The grating has a pitch of approximately $5 \mu\text{m}$ and the laser light has a wavelength of 780 nm.

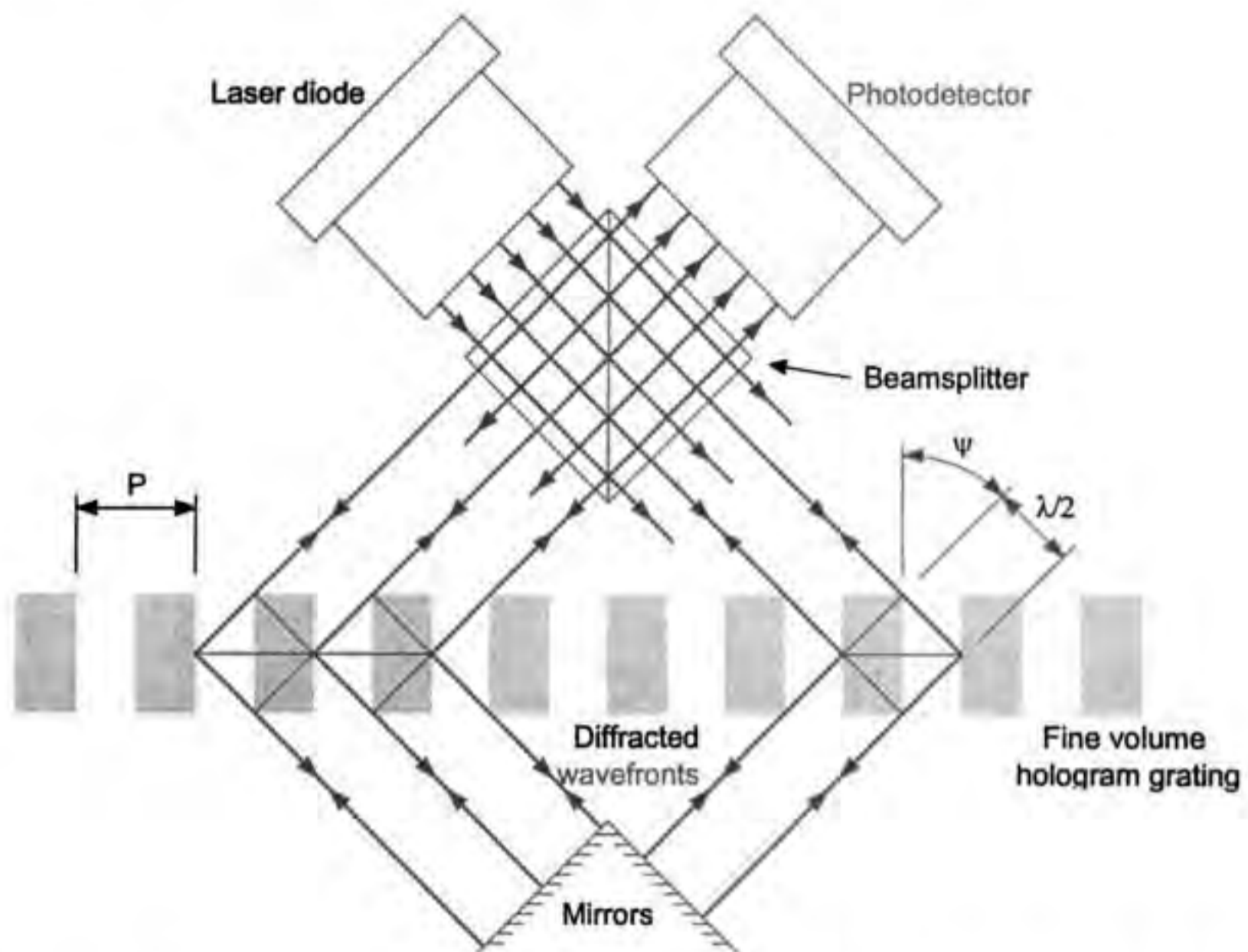


FIGURE 6.79 The coherent planar wavefront of the laser diode is split into two at the beam splitter. Both beams can be regarded as being partially reflected at the grating by minute mirrors separated by a pitch, P . Maximum reflection is achieved when $2P \sin \psi = \lambda$. Each reflected portion of the incident beam is phase shifted by one λ with respect to that reflected by an adjacent mirror so that the reflected wavefront remains coherent. The reflected wavefronts are truly reflected at the orthogonal mirrors for a second pass through the grating, after which they are recombined by the splitter to interfere at the photodetector. A displacement of the grating by x along its axis causes one beam path to be shortened by $2x\lambda/P$, while the other is lengthened by the same amount. A grating motion by P causes four fringe cycles.

Linear Encoders

Sony markets a laser linear encoder with a volume hologram grating of $0.55 \mu\text{m}$ pitch [24]. Figure 6.79 gives some insight into the principle of operation of this type of encoder. In theory, for maximum intensity of the reflected light beam through the hologram grating, the path length difference from successive plane mirrors must be equal to λ [20]. The Sony Laserscale encoder uses a semiconductor laser as the coherent light source. The beam splitter produces two beams to interrogate the grating at two separate locations. The beams are then diffracted by the hologram grating, followed by a reflection at the mirrors and pass a second time through the grating. They are finally recombined at the photodetector to produce interference fringes. Because the two beams are diffracted by the grating in opposite directions, and because they pass twice through the grating, four signal cycles are obtained when the grating moves by one pitch.

The Heidenhain company produces both linear and rotary encoders of very high resolution, also using the principles of diffraction. Their approach, unlike those previously mentioned, uses an index grating and a reflecting scale grating [17, 25].

The Renishaw company uses a fine reflective grating of $20 \mu\text{m}$ pitch that diffuses the light from an infrared light-emitting diode. In order to avoid the problems caused by diffraction effects, the index

TABLE 6.19 Commercial Optical Rotary Incremental Encoders

Manufacturer	Model No.	Output type	Counts ^a	Resolution ^b	Price ^c
BEI	H25	Square wave	2540	50,800	\$340
Canon	K1	Sine wave	81,000	1,296,000	\$2700
Canon	X-1M	Sine wave	225,000	18,000,000	\$14,000
Dynamics	25	Sine wave	3000	60,000	\$290
Dynamics	35	Sine wave	9000	360,000	\$1150
Gurley	911	Square wave	1800	144,000	\$1300
Gurley	920	Square wave	4500	144,000	\$500
Gurley	835	Square wave	11,250	360,000	\$1500
Heidenhain	ROD 426	Square wave	50 to 10,000	200 to 40,000	\$370
Heidehain	ROD 905	Analog	36,000	0.035 arcsec	\$12,600
Lucas Ledex	LD20	Square wave	100	—	\$180
Lucas Ledex	LD20	Square wave	1000	—	\$195
Lucas Ledex	DG60L	Square wave	5000	—	\$245
Renco	RM21	Square wave	2048	—	\$176
TR Electronic	IE58	Square wave	10,000	—	\$428

^a Number of quadrature cycles per revolution without electronic divide-by-four or interpolation.

^b Unless otherwise specified, is the number of counts per revolution with electronic interpolation, either internal or external to the encoder, supplied by the manufacturer as an option.

^c Based on orders of one unit. Must not be used to compare products since many other characteristics, not listed in this table, determine the price.

grating of a conventional Moiré approach would require a very small gap between the index and main gratings. Instead, the index grating is located at a distance of 2.5 mm. The index grating is then able to diffract the diffused light from the main grating and to image it 2.5 mm further. There, a fringe pattern of the Moiré type is produced and photoelectrically analyzed [26].

Components and Technology

The choice of an encoder demands careful consideration of a number of factors, such as: (1) the required resolution, repeatability, and accuracy (linearity); (2) the maximum and minimum operating speeds; (3) the environmental conditions: temperature range, relative humidity (condensing or not condensing), contaminants such as dust, water, oil, etc.; (4) minimum friction torque (or force) acceptable; (5) maximum inertia acceptable; (6) available space; (7) position to be known immediately after a power loss; (8) range in degrees (or mm); (9) mounting and shaft loading; and (10) price.

Rotary encoders also sometimes employ reflective tapes attached to, or markings directly etched into, the surface of a drum or spindle to be read from the side. In special cases, tapes or engravings of this type can be made to conform to the surface of an elliptical cross-section cam or other irregular contour. Cylindrical primary gratings employing transmissive readout have also been produced.

The main gratings (scale) come in a wide variety of materials, both for the substrate and for the marking. Flexible scales based on a metal tape substrate are also available from Renishaw, allowing very long (tens of meters) continuous reading of linear motion. Tables 6.19 through 6.21 list a number of encoders currently available on the market. The list is not exhaustive in terms of suppliers and does not cover all the models of the suppliers listed. Table 6.22 gives the address and telephone numbers of the suppliers. Tables 6.23 and 6.24 list some suppliers of quadrature decoding circuits.

TABLE 6.20 Commercial Optical Rotary Absolute Encoders

Manufacturer	Model Number	Steps per turn	No. of turn	Price ^a
BEI	M25	65,536	1	\$2130
BEI	MT40	512	16	\$1240
BEI	MT40	65,530	512	\$5000
Gurley	25/04S	131,072	1	\$1900
Heidenhain	ROC 424	4096	4096	
Lucas Ledex	AG60E	360 or 512	1	\$486
Lucas Ledex	AG661	4096	4096	\$1260
TR Electronic	CE65 ^b	8192	4096	\$1408

^a Based on orders of one unit. Must not be used to compare products since many other characteristics, not listed in this table, determine the price.

^b Programmable output.

TABLE 6.21 Commercial Optical Linear Incremental Encoders

Manufacturer	Model No.	Output type	Pitch ^a	Resolution ^b	Length (mm)	Price ^c (length)
Canon	ML-16+	Sine wave	1.6 μm	0.4 μm	To 300	\$1525 (50 mm)
Canon	ML-08+	Sine wave	0.8 μm	0.2 μm	To 150	\$3100
Gurley	LE18	Square wave	20 μm	0.1 μm	To 1500	\$750 (1000 mm)
Gurley	LE25	Square wave	20 μm	0.1 μm	To 3000	\$800 (1000 mm)
Heidenhain	LS603	Sine wave	20 μm	5 μm	To 3040	\$932 (1020 mm)
Heidenhain	LIP401	Sine wave	2 μm	0.005 μm	To 420	\$4000 (100 mm)
Renishaw	RG2	RS422A	20 μm	0.5 μm	To 60,000	\$640 + \$360/1000 mm
Sony	BS75A-30NS	Square wave	0.14 μm	0.05 μm	30	\$2628

^a Period of the quadrature cycle without electronic divide-by-four or interpolation.

^b With electronic interpolation supplied by the manufacturer.

^c Based on orders of one unit. Must not be used to compare products since many other characteristics, not listed in this table, determine the price.

TABLE 6.22 Companies that Make Optical Encoders

BEI Sensors and Motion Systems Company Encoder Systems Division 13100 Telfair Avenue Sylmar, CA Tel: (848) 341-6161	Renco Encoders Inc. 26 Coromar Drive Goleta, CA 93117 Tel: (805) 968-1525
Canon USA Inc. Components Division New York Headquarters : One Canon Plaza Lake Success, NY 11042 Tel: (516) 488-6700	Renishaw plc, Transducer Systems Division Old Town, Wotton-under-Edge Gloucestershire GL12 7DH United Kingdom Tel: +44 1453 844302
DR. JOHANNES HEIDENHAIN GmbH DR.-Johannes-Heidenhain-Strasse 5 D83301 Traunreut, Deutschland Tel: (08669)31-0	TR Electronic GmbH Eglishalde 6 Postfach 1552 D-7218 Trossingen Germany Tel: 0 74 25/228-0
Gurley Precision Instruments Inc. 514 Fulton Street Troy, NY 12181-0088 Tel: (518) 272-6300	Sony Magnescale Inc. Toyo Building, 9-17 Nishigotanda 3-chome Shinagawa-ku, Tokyo 141 Japan Tel: (03)-3490-9481
Ledex Products Lucas Control Systems Products 801 Scholz Drive P.O. Box 427 Vandalia, OH 45377-0427 Tel: (513) 454-2345	

TABLE 6.23 Commercial Digital Quadrature Signal Decoder Circuits

Manufacturer	Model No.	Output	Decoding factor	Counter	Price
Hewlett Packard	HCTL-2000	Count	12-bit	×4	\$12.75
Hewlett Packard	HCTL-2016	Count	16-bit	×4	\$12.75
Hewlett Packard	HCTL-2020	Count	16-bit & cascade o/p	×4	\$14.55
U.S. Digital Corp.	LS7083	Up and Down clock		×1 or ×4	\$3.05
U.S. Digital Corp.	LS7084	Count and direction		×1 or ×4	\$3.60

TABLE 6.24 Companies that Make Divide-by-Four Decoders

U.S. Digital Corporation 3800 N.E. 68 th Street, Suite A3 Vancouver, WA 98661-1353 Tel: (360) 696-2468
Hewlett-Packard Company Direct Marketing Organization 5301 Stevens Creek Boulevard P.O. Box 58059, MS 51LSJ Santa Clara, CA 95052-8059 Tel: (408) 246-4300

References

1. P. E. Stephens and G. G. Davies, New developments in optical shaft-angle encoder design, *Marconi Rev.*, 46 (228), 26-42, 1983.
2. P. Sente and H. Buyse, From smart sensors to smart actuators: application of digital encoders for position and speed measurements in numerical control systems, *Measurement*, 15(1), 25-32, 1995.
3. E. M. Petriu, Absolute-type position transducers using a pseudorandom encoding, *IEEE Trans. Instrum. Meas.*, IM-36, 950-955, 1987.
4. E. M. Petriu, Scanning method for absolute pseudorandom position encoders, *Electron. Lett.*, 24, 1236-1237, 1988.
5. G. H. Tomlinson, Absolute-type shaft encoder using shift register sequences, *Electron. Lett.*, 23, 398-400, 1987.
6. J. N. Ross and P. A. Taylor, Incremental digital position encoder with error detection and correction, *Electron. Lett.*, 25, 1436-1437, 1989.
7. M. Arsic and D. Denic, New pseudorandom code reading method applied to position encoders, *Electron. Lett.*, 29, 893-894, 1993.
8. B. Arazi, Position recovery using binary sequences, *Electron. Lett.*, 20, 61-62, 1984.
9. D. Conner, Long-lived devices offer high resolution, *EDN*, 35 (9), 57-64, 1990.
10. J. A. Kuzdrall, Build an error-free encoder interface, *Electron. Design*, September 17, 81-86, 1992.
11. P. Venugopal, Reflective optical SMT module reduces encoder size, *Power Conversion and Intelligent Motion*, 21(5), 60-62, 1995.
12. S. Holle, Incremental encoder basics, *Sensors*, 7(4), 22-30, 1990.
13. T. Wigmore, Optical shaft encoder from sharp, *Elektor Electron.*, 15(169), 60-62, 1989.
14. M. M. Butler, Simplified multiplier improves standard shaft encoder, *Electronics*, November 20, 128-129, 1980.
15. B. Marty, Design a robust quadrature encoder, *Electron. Design*, June 24, 71-72, 74, 76, 1993.
16. O. Benzaid and B. M. Bird, Interpolation techniques for incremental encoders, *Proc. 23rd Int. Intelligent Motion Conf.*, Jun 22-24, 165-172, 1993.
17. Heidenhain General Catalog, Dr. Johannes Heidenhain GmbH, DR.-Johannes-Heidenhain-Strasse 5, D83301 Traunreut, Deutschland, November 1993, 8.

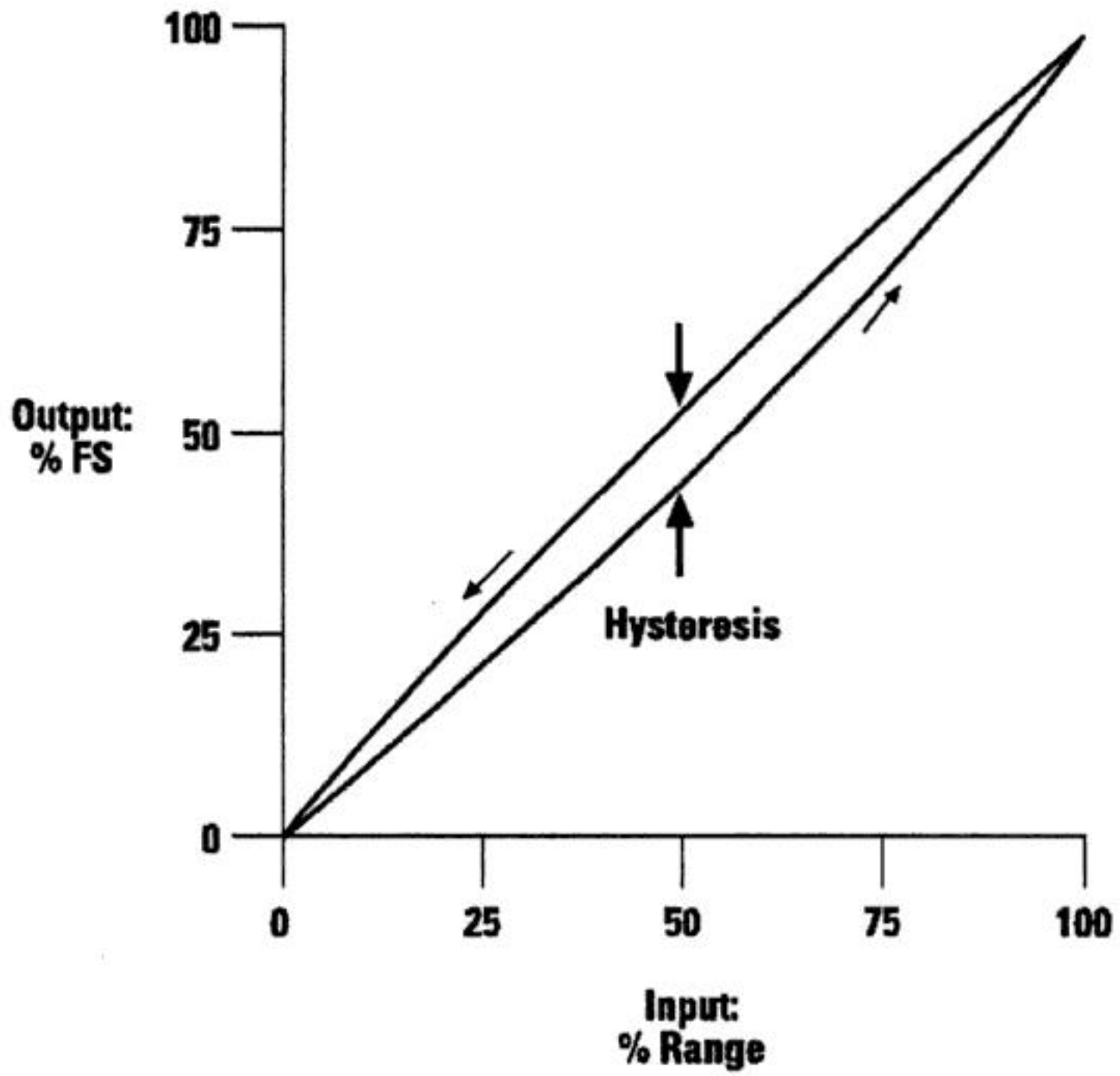


FIGURE 6.80 Hysteresis: output vs. input.

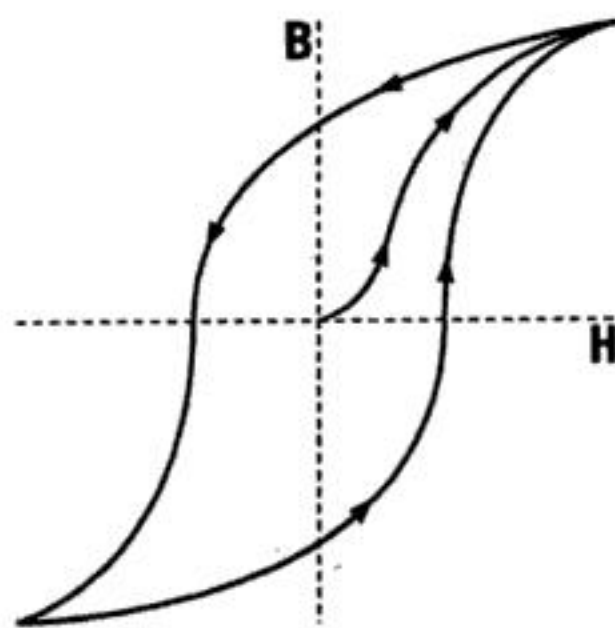


FIGURE 6.81 Magnetic hysteresis.

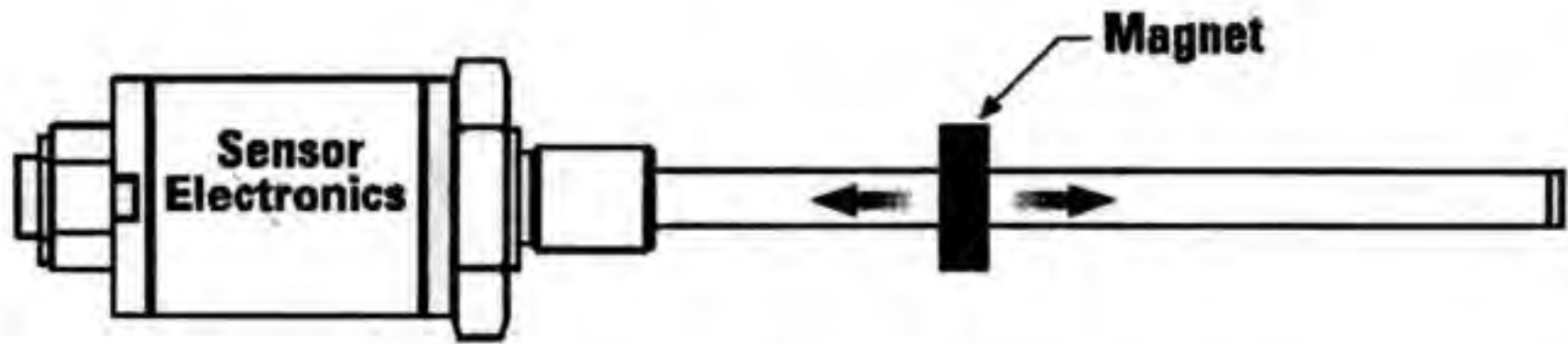


FIGURE 6.82 Magnetostrictive sensor with position magnet.

is to be sensed, and the sensor body remains stationary, see Figure 6.82. The position magnet moves along the measuring area without contacting the sensing element.

Ferromagnetic materials such as iron and nickel display the property called *magnetostriction*. Application of a magnetic field to these materials causes a strain in the crystal structure, resulting in a change in the size and shape of the material. A material exhibiting positive magnetostriction will expand when magnetized. Conversely, with negative magnetostriction, the material contracts when magnetized [4].

The ferromagnetic materials used in magnetostrictive displacement sensors are transition metals, such as iron, nickel, and cobalt. In these metals, the $3d$ electron shell is not completely filled, which allows the formation of a magnetic moment (i.e., the shells closer to the nucleus are complete, and they do not contribute to the magnetic moment). As electron spins are rotated by a magnetic field, coupling between the electron spin and the electron orbit causes electron energies to change. The crystal strains so that electrons at the surface can relax to states of lower energy [5].

This physical response of a ferromagnetic material is due to the presence of magnetic moments, and can be understood by considering the material as a collection of tiny permanent magnets, called *domains*. Each domain consists of many atoms. When a material is not magnetized, the domains are randomly arranged. However, when the material is magnetized, the domains are oriented with their axes approximately parallel to each other. Interaction of an external magnetic field with the domains causes the magnetostrictive effect. See Figure 6.83. This effect can be optimized by controlling the ordering of domains through alloy selection, thermal annealing, cold working, and magnetic field strength.

While application of a magnetic field causes the physical strain, as described above, the reverse is also true: exerting stress causes the magnetic properties (permeability, susceptibility) to change. This is called the *Villari effect*.

In magnetostrictive sensors, uniform distortions of length, as shown in Figure 6.83, offer limited usefulness. Usually, the magnetization is rotated with a small field to induce a local distortion, using the

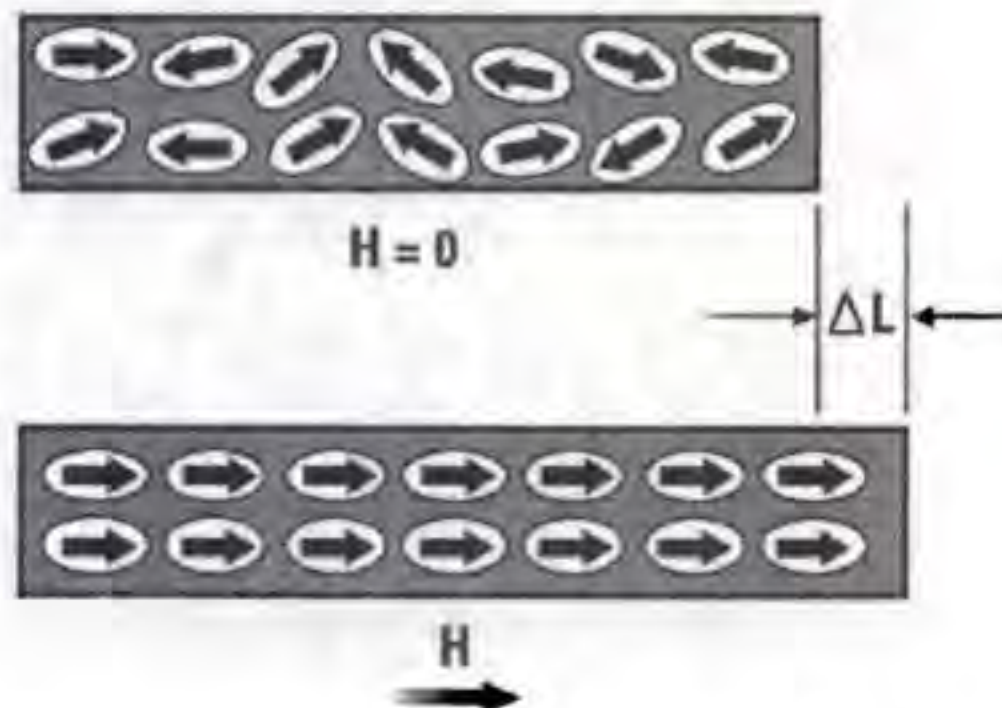


FIGURE 6.83 Magnetic domains: alignment with magnetic field, " H ", causes dimensional changes.

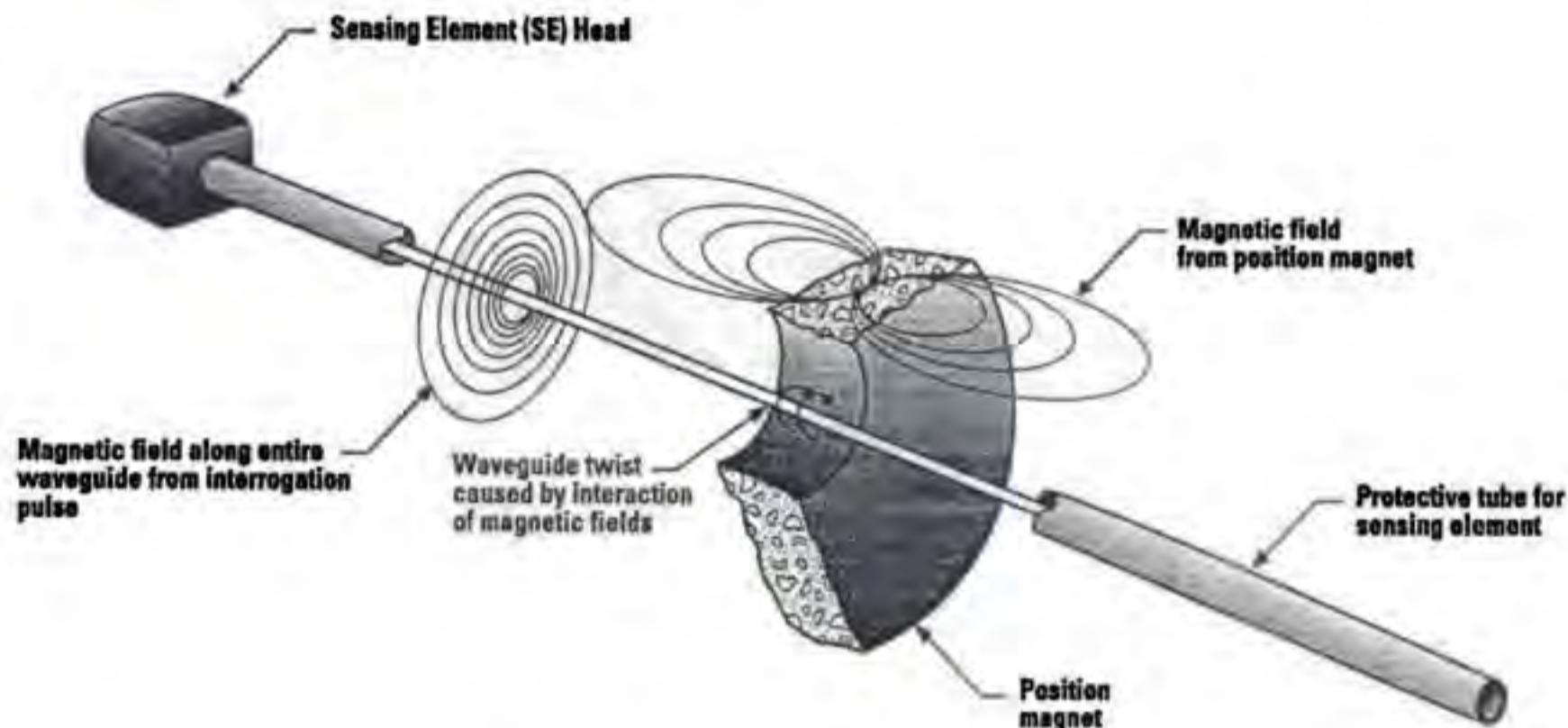


FIGURE 6.84 Operation of magnetostrictive position sensor.

Wiedemann effect. This is a mechanical torsion that occurs at a point along a magnetostrictive wire when an electric current is passed through the wire while it is subjected to an axial magnetic field. The torsion occurs at the location of the axial magnetic field, which is usually provided by a small permanent magnet called the position magnet.

In a displacement sensor, a ferromagnetic wire or tube called the waveguide is used as the sensing element, see Figure 6.84. The sensor measures the distance between the position magnet and the pickup. To start a measurement, a current pulse I (called the interrogation pulse), is applied to the waveguide. This causes a magnetic field to instantly surround it along its full length.

In a magnetostrictive position sensor, the current is a pulse of approximately 1 to 2 μs duration. A torsional mechanical wave is launched at the location of the position magnet due to the Wiedemann effect. Portions of this wave travel both toward and away from the pickup. The wave traveling along the waveguide toward the pickup is detected when it arrives at the pickup. The time measurement between application of the current pulse (launching of the torsion wave at the position magnet) until its detection by the pickup represents the location of the position magnet. The speed of the wave is typically about 3000 m s^{-1} . The portion of the wave traveling away from the pickup could act as an interfering signal after it is reflected from the waveguide tip. So instead, it is damped by a damping element when it reaches the end of the waveguide opposite the pickup. Damping is usually accomplished by attaching one of various configurations of elastomeric materials to the end of the waveguide. The end of the waveguide within the damping element is unusable for position determination, and therefore called the "dead zone."

The time measurement can be buffered and used directly as the sensor output, or it can be conditioned inside the sensor to provide various output types, including analog voltage or current, pulse width modulation, CANbus, SSI, HART, Profibus, etc. Magnetostrictive position sensors can be made as short as 1 cm long or up to more than 30 m long. Resolution of those produced by MTS Systems Corp. is as fine as $1 \mu\text{m}$. Temperature coefficients of 2 to 5 $\text{ppm } ^\circ\text{C}^{-1}$ can be achieved. The sensors are inherently stable, since the measurement relies on the physical properties of the waveguide material. Longer sensors become very cost effective because the same electronics package can drive sensors of varying length; only the waveguide and its packaging are increased in length to make the sensor longer.

The magnetostrictive wire can be straight for a linear sensor, or shaped to provide curved or rotary measurements. Curved sensors are often used to measure angular or nonlinear motion in industrial applications, although rotary magnetostrictive sensors are not yet very popular.

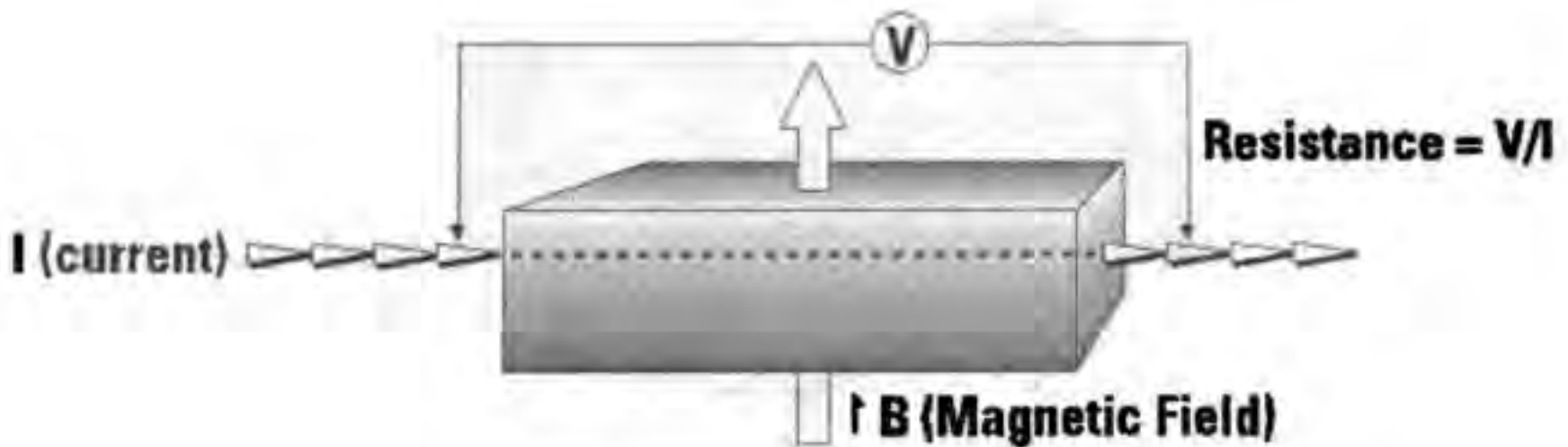


FIGURE 6.85 Magnetoresistance.

Magnetoresistive Sensors

In most magnetic materials, electrical resistance decreases when a magnetic field is applied and the magnetization is *perpendicular* to the current flow (a current will be flowing any time electrical resistance is measured) (see Figure 6.85). The resistance decreases as the magnetic flux density increases, until the material reaches magnetic saturation. The rate of resistance decrease is less as the material nears saturation. The amount of resistance change is on the order of about 1% at room temperature (0.3% in iron, 2% in nickel). When the magnetic field is *parallel* to the current, the resistance increases with increasing magnetic field strength. Sensitivity is greatest when the magnetic field is perpendicular to the current flow. These are properties of the phenomenon called *magnetoresistance* (MR). The MR effect is due to the combination of two component parts. These are: a reduction in forward carrier velocity as a result of the carriers being forced to move sideways as well as forward, and a reduction in the effective cross-sectional area of the conductor as a result of the carriers being crowded to one side [6].

When a position magnet is brought close to a single MR sensing element, the resistance change is maximum as the magnet passes over the approximate center of the element and then reduces until the magnet is past the element. The resistance changes according to:

$$\text{Resistivity} = \text{Voltage} / (\text{carrier density} \times \text{carrier velocity}) \quad (6.111)$$

By using multiple MR elements arranged along a line, a longer displacement measuring device can be fashioned. The signals from the string of sensors are decoded to find which elements are being affected by the magnet. Then the individual readings are used to determine the magnet position more precisely. Relatively high-performance sensors can be manufactured. Temperature sensitivity of the MR elements needs to be compensated, and longer sensors contain many individual sensing elements. Because of this, longer sensors become more difficult to manufacture, and are expensive.

Anisotropic MR materials are capable of resistance changes in the range of 1% or 2%. The MR of a conductor body can be increased by making it a composite of two or more layers of materials having different levels of magnetoresistance. Multilayered structures of exotic materials (sometimes more than 10 layers) have enabled development of materials that exhibit much greater magnetoresistive effect, and saturate at larger applied fields. This has been named Giant MagnetoResistance (GMR). Some commercial sensors based on GMR are currently available. The GMR elements can be arranged in a four-element bridge connection for greater sensitivity. In this arrangement, two of the elements are shielded from the applied magnetic field. The other two elements are sensitive to the applied field. Sensitivity can also be increased by incorporating flux concentrators on the sensitive elements. In a bridge connection, the output voltage can vary by more than 5% of the supply voltage [7]. Rotary sensors can be constructed by attaching a pole piece to a rotating shaft. One or more permanent magnets and the pole piece are arranged to cause the magnetic field around the MR element to change with angular displacement.

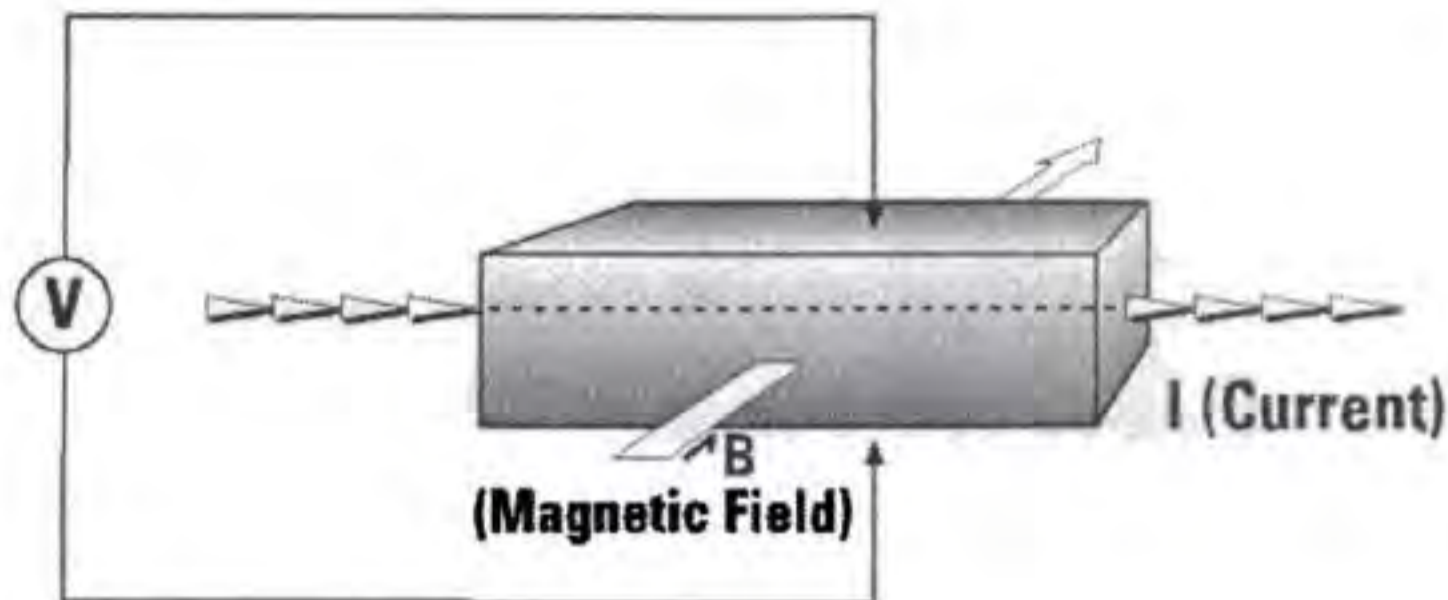


FIGURE 6.86 Hall effect.

Further research is being conducted on MR materials to improve the sensitivity by lowering the strength of magnetic field needed, and increasing the amount of resistance change. The next higher level of MR performance is being called Colossal MagnetoResistance (CMR). CMR is not yet practical for industrial sensors because of severe limitations on the operating temperature range.

Although MR, GMR, and CMR are limited for use in displacement sensors at this time by cost, temperature, and fabrication constraints, much research is in progress. Maybe Humongous Magneto-Resistance (HMR) is next?

Hall Effect Sensors

The Hall effect is a property exhibited in a conductor affected by a magnetic field. A voltage potential V_H , called the Hall voltage, appears across the conductor when a magnetic field is applied at right angles to the current flow. Its direction is perpendicular to both the magnetic field and current. The magnitude of the Hall voltage is proportional to both the magnetic flux density and the current. The magnetic field causes a gradient of carrier concentration across the conductor. The larger number of carriers on one side of the conductor, compared to the other side, causes the voltage potential V_H . A pictorial representation is shown in Figure 6.86. The amplitude of the voltage varies with the current and magnetic field according to: [8]

$$V_H = K_H \beta I / z \quad (6.112)$$

where V_H = Hall voltage
 K_H = Hall constant
 β = magnetic flux density
 I = current flowing through the conductor
 z = thickness of the conductor

Sensors utilizing the Hall effect typically are constructed of semiconductor material, giving the advantage of allowing conditioning electronics to be deposited right on the same material. Either p - or n -type semiconductor material can be used, with the associated polarity of current flow. The greatest output is achieved with a large Hall constant, which requires high carrier mobility. Low resistivity will limit thermal noise voltage, for a more useful signal-to-noise ratio (SNR). These conditions are optimized using an n -type semiconductor [6].

A displacement sensor can be made with a Hall sensing element and a movable magnet, with an output proportional to the distance between the two. Two magnets can be arranged with one Hall sensor as in Figure 6.87 to yield a near-zero field intensity when the sensor is equidistant between the magnets. These

TABLE 6.26 Sensors and Manufacturers

Technology	Manufacturers	Description	Price
Magnetostrictive	MTS Systems Corp. Cary, NC & Germany	Lengths to 20 m; 2 μ m resolution; CAN, SSI, Profibus, HART	\$150–\$3000
	Balluff Germany	Lengths to 3.5 m, 20 μ m resolution no standard interfaces in head	\$400–\$2300
Magnetoresistive	Nonvolatile Electronics Eden Prairie, MN	GMR sensors with flux concentrator and shield	\$2.50–\$6.00
	Midori America Fullerton, CA	Rotary MR sensors Linear MR up to 30 mm	\$64–\$500 \$67–\$200
	Hall Effect Optec Technology, Inc. Carrollton, TX	Linear position	\$5–\$50
Magnetic encoder	Spectec Emigrant, MT	Standard and custom sensors	Approx. \$90
	Heidenhain Schaumburg, IL	Rotary and linear encoders	\$300–\$2000
	Sony Precision Technology America Orange, CA	Rotary and linear encoders	\$100–\$2000

7. Nonvolatile Electronics Inc. NVSB series datasheet. March 1996.
8. J. R. Carstens, *Electrical Sensors and Transducers*, Englewood Cliffs, NJ: Regents/Prentice-Hall, 1992, p. 125.
9. H. Norton, *Handbook of Transducers*, Englewood Cliffs, NJ: Prentice-Hall, 1989, 106-112.

Further Information

- B. D. Cullity, *Introduction to Magnetic Materials*, Reading, MA: Addison-Wesley, 1972.
- D. Craik, *Magnetism Principles and Applications*, New York: John Wiley & Sons, 1995.
- P. Lorrain and D. Corson, *Electromagnetic Fields and Waves*, San Francisco: W.H. Freeman, 1962.
- R. Boll, *Soft Magnetic Materials*, London: Heyden & Son, 1977.
- H. Olson, *Dynamical Analogies*, New York: D. Van Nostrand, 1943.
- D. Askeland, *The Science and Engineering of Materials*, Boston: PWS-Kent Publishing, 1989.
- R. Rose, L. Shepard, and J. Wulff, *The Structure and Properties of Materials*, New York: John Wiley & Sons, 1966.
- J. Shackelford, *Introduction to Materials Science for Engineers*, New York: Macmillan, 1985.
- D. Jiles, *Introduction to Magnetism and Magnetic Materials*, London: Chapman and Hall, 1991.
- F. Mazda, *Electronics Engineer's Reference Book, 6th ed.*, London: Butterworth, 1989.
- E. Herceg, *Handbook of Measurement and Control*, New Jersey: Schaevitz Engineering, 1976.

6.10 Synchro/Resolver Displacement Sensors

Robert M. Hyatt, Jr. and David Dayton

Most electromagnetic position transducers are based on transformer technology. Transformers work by exciting the primary winding with a continuously changing voltage and inducing a voltage in the secondary winding by subjecting it to the changing magnetic field set up by the primary. They are ac-only devices, which make all electromagnetically coupled position sensors ac transformer coupled. They are inductive by nature, consisting of wound coils. By varying the amount of coupling from the primary (excited) winding of a transformer to the secondary (coupled) winding with respect to either linear or rotary displacement, an analog signal can be generated that represents the displacement. This coupling variation is accomplished by moving either one of the windings or a core element that provides a flux path between the two windings.

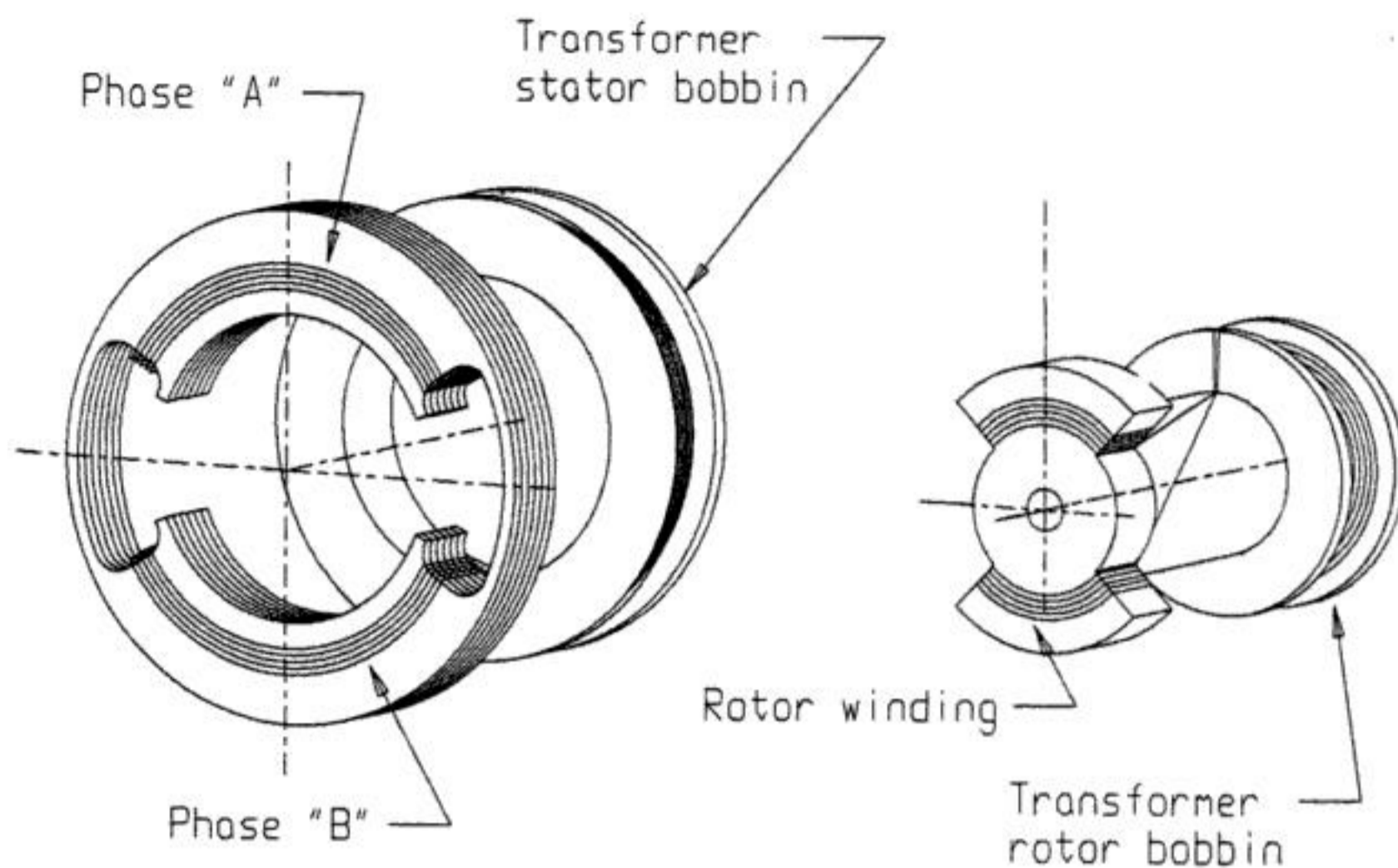


FIGURE 6.89 The induction potentiometer has windings on the rotor and the stator.

One of the simplest forms of electromagnetic position transducers is the LVDT, which is described in Section 6.2 on inductive sensors. If the displacement of the core in an LVDT-type unit is changed from linear to rotary, the device becomes an RVDT (rotary variable differential transformer).

Induction Potentiometers

The component designer can “boost” the output, increase accuracy, and achieve a slightly greater angular range if windings are placed on the rotor as shown in Figure 6.89. The disadvantages to this method are (1) additional windings, (2) more physical space required, (3) greater variation over temperature, and (4) greater phase shift due to the additional windings.

The advantage of the induction pot design is greater sensitivity (more volts per degree), resulting in better signal to noise and higher accuracy in most cases.

Resolvers

If the two-slot lamination in the stator stack shown in Figure 6.89 is replaced by a multislot lamination (see Figure 6.90), and two sets of windings are designed in concentric coil sets and distributed in each quadrant of the laminated stack; a close approximation to a sine wave can be generated on one of the secondary windings and a close approximation to a cosine wave can be generated on the other set of windings. Rotary transformers of this design are called *resolvers*. Using a multislot rotor lamination and distributing the windings in the rotor, the sine-cosine waveforms can be improved even further.

A resolver effectively amplitude modulates the ac excitation signal placed on the rotor windings in proportion to the sine and the cosine of the angle of mechanical rotation. This sine-cosine electrical output information measured across the stator windings may be used for position and velocity data. In this manner, the resolver is an analog trigonometric function generator. Most resolvers have two primary windings that are located at right angles to each other in the stator, and two secondary windings also at right angles to each other, located on the rotor (see Figure 6.91).

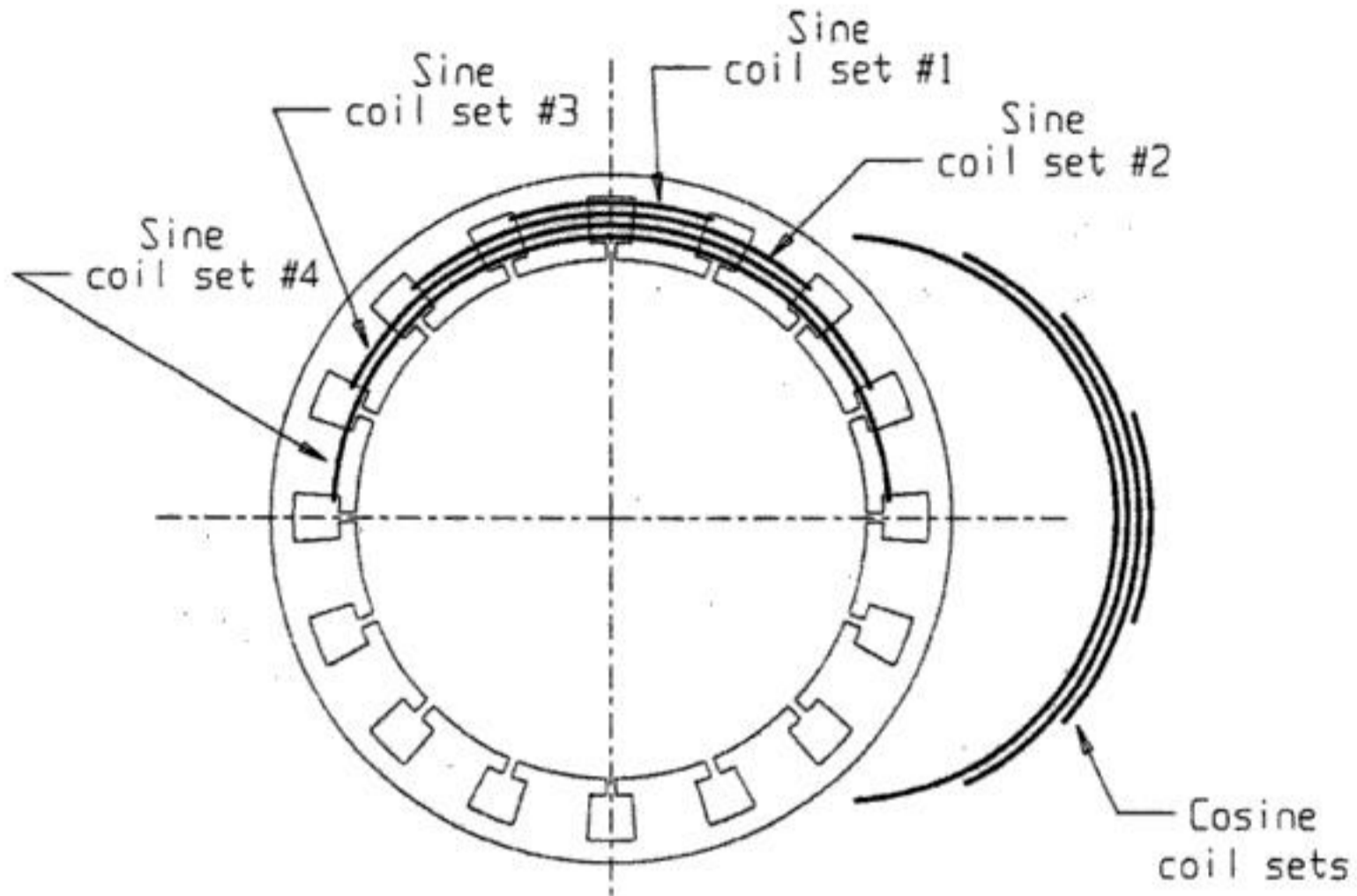


FIGURE 6.90 The resolver stator has distributed coil windings on a 16-slot lamination to generate a sine wave.

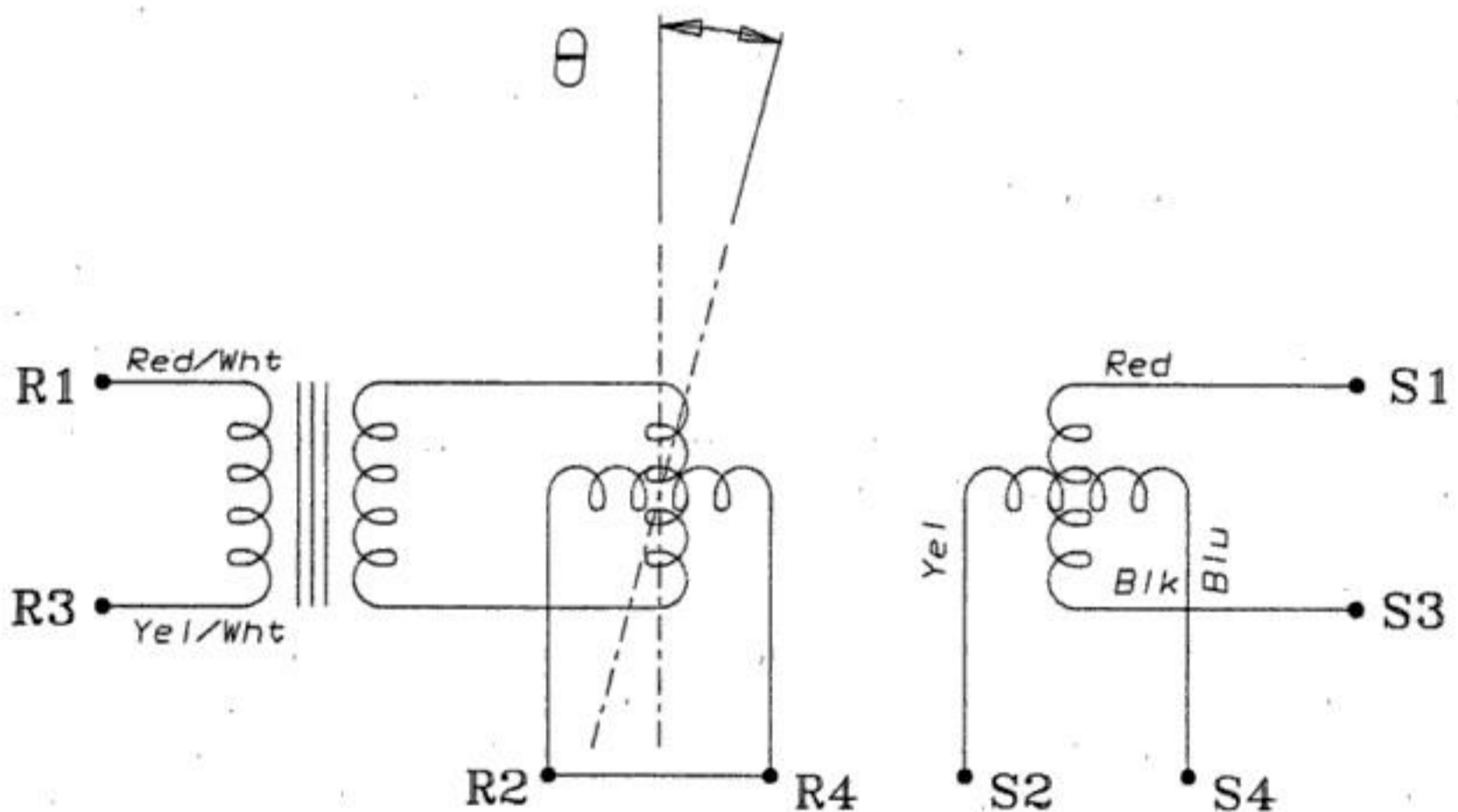


FIGURE 6.91 A brushless resolver modulates the ac excitation on the rotor by the rotation angle.

If the rotor winding (R1-R3) is excited with the rated input voltage (see Figure 6.92), the amplitude of the output winding of the stator (S2-S4) will be proportional to the sine of the rotor angle θ , and the amplitude of the output of the second stator winding (S1-S3) will be proportional to cosine θ . (See Figure 6.93.) This is commonly called the "control transmitter" mode and is used with most "state-of-the-art" resolver to digital converters.

In the control transmitter mode, electrical zero may be defined as the position of the rotor with respect to the stator at which there is minimum voltage across S2-S4 when the rotor winding R1-R3 is excited

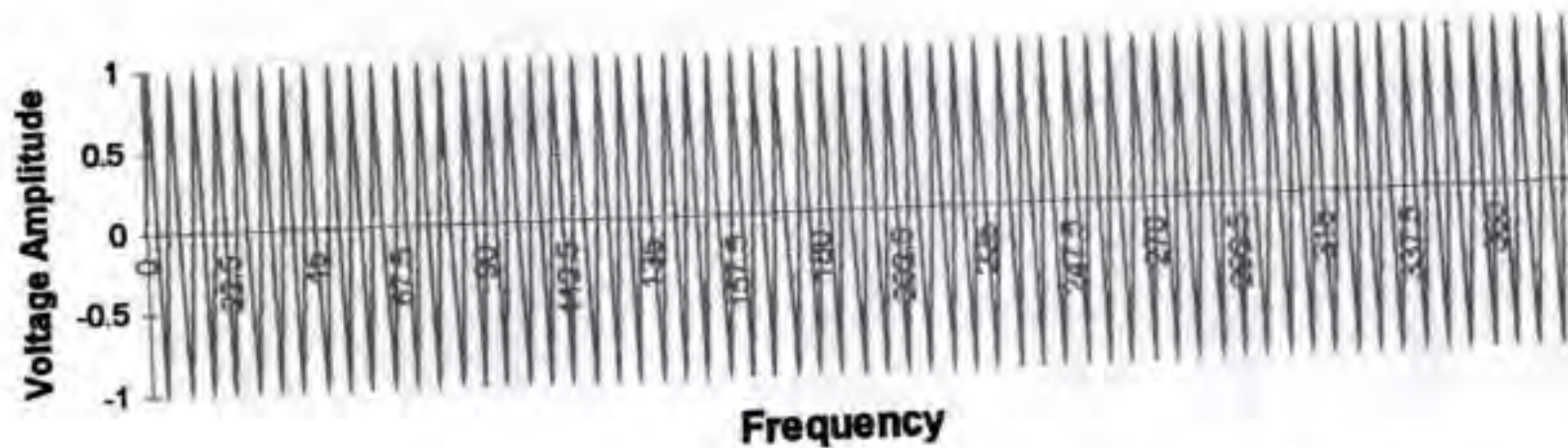


FIGURE 6.92 The resolver rotor winding is excited with the rated input voltage.

with rated voltage. Nulls will occur across S2-S4 at the 0° and 180° positions, and will occur across S1-S3 at the 90° and 270° positions.

If the stator winding S1-S3 is excited with the rated input voltage and stator winding S2-S4 is excited with the rated input voltage electrically shifted by exactly 90° , then the output sensed on the rotor winding R1-R3 does not vary with rotor rotation in amplitude or frequency from the input reference signal. It is the sum of both inputs. It does, however, vary in time phase from the rated input by the angle of the shaft from a referenced "zero" point (see Figure 6.94). This is a "phase analog" output and the device is termed a "control transformer." By measuring the time difference between the zero crossing of the reference voltage waveform and the output voltage waveform, the phase angle (which is the physical angular displacement of the output shaft) can be calculated.

Because the resolver is an analog device and the outputs are continuous through 360° , the theoretical resolution of a resolver is infinite. There are, however, ambiguities in output voltages caused by inherent variations in the transformation of the voltage from primary to secondary through 360° of rotation. These ambiguities result in inaccuracy when determining the true angular position. The types of error signals that are found in resolvers are shown in Figure 6.95.

As a rule, the larger the diameter of the stator laminations, the better the accuracy and the higher the absolute resolution of the device. This is a function of the number of magnetic poles that can be fit into the device, which is a direct function of the number of slots in the stator and rotor laminations. With multispeed units (see the section on multispeeds below), the resolution increases as a multiple of the speeds. For most angular excursions, the multispeed resolver can exceed the positioning accuracy capability of any other component in its size, weight, and price range.

Operating Parameters and Specifications for Resolvers

There are seven functional parameters that define the operation of a resolver in the analog mode. These are (1) accuracy, (2) operating voltage amplitude, (3) operating frequency, (4) phase shift of the output voltage from the referenced input voltage, (5) maximum allowable current draw, (6) the transformation ratio of output voltage over input voltage, and (7) the null voltage. Although impedance controls the functional parameters, it is transparent to the user. The lamination and coil design are usually developed to minimize null voltage and input current, and the impedance is a direct fallout of the inherent design of the resolver. The following procedure can be used to measure the seven values for most resolvers.

Equipment Needed for Testing Resolvers

A mechanical index stand that can position the shaft of the resolver to an angular accuracy that is an order of magnitude greater than the specified accuracy of the resolver.

An ac signal generator capable of up to 24 Vrms at 10 kHz.

A phase angle voltmeter (PAV) capable of measuring "in phase" and "quadrature" voltage components for determining the transformation ratio and the null voltage as well as the phase angle between the output voltage and the reference input voltage.

A $1\ \Omega$ resistor used to measure input current with the PAV.

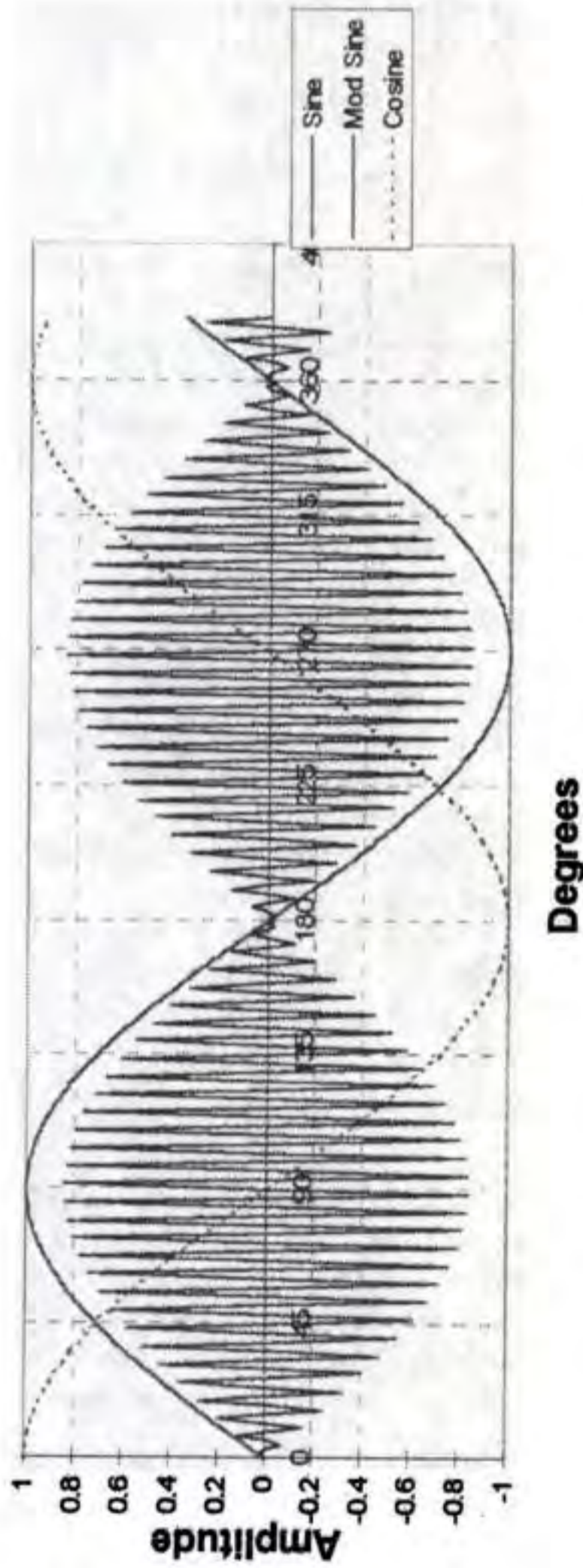


FIGURE 6.93 The single-speed resolver stator output is the sine or cosine of the angle.

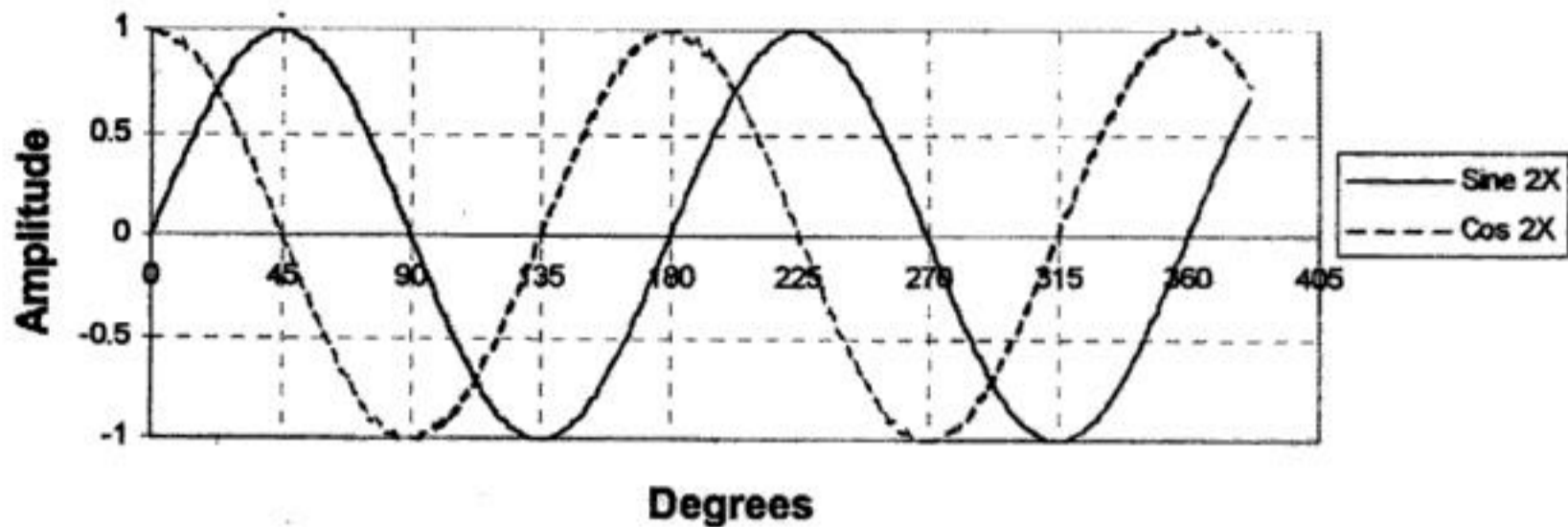


FIGURE 6.96 A two-speed resolver yields two electrical cycles for one rotation.

Multispeed Units

The relationship for multispeeds is that the speed ($2\times$, $3\times$, etc.) designates how many full sinusoidal cycles the resolver output electrically completes in 360° of mechanical rotation. The $2\times$ electrical output is such that the full sinusoidal cycle for a $2\times$ resolver occurs in 180° instead of 360° . A $2\times$ resolver output is shown in Figure 6.96. A full $3\times$ cycle is completed in 120° . The number of speeds selected for use is a function of the system requirements. Increasing the number of magnetic poles in the rotor and stator creates multispeed units. Each speed has several winding and slot combinations. The optimum combination is selected by the resolver designer based on system demands.

Applications

Resolvers are often used in conjunction with motors, and because of their inherent similarity of design (copper windings on iron lamination stacks), their environmental resistance is quite similar. They are ideal to design into industrial applications where dust and airborne liquids can obscure optical encoder signals. NC machines, coil winders, presses, and positioning tables are uses where resolvers excel. The resolver's inherent resistance to shock and vibration makes it uniquely suited to moving platforms, and their reliability under these conditions lends a welcome hand to the designers of robots, gantries, and automotive transfer lines.

Heat sensitivity is always a problem for motion control systems designers. Resolvers used for sensing the position of valves in high-temperature applications such as aircraft engines, petrochemical refining, and chemical processing have continually proven their reliability.

Moving devices to precise positions with smooth and accurate control can be a real challenge in the electromagnetic noise environment of the modern industrial facility. Emitted and conducted EMI from adjacent equipment, and input voltage variations with unwanted current spikes on input power lines can rob digital systems of their signal integrity. The analog resolver continues to function without information loss or signal interruption. Digitizing the signal can be done at a remote interface under more controlled conditions than on the factory floor. Only robust materials can perform well in harsh environments.

Synchros

As long ago as World War II, synchros were used in analog positioning systems to provide data and to control the physical position of mechanical devices such as radar antennae, indicator needles on instrumentation, and fire control mechanisms in military equipment. The term "synchro" defines an electromagnetic position transducer that has a set of three phase output windings that are electrically and mechanically spaced by 120° instead of the 90° spacing found in a resolver. In the rotor primary mode, the synchro is excited by a single-phase ac signal on the rotor. As the rotor moves 360° , the three amplitude modulated sine waves on the three phases of the output have a discrete set of amplitudes for each angular position. By interpreting these amplitudes, a table can be established to decode the exact rotary position.

In most applications, resolvers have replaced synchros because of the sophistication of the resolver-to-digital converters that are commercially available. Working with a sine and cosine is simpler and requires less conversion and decoding than using three 120° spaced signals. If conversion of a synchro output is desired in resolver format, a device known as a Scott “T” transformer can be used for conversion. In most synchro-to-digital processors, the first step is to convert the signal to a resolver format with a Scott “T” device.

A Modular Solution

The brushless resolver is a self-contained feedback device that, unlike optical encoders, provides an analog signal with infinite resolution. Not only can the output signal be converted to precise digital position information, but it also provides an accurate velocity signal, thus eliminating the need for using a separate tachometer. Reliability is enhanced using the same resolver for speed feedback and commutation. Piece part count can be reduced and the complexity of using Hall-effect devices for timing signals for commutation can be eliminated.

A modular approach allows the designer to easily select a single or multispeed resolver and appropriate electronics that will meet almost any desired level of resolution and accuracy. The resolvers are designed in the most commonly used frame sizes: 8, 11, 15, and 21. Housed models feature high-quality, motor-grade ball bearings. Heavy-duty industrial grade units are enclosed in rugged black painted aluminum housings with either flange, face, or servo-type mounting, and utilize MS-style connectors.

The Sensible Design Alternative for Shaft Angle Encoding

The requirement for velocity and position feedback plays an important role in today’s motion control systems. With the development of low-cost monolithic resolver-to-digital converters, a resolver-based system provides design engineers with the building blocks to handle a wide variety of applications. A resolver’s small size, rugged design, and the ability to provide a very high degree of accuracy under severe conditions, make this the ideal transducer for absolute position sensing. These devices are also well suited for use in extremely hostile environments such as continuous mechanical shock and vibration, humidity, oil mist, coolants, and solvents. Absolute position sensing vs. incremental position sensing is a necessity when working in an environment where there is the possibility of power loss. Whenever power is supplied to an absolute system, it is capable of reading its position immediately; this eliminates the need for a “go home” or reference starting point.

Resolver-to-Digital Converters

A monolithic resolver-to-digital converter requires only six external passive components to set the bandwidth and maximum tracking rate. The bandwidth controls how quickly the converter will react to a large change in position on the resolver output. The converter can also be programmed to provide either 10, 12, 14, or 16 bits of parallel data. A resolver-based system can provide high dynamic capability and high resolution for today’s motion control systems where precision feedback for both position and velocity is required.

Closed Loop Feedback

In a typical closed loop servo model as in Figure 6.97, the position sensor plays an important role by constantly updating the position and velocity information. Selection of a machine control strategy will often be based on performance, total application cost, and technology comfort. The accuracy of the system is determined by the smallest resolution of the position-sensing device. A resolver-to-digital converter in the 16-bit mode has 2^{16} (65,536) counts per revolution, which is equivalent to a resolution of 20 arc seconds. The overall accuracy of the resolver-to-digital converter is ± 2.3 arc minutes. An accuracy specification defines the maximum error in achieving a desired position. System accuracy must be smaller than the tolerance on the desired measurement. An important feature of the resolver-to-digital converter

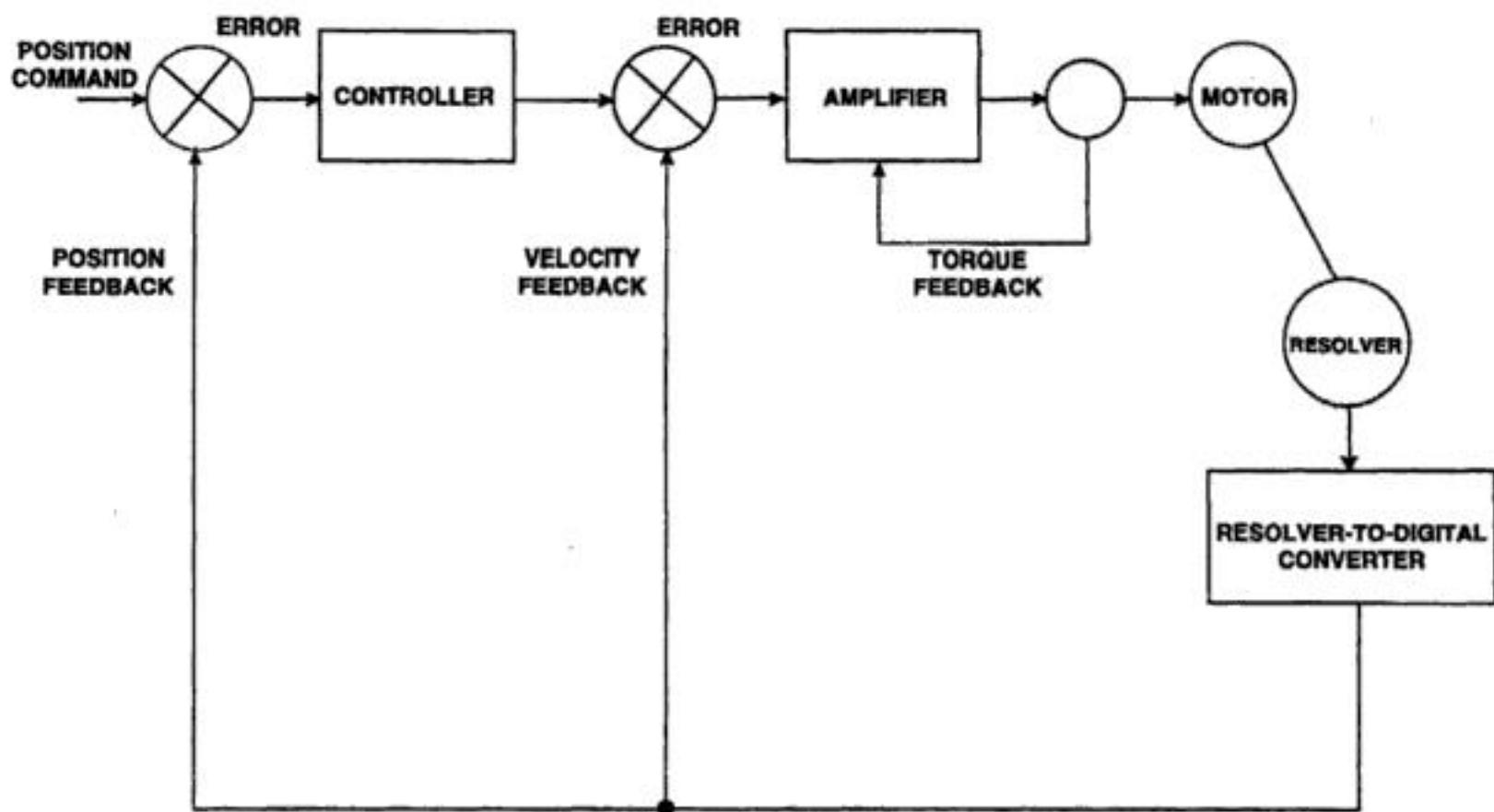


FIGURE 6.97 A closed-loop servo model uses a resolver-to-digital converter.

is repeatability. With a repeatability specification of ± 1 LSB (least significant bit) in the 16-bit mode, this provides an accurate measurement when determining position from point to point. For example, moving from point A to point B and back to point A, the converter in the 16-bit mode will be accurate within 20 arc seconds of the original position. The error curve of a resolver-to-digital converter is repeatable within ± 1 LSB. The combination of high precision resolvers (± 20 arc seconds) with a resolver-to-digital converter provides accurate absolute position information for precision feedback for motion control.

Type II Servo Loop

The motor speed is monitored using the velocity output signal generated by the resolver-to-digital converter. This signal is a dc voltage proportional to the rate of speed, positive for increasing angles and negative for decreasing angles, with a typical linearity specification of 0.25% and a typical reversal error of 0.75%. The error processing is performed using the industry standard technique for type II tracking, resolver-to-digital converters (see Figure 6.98).

The dc error is integrated, yielding a velocity voltage that drives a voltage-controlled oscillator (VCO). This VCO is an incremental integrator (constant voltage input to position rate output) that together with the velocity integrator, forms a type II critically damped, servo feedback loop. This information allows the motor to maintain constant speeds under varying loads when it is interfaced with a programmable logic controller (PLC). The PLC-based architecture is used for I/O intensive control applications. The PLC provides a low-cost option for those developers familiar with its ladder logic programming language. Integration of the motion, I/O, operator's interface, and communication are usually supported through additional cards that are plugged into the backplane.

Applications

Specific applications require unique profiles to control the speed and acceleration of the motor to perform the task at hand. By reducing the accelerations and decelerations that occur during each operation, it is possible to lower the cost and use more efficient motors. Industrial applications include the following:

- Ballscrew positioning
- Motor commutation
- Robotics positioning

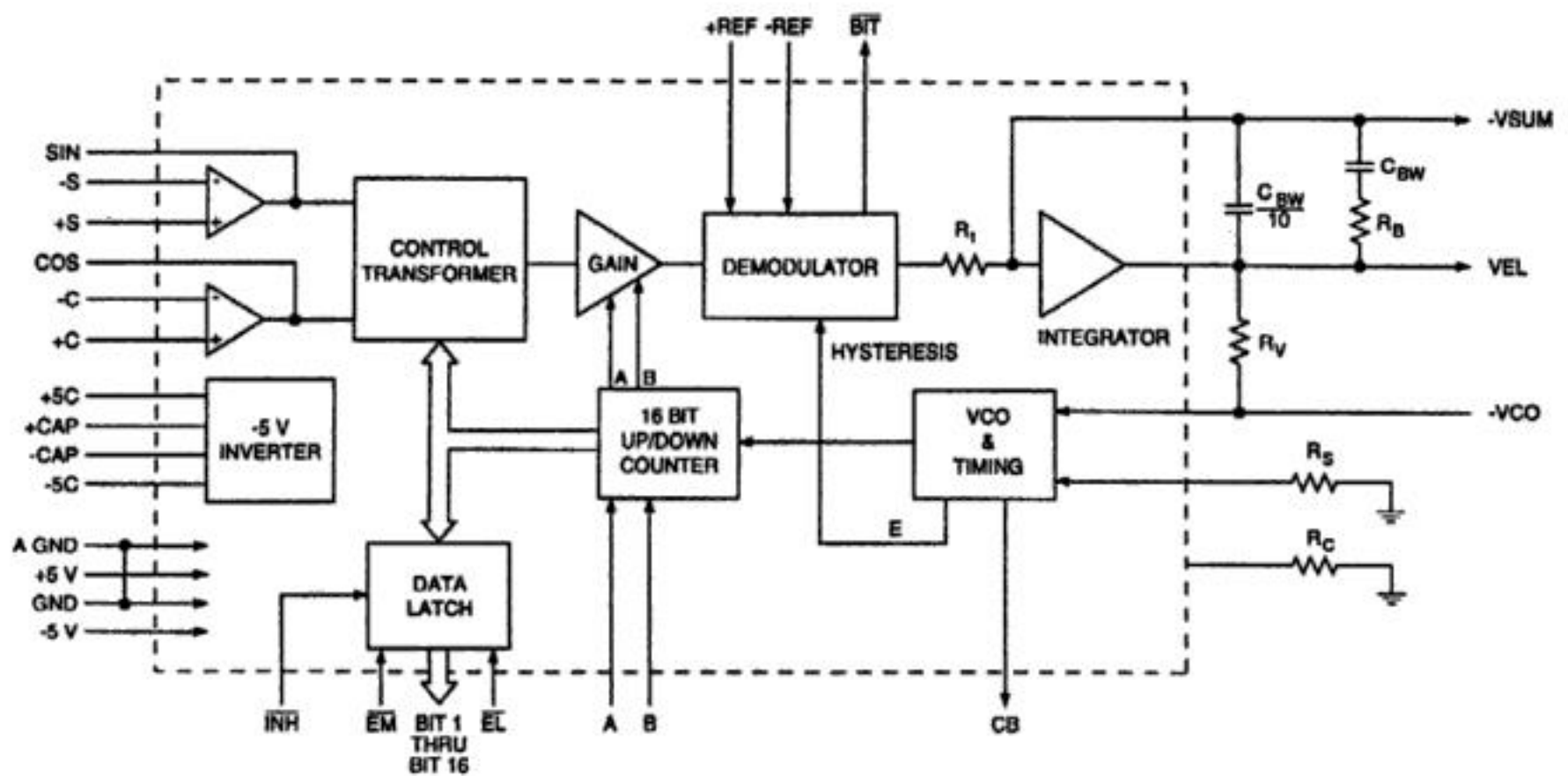


FIGURE 6.98 Error processing uses type II tracking resolver-to-digital converters.

Machine vision systems
 X-Y tables
 Component insertion
 Remote video controls
 Web guides
 Pick and place machines

Resolver-to-Digital Conversion

For a resolver-to-digital converter, the resolver information is presented to a solid-state resolver conditioner that reduces the signal amplitude to 2 V rms sine and cosine; the amplitude of one being proportional to the sine of θ (the angle to be digitized), and the amplitude of the other being proportional to the cosine of θ . (The amplitudes referred to are, of course, the carrier amplitudes at the reference frequency, i.e., the cosine wave is actually $\cos \theta \cos \omega t$; but the carrier term, $\cos \omega t$, will be ignored in this discussion because it will be removed in the demodulator, and at any rate contains no data). A quadrant selector circuit in the control transformer enables selection of the quadrant in which θ lies, and automatically sets the polarities of the sine θ and $\cos \theta$ appropriately, for computational significance. The $\sin \theta$, $\cos \theta$ outputs of the quadrant selector are then fed to the sine and cosine multipliers, also contained in the control transformer. These multipliers are digitally programmed resistive networks. The transfer function of each of these networks is determined by a digital input (which switches in proportioned resistors), so that the instantaneous value of the output is the product of the instantaneous value of the analog input and the sine (or cosine) of the digitally encoded angle. If the instantaneous value of the analog input of the sine multiplier is $\cos \theta$, and the digitally encoded "word" presented to the sine multiplier is ϕ , then the output code is $\cos \theta \sin \phi$. Thus, the two outputs of the multipliers are

From the sine multiplier: $\cos \theta \sin \phi$

From the cosine multiplier: $\sin \theta \cos \phi$

These outputs are fed to an operational subtractor, at the differencing junction shown, so that the input fed to the demodulator is

$$\sin \theta \cos \phi - \cos \theta \sin \phi = \sin (\theta - \phi) \quad (6.114)$$

The right-hand side of this trigonometric identity indicates that the differencing-junction output represents a carrier-frequency sine wave with an amplitude proportional to the sine of the difference between θ (the angle to be digitized) and ϕ (the angle stored in digital form in the up/down counter). This point is the ac error signal brought out as (e). The demodulator is also presented with the reference voltage, which has been isolated from the reference source, and appropriately scaled, by the reference conditioner. The output of the demodulator is then, an analog dc level, proportional to $\sin(\theta - \phi)$, in other words, to the sine of the "error" between the actual angular position of the resolver and the digitally encoded angle, ϕ , which is the output of the counter. This point dc error is sometimes brought out as (D) while an addition of a threshold detector will give a built-in-test (BIT) flag. When the ac error signal exceeds 100 LSBs, the BIT flag will indicate a tracking error. This angular error signal is then fed into the error processor and VCO. This circuit consists essentially of an analog integrator whose output (the time integral to the error) controls the frequency of a voltage-controlled oscillator (VCO). The VCO produces clock pulses that are counted by the up/down counter. The "sense" of the error (ϕ too high or ϕ too low) is determined by the polarity of (ϕ), and is used to generate a control counter signal (U), which determines whether the counter increments upward or downward. Finally, note that the up/down counter, like any counter, is functionally an incremental integrator; therefore, the tracking converter constitutes in itself a closed-loop servomechanism (continuously attempting to null the error to zero) with two integrators in series. This called a "Type II" servo loop, which has decided advantages over Type 1 or Type 0 loops. In order to appreciate the value of a Type II servo behavior of this tracking converter, consider first the shaft is not moving. Ignoring inaccuracies, drifts, and the inevitable quantizing error, the error should be zero ($\theta = \phi$), and the digital output represents the true shaft angle of the resolver. Now, start the resolver shaft moving, and allow it to accelerate uniformly, from $d\theta/dt = 0$ to $d\theta/dt = V$. During the acceleration, an error will develop because the converter cannot instantaneously respond to the change of angular velocity. However, since the VCO is controlled by an integrator, the output of which is the integral of the error, the greater the lag (between θ and ϕ), the faster the counter will be called on to catch up. When the velocity becomes constant at V , the VCO will have settled to a rate of counting that exactly corresponds to the rate of change in θ per unit time and instantaneously $\theta = \phi$. Therefore, $d\phi/dt$ will always track $d\theta/dt$ without a velocity or position error. the only error will be momentary (transient) error, during acceleration or deceleration. Furthermore, the information produced by the tracking converter is always "fresh," being continually updated, and always available at the output of the counter. Since $d\theta/dt$ tracks the input velocity it can be brought out as velocity, a dc voltage proportional to the rate of rotation, which is of sufficient linearity in modern converters to eliminate the need for a tachometer in many systems.

Bandwidth Optimization

When using a low-cost monolithic converter for position and velocity feedback, it is important to understand the dynamic response for a changing input. When considering what bandwidth to set the converter, several parameters must be taken into consideration. The ability to track step responses and accelerations will determine what bandwidth to select. The lower the bandwidth of the resolver-to-digital converter, the greater the noise immunity; high frequency noise will be rejected. The relationship between the maximum tracking rate and bandwidth determines the settling time for small and large steps. For a small step input, the bandwidth determines the converter settling time. When one has a large step, the maximum velocity slew rate and bandwidth together, determine the settling time.

Encoder Emulation

Today's resolver-to-digital converters also have the ability to emulate the output of an optical incremental encoder. By providing the outputs A, B, and Zero Index, the encoder can be replaced with a resolver and resolver-to-digital converter without changing the existing interface hardware.

on the basis of their complexity. Finally, a theoretical analysis of the cross-sensitivities of the four sensing schemes is presented and their performances are compared.

Strain measurements using optical fiber sensors in both embedded and surface-mounted configurations have been reported by researchers in the past [2]. Fiber optic sensors are small in size, immune to electromagnetic interference, and can be easily integrated with existing optical fiber communication links. Such sensors can typically be easily multiplexed, resulting in distributed networks that can be used for health monitoring of integrated, high-performance materials and structures. Optical fiber sensors for strain measurements should possess certain important characteristics. These sensors should either be insensitive to ambient fluctuations in temperature and pressure, or should have demodulation techniques that compensate for changes in the output signal due to the undesired perturbations. In the embedded configuration, the sensors for axial strain measurements should have minimum cross-sensitivity to other strain states. The sensor signal should itself be simple and easy to demodulate. Nonlinearities in the output demand expensive decoding procedures or require precalibrating the sensor. The sensor should ideally provide an absolute and real-time strain measurement in a form that can be easily processed. For environments where large strain magnitudes are expected, the sensor should have a large dynamic range while at the same time maintaining the desired sensitivity. A discussion of each of the four sensing schemes individually, along with their relative merits and demerits, follows.

Extrinsic Fabry–Perot Interferometric Sensor

The extrinsic Fabry–Perot interferometric (EFPI) sensor, proposed by Murphy et al., is one of the most popular fiber optic sensors used for applications in health monitoring of smart materials and structures [3]. As the name suggests, the EFPI is an interferometric sensor in which the detected intensity is modulated by the parameter under measurement. The simplest configuration of an EFPI is shown in Figure 6.99.

The EFPI system consists of a single-mode laser diode that illuminates a Fabry–Perot cavity through a fused biconical tapered coupler. The cavity is formed between an input single-mode fiber and a reflecting single-mode or multimode fiber. Since the cavity is external to the lead-in/lead-out fiber, the EFPI sensor is independent of transverse strain and small ambient temperature fluctuations. The input fiber and the reflecting fiber are aligned using a hollow-core silica fiber. For uncoated fiber ends, a 4% Fresnel reflection results at both ends. The first reflection, R_1 , called the reference reflection, is independent of the applied perturbation. The second reflection, R_2 , termed the sensing reflection, is dependent on the length of the cavity, d , which in turn is modulated by the applied perturbation. These two reflections interfere (provided $2d < L_c$, the laser diode's coherence length), and the intensity I at the detector varies as a function of the cavity length:

$$I = I_0 \cos\left(\frac{4\pi}{\lambda} d\right) \quad (6.117)$$

where, I_0 is the maximum value of the output intensity and λ is the laser diode center wavelength.

The typical EFPI transfer function curve is shown in Figure 6.100. Small perturbations that result in operation around the quiescent-point or Q-point of the sensor lead to a linear variation in output intensity. A fringe in the output signal is defined as the change in intensity from a maximum to a maximum or from a minimum to a minimum. Each fringe corresponds to a change in the cavity length by one half of the operating wavelength, λ . The change in the cavity length, Δd , is then employed to calculate the strain ϵ using the expression:

$$\epsilon = \frac{\Delta d}{L} \quad (6.118)$$

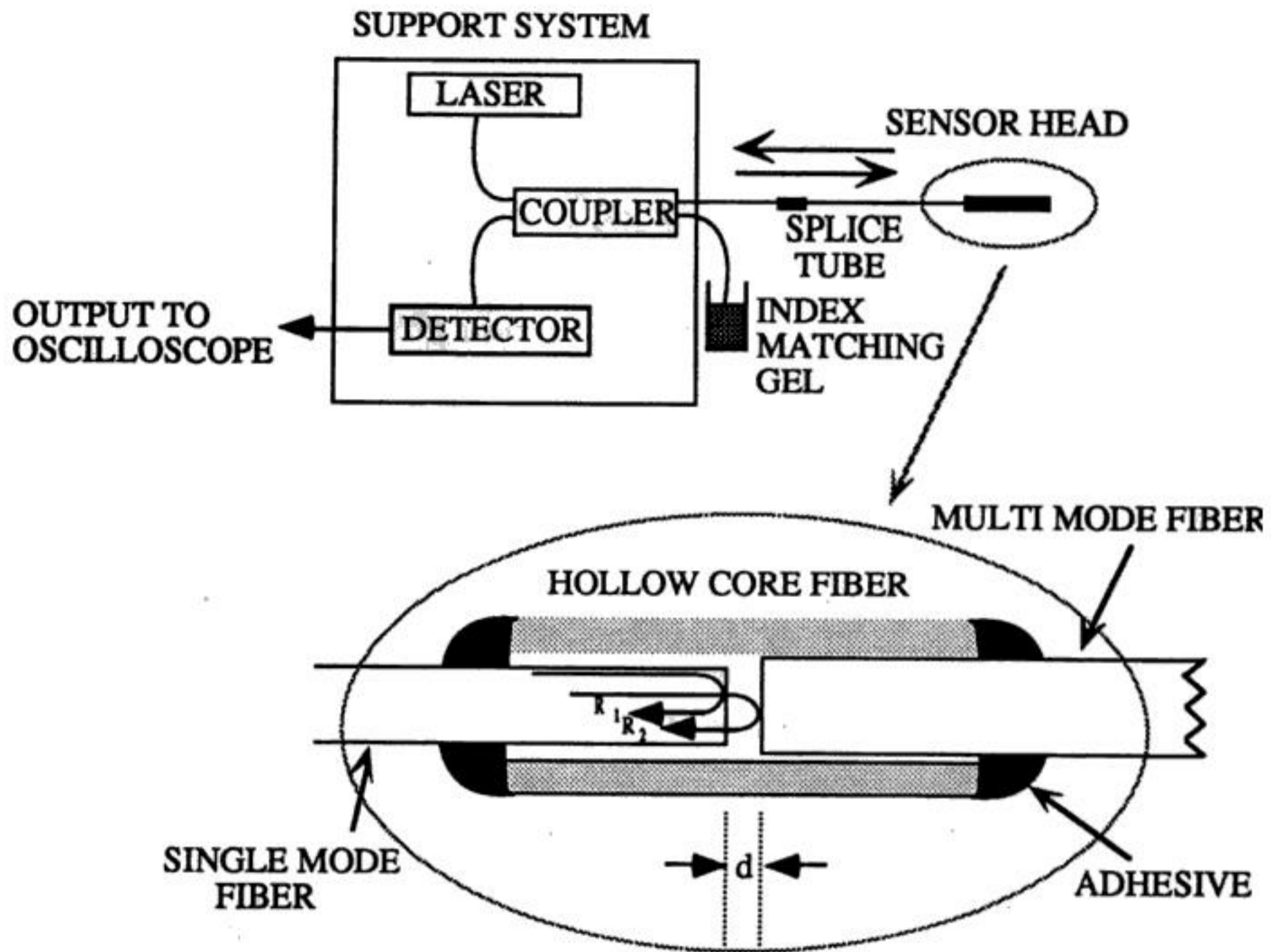


FIGURE 6.99 A simple configuration of an extrinsic Fabry-Perot interferometric (EFPI) sensing system.

where, L is defined as the gage length of the sensor and is typically the distance between two points where the input and reflecting fibers are bonded to the hollow-core fiber. Matching of the two reflection signal amplitudes allows good fringe visibility in the output signal.

The EFPI sensor has been extensively used for measuring fatigue loading on F-15 aircraft wings, detection of crack formation and propagation in civil structures, and cure and lifetime monitoring in concrete and composite specimens [2, 4]. The temperature insensitivity of this sensor makes it attractive for a large number of applications. The EFPI sensor is capable of measuring sub-Angstrom displacements with strain resolution better than 1 microstrain and a dynamic range greater than 10,000 $\mu\epsilon$. Although the change in output intensity of the EFPI is nonlinear corresponding to the magnitude of the parameter being measured, for small perturbations its operation can be limited to that around the Q-point of the transfer function curve. Moreover, the large bandwidth available with this sensor simplifies the measurement of highly cyclical strain. The EFPI sensor is capable of providing single-ended operation and is hence suitable for applications where access to the test area is limited. The sensor requires simple and inexpensive fabrication equipment and an assembly time of less than 10 min. Additionally, since the cavity is external to the fibers, transverse strain components that tend to influence intrinsic sensors through the Poisson's effect have negligible effect on the EFPI sensor output. The sensitivity to only axial strain and insensitivity to input polarization state have made the EFPI sensor the most preferred fiber optic sensor for embedded applications [1]. Thus, overall, the EFPI sensing system is very well suited to measurement of small magnitudes of cyclical strain.

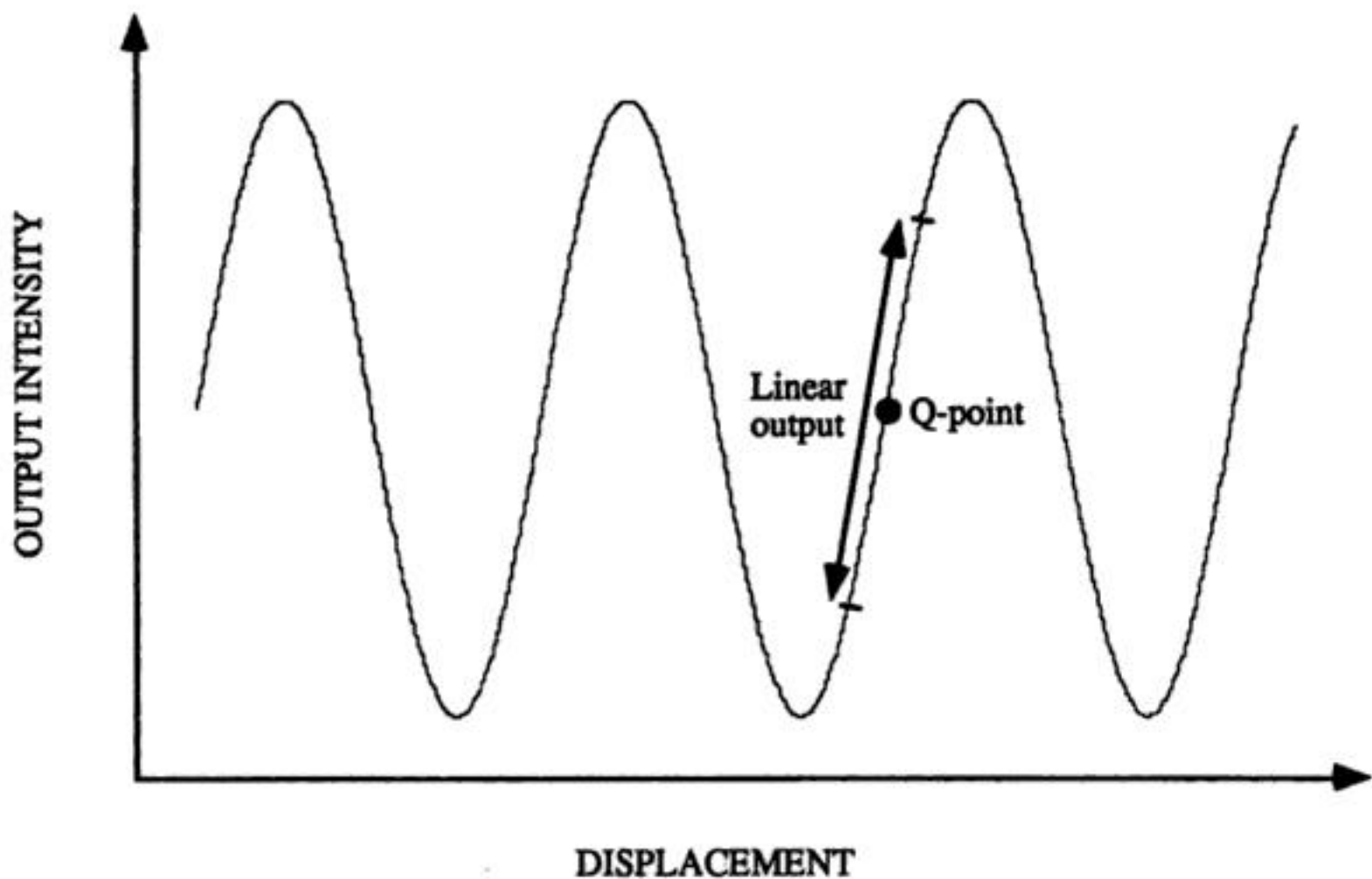


FIGURE 6.100 A typical EFPI transfer function curve.

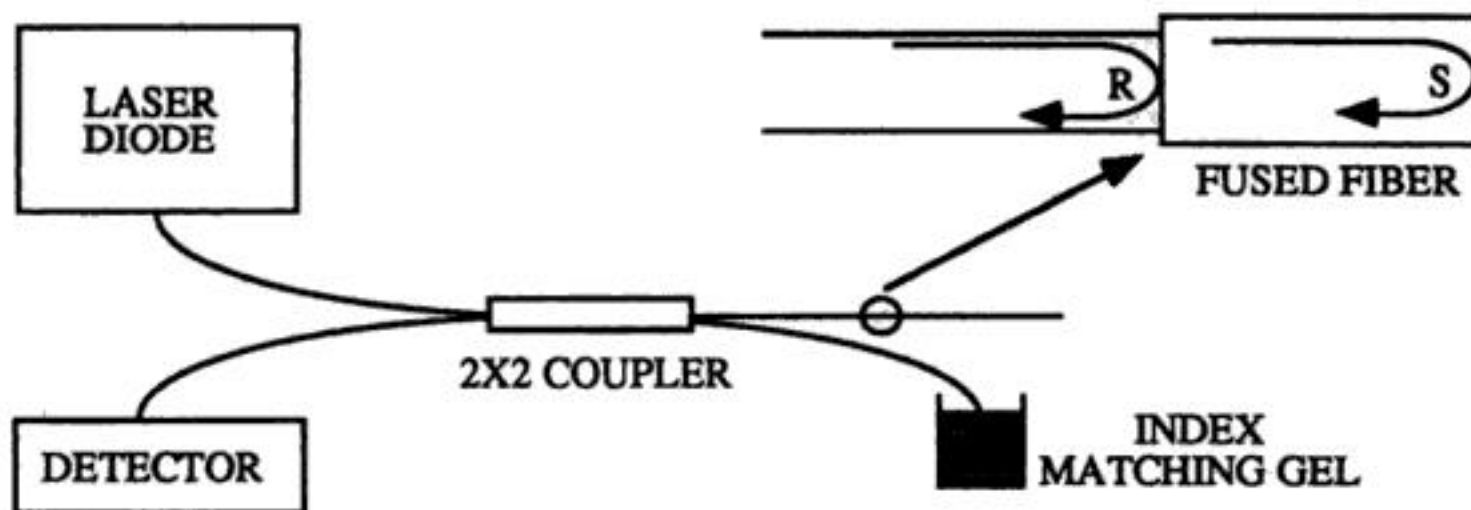


FIGURE 6.101 An intrinsic Fabry-Perot interferometric sensor (IFPI).

Although a version of the EFPI sensor that provides absolute output has been demonstrated, it lacks the bandwidth typically desired during the measurement of cyclical strain [5]. We have also recently proposed a small cavity length/high finesse EFPI sensor for measurement of small perturbations [6]. This configuration has a simple output that can be demodulated using an optical filter/photodetector combination.

Intrinsic Fabry-Perot Interferometric Sensor

The intrinsic Fabry-Perot interferometric (IFPI) sensor is similar in operation to its extrinsic counterpart but significant differences exist in the configurations of the two sensors [7]. The basic IFPI sensor is shown in Figure 6.101. An optically isolated laser diode is used as the optical source to one of the input arms of a bidirectional 2×2 coupler. The Fabry-Perot cavity is formed by fusing a small length of a single-mode fiber to one of the output legs of the coupler. As shown in Figure 6.101, the reference (R) and sensing (S) reflections interfere at the detector face to provide a sinusoidal intensity variation. The cavity can also be obtained by introducing two Fresnel reflectors — discontinuities in refractive index — along the length of a single fiber. Photosensitivity in germanosilicate fibers has been used in the past to fabricate broadband reflectors that enclose an IFPI cavity [8]. Since the cavity is formed within an optical

fiber, changes in the refractive index of the fiber due to the applied perturbation can significantly alter the phase of the sensing signal, S . Thus, the intrinsic cavity results in the sensor being sensitive to ambient temperature fluctuations and all states of strain.

The IFPI sensor, like all other interferometric signals, has a nonlinear output that complicates the measurement of large magnitude strain. This can again be overcome by operating the sensor in the linear regime around the Q-point of the sinusoidal transfer function curve. The main limitation of the IFPI strain sensor is that the photoelastic effect-induced change in index of refraction results in a nonlinear relationship between the applied perturbation and the change in cavity length. In fact, for most IFPI sensors, the change in propagation constant of the fundamental mode dominates the change in cavity length. Thus, IFPIs are highly susceptible to temperature changes and transverse strain components [1]. In embedded applications, the sensitivity to all the strain components can result in erroneous outputs. The fabrication process of an IFPI strain sensor is more complicated than that of the EFPI sensor since the sensing cavity must be formed within the optical fiber by some special procedure. The strain resolution of the IFPIs is also expected to be around $1 \mu\epsilon$ with an operating range greater than $10,000 \mu\epsilon$. IFPI sensors also suffer from drift in the output signal due to variations in the polarization state of the input light.

Thus, the preliminary analysis shows that the extrinsic version of the Fabry–Perot optical fiber sensor seems to have an overall advantage over its intrinsic version. The extrinsic sensor has negligible cross-sensitivity to temperature and transverse strain. Although the strain sensitivity, dynamic range, and bandwidth of the two sensors are comparable, the IFPIs can be expensive and cumbersome to fabricate due to the intrinsic nature of the sensing cavity.

The extrinsic and intrinsic Fabry–Perot interferometric sensors possess nonlinear sinusoidal outputs that complicate the signal processing at the detection end. Although intensity-based sensors have a simple output variation, they suffer from limited sensitivity to strain or other perturbations of interest. Grating-based sensors have recently become popular as transducers that provide wavelength-encoded output signals that can typically be easily demodulated to derive information about the perturbation under investigation. The advantages and drawbacks of Bragg grating sensing technology are discussed first. The basic operating mechanism of the Bragg grating-based strain sensor is elucidated and the expressions for strain resolution is obtained. These sensors are then compared to the recently developed long-period gratings in terms of fabrication process, cross-sensitivity to other parameters, and simplicity of signal demodulation.

Fiber Bragg Grating Sensor

The phenomenon of photosensitivity in optical fibers was discovered by Hill and co-workers in 1978 [9]. It was found that permanent refractive index changes could be induced in fibers by exposing the germanium-doped core to intense light at 488 or 514 nm. The sinusoidal modulation of index of refraction in the core due to the spatial variation in the writing beam gives rise to a refractive index grating that can be used to couple the energy in the fundamental guided mode to various guided and lossy modes. Later Meltz et al. proposed that photosensitivity is more efficient if the fiber is side-exposed to fringe pattern at wavelengths close to the absorption wavelength (242 nm) of the germanium defects in the fiber [10]. The side-writing process simplified the fabrication of Bragg gratings, and these devices have recently emerged as highly versatile components for communication and sensing systems. Recently, loading the fibers with hydrogen has been reported to result in two orders of magnitude higher index change in germanosilicate fibers [11].

Principle of Operation

Bragg gratings are based on the phase-matching condition between spatial modes propagating in optical fibers. This phase-matching condition is given by:

$$k_g + k_c = k_B \quad (6.119)$$

where, k_g , k_c , and k_B are, respectively, the wave-vectors of the coupled guided mode, the resulting coupling mode, and the grating. For a first-order interaction, $k_B = 2\pi/\Lambda$, where Λ is the grating periodicity. Since it is customary to use propagation constants while dealing with optical fiber modes, this condition reduces to the widely used equation for mode coupling due to a periodic perturbation:

$$\Delta\beta = \frac{2\pi}{\Lambda} \quad (6.120)$$

where, $\Delta\beta$ is the difference in the propagation constants of the two modes involved in mode coupling (both assumed to travel in the same direction).

Fiber Bragg gratings (FBGs) involve the coupling of the forward-propagating fundamental LP_{01} optical fiber waveguide propagation mode to the reverse-propagating LP_{01} mode [12]. Consider a single mode fiber with β_{01} and $-\beta_{01}$ as the propagation constant of the forward- and reverse-propagating fundamental LP_{01} modes. To satisfy the phase-matching condition,

$$\Delta\beta = \beta_{01} - (-\beta_{01}) = \frac{2\pi}{\Lambda} \quad (6.121)$$

where, $\beta_{01} = 2\pi n_{\text{eff}}/\lambda$ (n_{eff} is the effective index of the fundamental mode and λ is the free-space wavelength). Equation 6.121 reduces to [12]:

$$\lambda_B = 2\Lambda n_{\text{eff}} \quad (6.122)$$

where λ_B is termed the Bragg wavelength. The Bragg wavelength is the wavelength at which the forward-propagating LP_{01} mode couples to the reverse-propagating LP_{01} mode. This coupling is wavelength dependent since the propagation constants of the two modes are a function of the wavelength. Hence, if an FBG is interrogated by a broadband optical source, the wavelength at which phase-matching occurs is found to be reflected back. This wavelength is a function of the grating periodicity (Λ) and the effective index (n_{eff}) of the fundamental mode (Equation 6.122). Since strain and temperature effects can modulate both these parameters, the Bragg wavelength shifts with these external perturbations. This spectral shift is utilized to fabricate FBGs for sensing applications.

Figure 6.102 shows the mode coupling mechanism in fiber Bragg gratings using the β -plot. Since the difference in propagation constants ($\Delta\beta$) between the modes involved in coupling is large, Equation 6.120 reveals that only a small value of periodicity, Λ , is needed to induce this mode coupling. Typically for telecommunication applications, the value of λ_B is around $1.5 \mu\text{m}$. From Equation 6.122, Λ is determined to be $0.5 \mu\text{m}$ (for $n_{\text{eff}} = 1.5$). Due to the small periodicities (of the order of $1 \mu\text{m}$), FBGs are classified as short-period gratings.

Fabrication Techniques

Fiber Bragg gratings have commonly been manufactured using two side-exposure techniques: the interferometric method and the phase mask method. The interferometric method, depicted in Figure 6.103

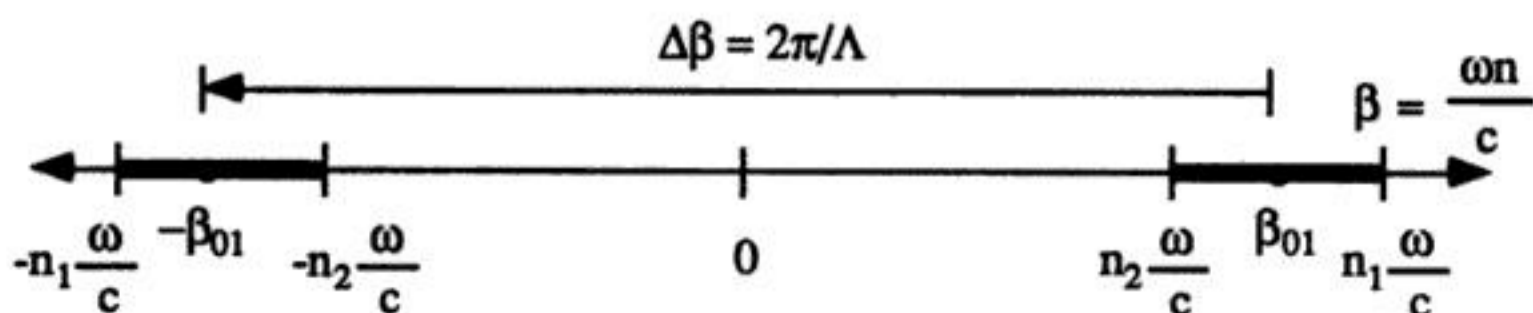


FIGURE 6.102 Mode coupling mechanism in a fiber Bragg grating.

where, $\Delta\Lambda/\Lambda$ and $\Delta n_{\text{eff}}/n_{\text{eff}}$ are the fractional changes in the periodicity and the effective index, respectively. The relative magnitudes of the two changes depend on the type of perturbation the grating is subjected to; for most applications, the effect due to change in effective index is the dominating mechanism.

Any axial strain, ϵ , applied to the grating changes the periodicity and the effective index and results in a shift in the Bragg wavelength, given by:

$$\frac{1}{\lambda} \frac{\Delta\lambda}{\epsilon} = \frac{1}{\Lambda} \frac{\Delta\Lambda}{\epsilon} + \frac{1}{n_{\text{eff}}} \frac{\Delta n_{\text{eff}}}{\epsilon} \quad (6.124)$$

The first term on the right-hand side is unity, while the second term has its origin in the photoelastic effect. An axial strain on the fiber serves to change the refractive index of both the core and the cladding. This results in the variation in the value of the effective index of glass. The photoelastic or strain-optic coefficient that relates the change in index of refraction due to mechanical displacement is about -0.27 . Thus, the variation in n_{eff} and Λ due to strain have contrasting effects on the Bragg peak. The fractional change in the Bragg wavelength due to axial strain is 0.73ϵ or 73% of the applied strain. At 1550 and 1300 nm, the shifts in the resonance wavelength are 11 nm/% ϵ and 9 nm/% ϵ , respectively. With temperature, a FBG at 1500 nm shifts by 1.6 nm for every 100°C rise in temperature [9].

Limitations of Bragg Grating Strain Sensors

The major limitation of Bragg grating sensors is the complex and expensive fabrication technique. Although side-writing is commonly used to manufacture these gratings, the requirement of expensive phase masks increases the cost of the sensing system. In the interferometric technique, stability of the setup is a critical factor in obtaining high-quality gratings. Since index changes of the order of 10^{-3} are required to fabricate these gratings, laser pulses of high energy levels are necessary. This might reduce laser operating lifetime and lead to increased maintenance expense. Additionally, introducing hydrogen or deuterium into the fiber allows increased index modulation as a result of the irradiation process.

The second major limitation of Bragg gratings is their limited bandwidth. The typical value of the full-width at half maximum (FWHM) is between 0.1 and 1 nm. Although higher bandwidths potentially can be obtained by chirping the index or periodicity along the grating length, this adds to the cost of the grating fabrication. The limited bandwidth requires high-resolution spectrum analyzers to monitor the grating spectrum. Kersey et al. have proposed an unbalanced Mach-Zender interferometer to detect the perturbation-induced wavelength shift [14]. Two unequal arms of the Mach-Zender interferometer are excited by the backreflection from a Bragg grating sensor element. Any change in the input optical wavelength modulates the phase difference between the two arms and results in a time-varying sinusoidal intensity at the output. This interference signal can be related to the shift in the Bragg peak and, hence, the magnitude of the perturbation can be obtained. Recently, modal interferometers have also been proposed to demodulate the output of a Bragg grating sensor [15]. The unbalanced interferometers are also susceptible to external perturbations and hence need to be isolated from the parameter under investigation. Moreover, the nonlinear output might require fringe counting equipment, which can be complex and expensive. Additionally, a change in the perturbation polarity at the maxima or minima of the transfer function curve will not be detected by this demodulation scheme. To overcome this limitation, two unbalanced interferometers can be employed for dynamic measurements.

The cross-sensitivity to temperature fluctuations leads to erroneous strain measurements in applications where the ambient temperature has a temporal variation. Thus, a reference grating that measures the temperature change must be utilized to compensate for the output of the strain sensor. Recently, temperature-independent sensing has been demonstrated using chirped gratings written in tapered optical fibers [16].

Last, the sensitivity of fiber Bragg grating strain sensors might not be adequate for certain applications. This sensitivity of the sensor depends on the minimum detectable wavelength shift at the detection end. Although excellent wavelength resolution can be obtained with unbalanced interferometric detection

techniques, standard spectrum analyzers typically provide a resolution of 0.1 nm. At 1300 nm, this minimum detectable change in wavelength corresponds to a strain resolution of 111 $\mu\epsilon$. Hence, in applications where strain smaller than 100 $\mu\epsilon$ is anticipated, Bragg grating sensors might not be practical. The dynamic range of strain measurement can be as much as 15,000 $\mu\epsilon$.

Long-Period Grating Sensor

This section discusses the use of novel long-period gratings as strain sensing devices. The principle of operation of these gratings, their fabrication process, preliminary strain tests, demodulation process, and cross-sensitivity to ambient temperature are analyzed.

Principle of Operation

Long-period gratings that couple the fundamental guided mode to different guided modes have been demonstrated in the past [17, 18]. Gratings with longer periodicities that involve coupling of a guided mode to forward-propagating cladding modes were recently proposed by Vengsarkar et al. [19, 20]. As stated previously, fiber gratings satisfy the Bragg phase-matching condition between the guided and cladding or radiation modes or, another guided mode. This wavelength-dependent phase-matching condition is given by:

$$\beta_{01} - \beta = \Delta\beta = \frac{2\pi}{\Lambda} \quad (6.125)$$

where Λ is the periodicity of the grating, β_{01} and β are the propagation constants of the fundamental guided mode and the mode to which coupling occurs, respectively.

For conventional fiber Bragg gratings, the coupling of the forward propagating LP_{01} mode occurs to the reverse propagating LP_{01} mode ($\beta = -\beta_{01}$). Since $\Delta\beta$ is large in this case (Figure 6.107(a)), the grating periodicity is small, typically of the order of 1 μm . Unblazed long-period gratings having index variations parallel to the long axis of the fiber couple the fundamental mode to the discrete and circularly-symmetric, forward-propagating cladding modes ($\beta = \beta^n$), resulting in smaller values of $\Delta\beta$ (Figure 6.107(b)) and

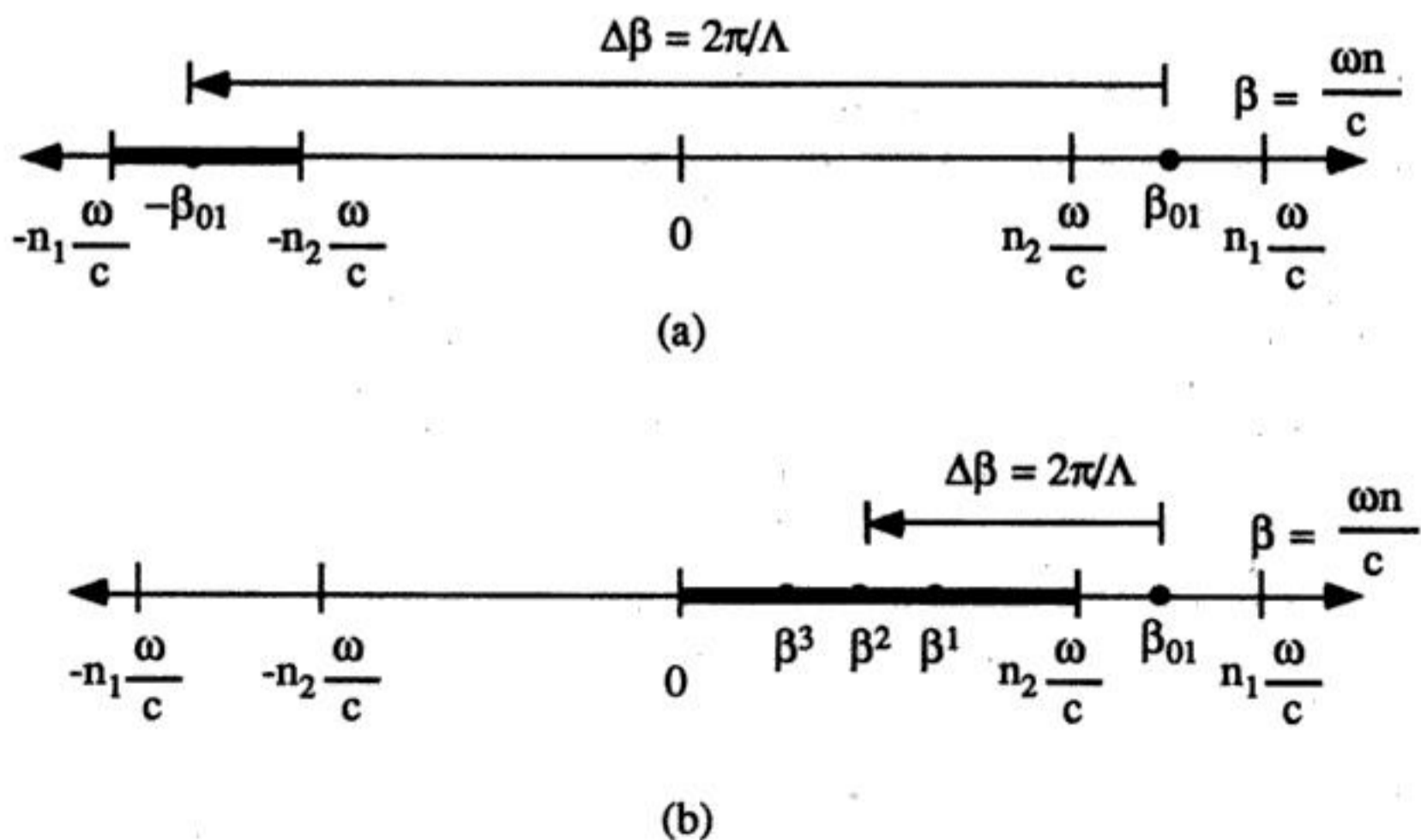


FIGURE 6.107

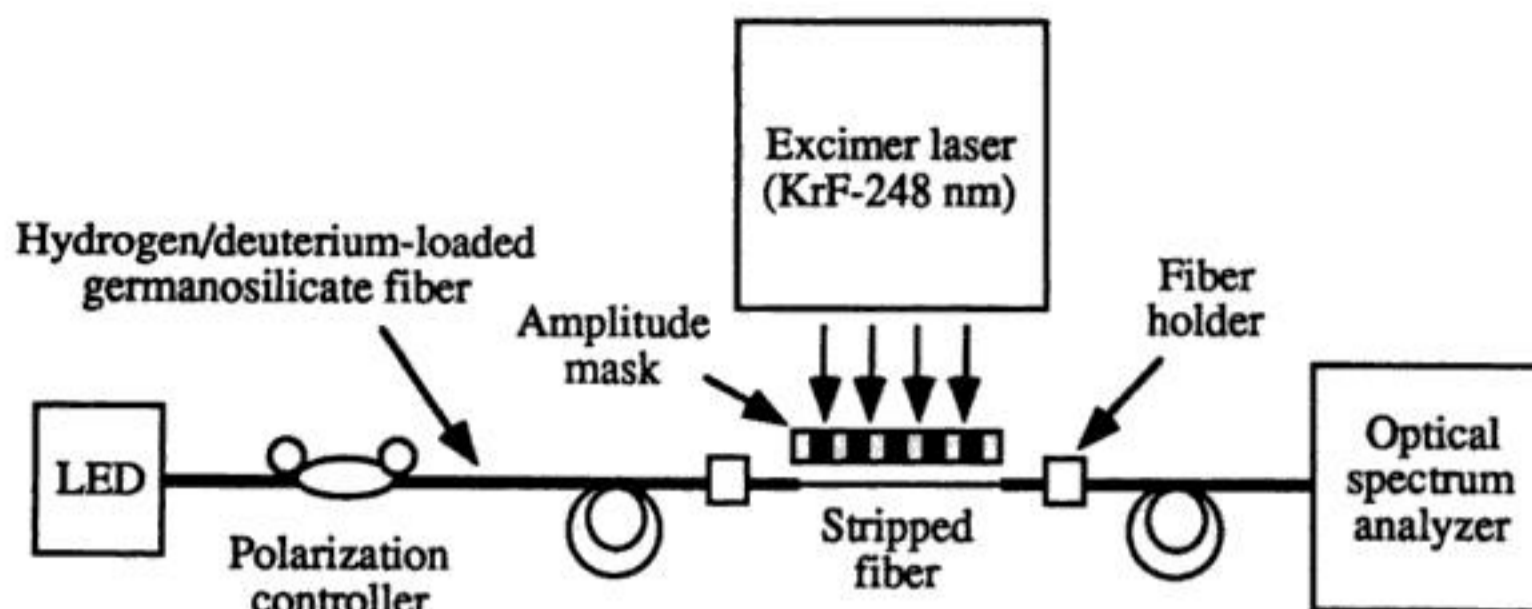


FIGURE 6.108 Setup used to fabricate long-period gratings.

hence periodicities ranging in hundreds of micrometers [19]. The cladding modes attenuate rapidly as they propagate along the length of the fiber due to the lossy cladding-coating interface and bends in the fiber. Since $\Delta\beta$ is discrete and a function of the wavelength, this coupling to the cladding modes is highly selective, leading to a wavelength-dependent loss. As a result, any modulation of the core and cladding guiding properties modifies the spectral response of long-period gratings, and this phenomenon can be utilized for sensing purposes. Moreover, since the cladding modes interact with the fiber jacket or any other material surrounding the cladding, changes in the properties of these ambient materials can also be detected.

Fabrication Procedure

To fabricate long-period gratings, hydrogen-loaded (3.4 mol%) germanosilicate fibers are exposed to 248 nm UV radiation from a KrF excimer laser, through a chrome-plated amplitude mask possessing a periodic rectangular transmittance function. Figure 6.108 shows the setup used to fabricate the gratings. The laser was pulsed at 20 Hz with a 8 ns pulse duration. The typical writing times for an energy of $100 \text{ mJ cm}^{-2} \text{ pulse}^{-1}$ and a 2.5 cm exposed length vary between 6 to 15 min for different fibers. The coupling wavelength, λ_p , shifts to higher values during exposure, due to the photoinduced enhancement of the refractive index of the fiber core and the resulting increase in β_{01} . After writing, the gratings are annealed at 150°C for 10 h to remove the unreacted hydrogen. This high-temperature annealing causes λ_p to move to shorter wavelengths due to the decay of UV-induced defects and diffusion of molecular hydrogen from the fiber. Figure 6.109 depicts the typical transmittance of a grating. Various attenuation bands correspond to coupling to discrete cladding modes of different orders. A number of gratings can be fabricated at the same time by placing more than one fiber behind the amplitude mask. Moreover, the stability requirements during the writing process are not as severe as those for short-period Bragg gratings.

For coupling to the highest-order cladding-mode, the maximum isolation (loss in transmission intensity) is typically in the 5 to 20 dB range on wavelengths, depending on fiber parameters, duration of UV exposure, and mask periodicity. The desired fundamental coupling wavelength can easily be varied using inexpensive amplitude masks of different periodicities. The insertion loss, polarization-mode dispersion, backreflection, and polarization-dependent loss of a typical grating are 0.2 dB, 0.01 ps, -80 dB, and 0.02 dB, respectively. The negligible polarization sensitivity and backreflection of these devices eliminate the need for expensive polarizers and isolators.

Preliminary experiments were performed to examine the strain sensitivity of long-period gratings written in different fibers [21, 22]. Gratings were fabricated in four different types of fibers: standard dispersion-shifted fiber (DSF), standard 1550 nm fiber, and conventional 980 and 1050 nm single-mode fibers, which for the sake of brevity are referred to as fibers A, B, C, and D, respectively. The strain sensitivity of gratings written in different fibers was determined by axially straining the gratings between

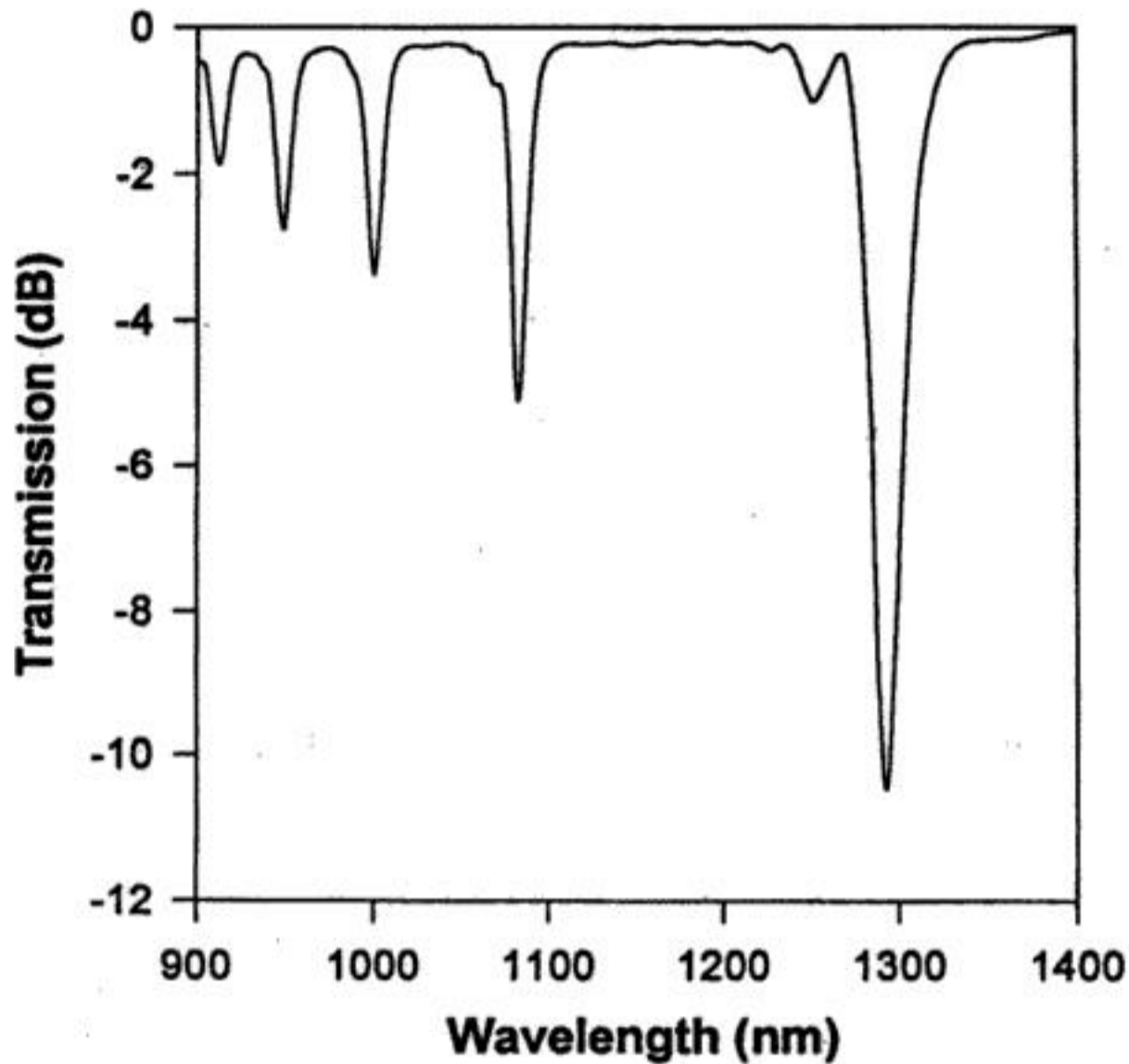


FIGURE 6.109 Typical transmission of a grating.

two longitudinally separated translation stages. The shift in the peak loss wavelength of the grating in fiber D as a function of the applied strain is depicted in Figure 6.110, along with that for a Bragg grating (about $9 \text{ nm } \% \epsilon^{-1}$, at 1300 nm) [9]. The strain coefficients of wavelength shift (β) for fibers A, B, C, and D are shown in Table 6.28. Fiber D has a coefficient $15.2 \text{ nm } \% \epsilon^{-1}$, which gives it a strain-induced shift that is 50% larger than that for a conventional Bragg grating. The strain resolution of this fiber for a 0.1 nm detectable wavelength shift is $65.75 \mu \epsilon$.

The demodulation scheme of a sensor determines the overall simplicity and sensitivity of the sensing system. Short-period Bragg grating sensors were shown to possess signal processing techniques that are complex and expensive to implement. A simple demodulation method to extract information from long-period gratings is possible. The wide bandwidth of the resonance bands enables the wavelength shift due to the external perturbation to be converted into an intensity variation that can be easily detected.

Figure 6.111 shows the shift induced by strain in a grating written in fiber C. The increase in the loss at 1317 nm is about 1.6 dB . A laser diode centered at 1317 nm was used as the optical source, and the change in transmitted intensity was monitored as a function of applied strain. The transmitted intensity is plotted in Figure 6.112 for three different trials. The repeatability of the experiment demonstrates the feasibility of using this simple scheme to utilize the high sensitivity of long-period gratings. The transmission of a laser diode centered on the slope of the grating spectrum on either side of the resonance wavelength can be used as a measure of the applied perturbation. A simple detector and amplifier combination at the output can be used to determine the transmission through the detector. On the other hand, a broadband source can also be used to interrogate the grating. At the output, an optical bandpass filter can be used to transmit only a fixed bandwidth of the signal to the detector. The bandpass filter should again be centered on either side of the peak loss band of the resonance band. These schemes are easy to implement, and unlike conventional Bragg gratings, the requirement of complex and expensive interferometric demodulation schemes is not necessary [22].

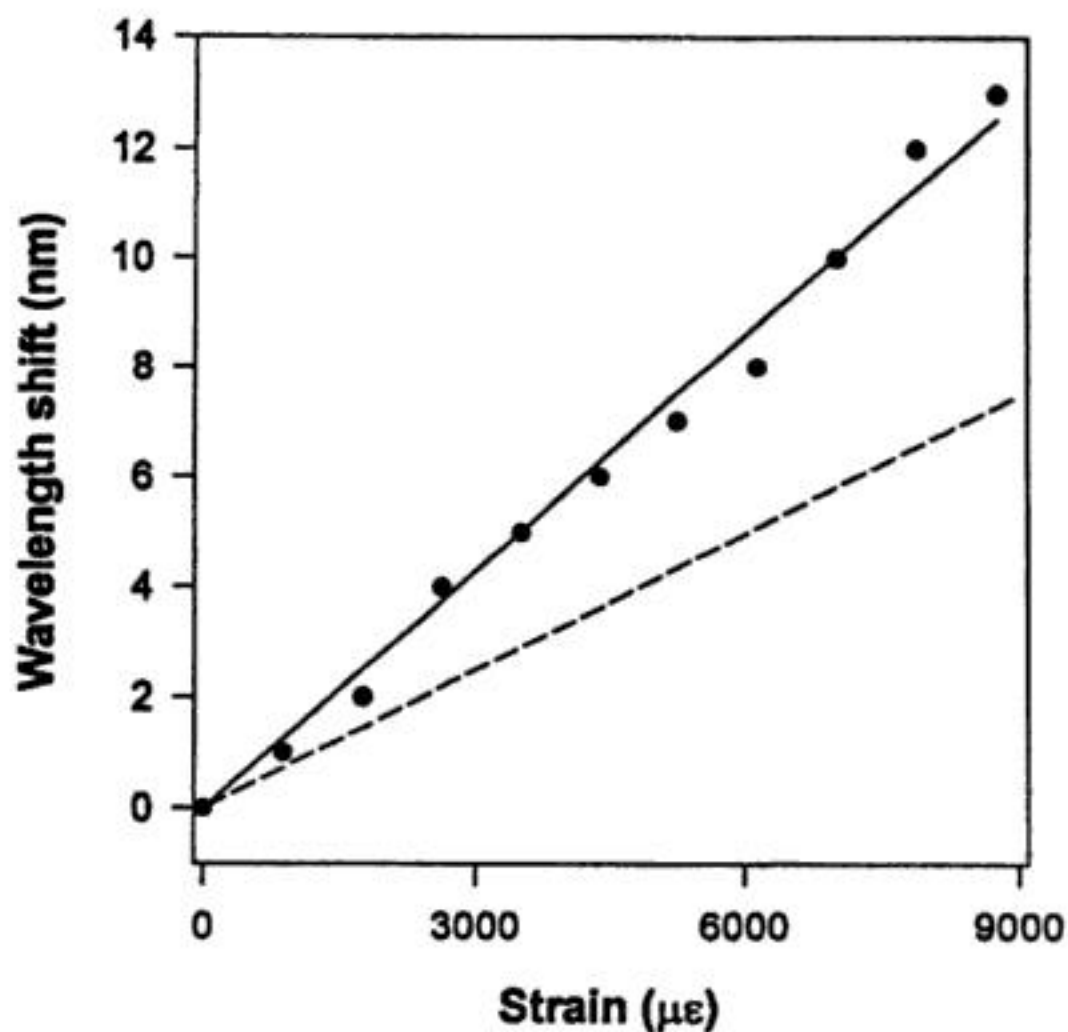


FIGURE 6.110 Shift in peak loss wavelength as a function of the applied strain.

TABLE 6.28 Strain Sensitivity of Long-Period Gratings Written in Four Different Types of Fibers

Type of fiber	Strain sensitivity (nm % ϵ^{-1})
A — Standard dispersion-shifted fiber (DSF)	-7.27
B — Standard 1550 nm communication fiber	4.73
C — Conventional 980 nm single-mode fiber	4.29
D — Conventional 1060 nm single-mode fiber	15.21

Note: The values correspond to the shift in the highest order resonance wavelength.

Temperature Sensitivity of Long-Period Gratings

Gratings written in different fibers were also tested for their cross-sensitivity to temperature [22]. The temperature coefficients of wavelength shift for different fibers are shown in Table 6.29. The temperature sensitivity of a fiber Bragg grating is $0.014 \text{ nm } ^\circ\text{C}^{-1}$. Hence, the temperature sensitivity of a long-period grating is typically an order of magnitude higher than that of a Bragg grating. This large cross-sensitivity to ambient temperature can degrade the strain sensing performance of the system unless the output signal is adequately compensated. Multiparameter sensing using long-period gratings has been proposed to obtain precise strain measurements in environments with temperature fluctuations [21].

In summary, long-period grating sensors are highly versatile. These sensors can easily be used in conjunction with simple and inexpensive detection techniques. Experimental results prove that these methods can be used effectively without sacrificing the enhanced resolution of the sensors. Long-period grating sensors are insensitive to the input polarization and do not require coherent optical sources. The cross-sensitivity to temperature is a major concern while using these gratings for strain measurements.

11. P. J. Lemaire, A. M. Vengsarkar, W. A. Reed, V. Mizrahi, and K. S. Kranz, Refractive index changes in optical fibers sensitized with molecular hydrogen, in *Proc. Conf. Optical Fiber Communications, OFC'94*, Technical Digest, paper TuL1, 47, 1994.
12. R. Kashyap, Photosensitive optical fibers: devices and applications, *Optical Fiber Technol.*, 1, 17-34, 1994.
13. D. Z. Anderson, V. Mizrahi, T. Ergodan, and A. E. White, Phase-mask method for volume manufacturing of fiber phase gratings, in *Proc. Conf. Optical Fiber Communication*, post-deadline paper PD16, 1993, p. 68.
14. A. D. Kersey and T. A. Berkoff, Fiber-optic Bragg-grating differential-temperature sensor, *IEEE Photonics Technol. Lett.*, 4, 1183-1185, 1992.
15. V. Bhatia, M. B. Sen, K. A. Murphy, A. Wang, R. O. Claus, M. E. Jones, J. L. Grace, and J. A. Greene, Demodulation of wavelength-encoded optical fiber sensor signals using fiber modal interferometers, *SPIE Photonics East*, Philadelphia, PA, paper 2594-09, October 1995.
16. M. G. Xu, L. Dong, L. Reekie, J. A. Tucknott, and J. L. Cruz, Chirped fiber gratings for temperature-independent strain sensing, in *Proc. First OSA Topical Meet. Photosensitivity and Quadratic Non-linearity in Glass Waveguides: Fundamentals and Applications*, paper PMB2, 1995.
17. K. O. Hill, B. Malo, K. Vineberg, F. Bilodeau, D. Johnson, and I. Skinner, Efficient mode-conversion in telecommunication fiber using externally written gratings, *Electron. Lett.*, 26, 1270-1272, 1990.
18. F. Bilodeau, K. O. Hill, B. Malo, D. Johnson, and I. Skinner, Efficient narrowband $LP_{01} \leftrightarrow LP_{02}$ mode converters fabricated in photosensitive fiber: spectral response, *Electron. Lett.*, 27, 682-684, 1991.
19. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, V. Bhatia, J. E. Sipe, and T. E. Ergodan, Long-period fiber gratings as band-rejection filters, *Proc. Conf. Optical Fiber Communications, OFC '95*, post-deadline paper, PD4-2, 1995.
20. A. M. Vengsarkar, P. J. Lemaire, J. B. Judkins, V. Bhatia, J. E. Sipe, and T. E. Ergodan, Long-period fiber gratings as band-rejection filters, *J. Lightwave Technol.*, 14, 58-65, 1996.
21. V. Bhatia, M. B. Burford, K. A. Murphy, and A. M. Vengsarkar, Long-period fiber grating sensors, *Proc. Conf. Optical Fiber Communication*, paper ThP1, February 1996.
22. V. Bhatia and A. M. Vengsarkar, Optical fiber long-period grating sensors, *Optics Lett.*, 21, 692-694, 1996.
23. C. D. Butter and G. B. Hocker, Fiber optics strain gage, *Appl. Optics*, 17, 2867-2869, 1978.
24. J. S. Sirkis and H. W. Haslach, Interferometric strain measurement by arbitrarily configured, surface mounted, optical fiber, *J. Lightwave Technol.*, 8, 1497-1503, 1990.

6.12 Optical Beam Deflection Sensing

Grover C. Wetsel

Measurements of the intensity of the light reflected and transmitted by a sample have been sources of information concerning the structure of matter for over a century. In recent decades, it has been found that measurement of the position of an optical beam that has scattered from a sample is an important and versatile means of characterizing materials and the motion of devices. Surely, the availability of a well-collimated beam from a laser has been crucial in the development of techniques and applications of *optical beam deflection* (OBD) sensing; however, the development and ready availability of various types of *position sensing detectors* (PSDs) have also been important factors. Optical beam deflection may be caused, for example, by propagation of a laser beam through a refractive-index gradient or by reflection from a displaced surface. A PSD provides an electronic signal that is a function of the laser beam position on the detector.

In this section, applications of optical beam deflection sensing are reviewed, the theories of operation of the three most common types of OBD sensors are developed, and typical operational characteristics of the devices are presented. The advantages and disadvantages of the various PSDs are also discussed.

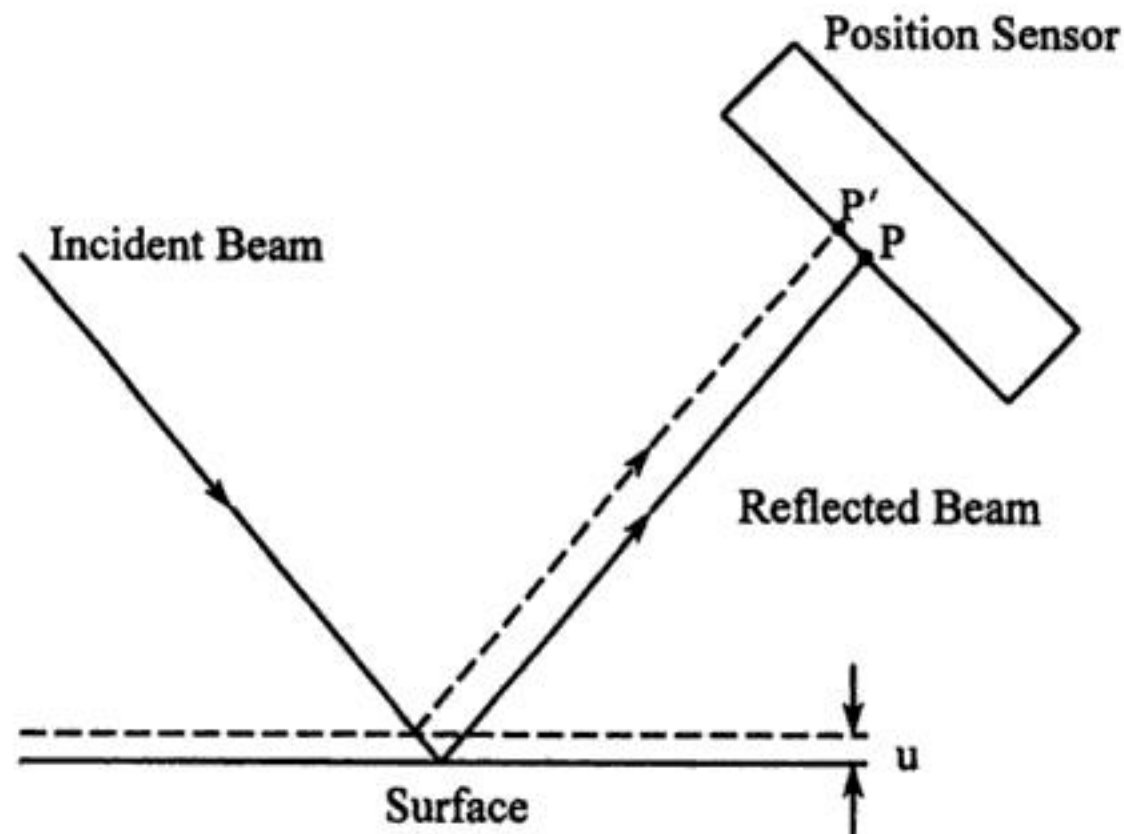


FIGURE 6.113 A schematic diagram of the basic optical-beam-deflection (OBD) sensing configuration.

A schematic diagram of the basic OBD sensing configuration is illustrated in Figure 6.113. In this case, the displacement, u , of the surface causes the position of the reflected beam on the PSD to move from point P to point P'; the positional change produces a change in the output voltage of the PSD. The output voltage, V , of the PSD electronics can be calibrated in terms of the actual displacement, u , by measuring V versus u for known displacements.

OBD sensing has been used in a variety of applications, including photothermal optical beam deflection (PTOBD) spectroscopy [1], absolute measurement of optical attenuation [2], PTOBD imaging of surface and subsurface structure [3], photothermal displacement spectroscopy [4], atomic-force microscopy [5], and materials characterization [6]. It has also been used as an uncomplicated, sensitive, and accurate method of measurement of surface motion for scanning tunneling microscope scanner transducers [7] and ultrasonic transducer imaging [8].

Theory

The three basic types of devices for OBD sensing are (1) a photodetector behind a sharply edged screen (a knife edge); (2) a small array of photodetectors separated by relatively small, insensitive areas (bicell, quadcell); and (3) a continuous solid-state position sensor (one or two dimensional). Sensing characteristics of a device are determined by the effect of optical beam displacement on the photodetector power distribution. Since laser beams are commonly used in OBD sensing, the analysis involves the assumption that the spatial distribution of the intensity (I) in the plane perpendicular to the direction of wave propagation is axially symmetric with a Gaussian radial variation.

Knife-Edge Photodetector

The essential features of a PSD are represented by a photodetector shadowed by a semi-infinite knife edge, $y < 0$, as illustrated in Figure 6.114. As can be anticipated from the symmetry of the arrangement and proved mathematically, the maximal deflection sensitivity occurs when the undeflected beam is centered on the knife edge. The intensity of the light reaching the photodetector due to the displacement (u) of the center of the beam is given in the reference frame of the displaced beam by:

$$I(r') = \frac{aP}{\pi} e^{-ar'^2} \quad (6.126)$$

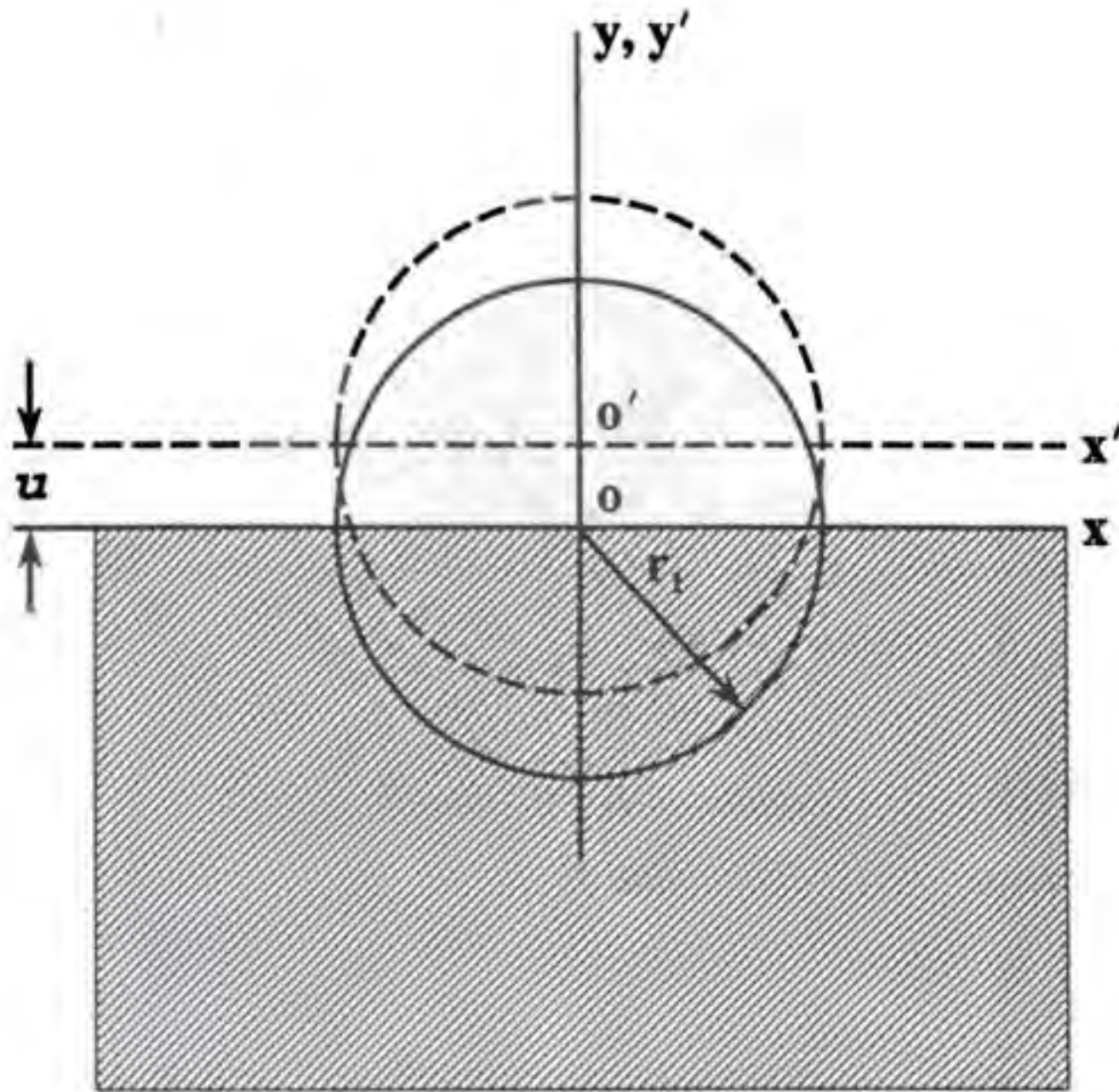


FIGURE 6.114 Essential features of a position-sending detector (PSD), as represented by a photodetector shadowed by a semiinfinite knife edge.

where P is the total incident beam power, $a = 2/r_1^2$, r_1 is the Gaussian beam radius, and $r'^2 = x'^2 + y'^2$. In terms of the coordinates (x, y) of the undeflected beam, the rectangular coordinates of the deflected beam are $x' = x$ and $y' = y - u$. The power (P_d) on the detector is thus given by:

$$P_d = \frac{aP}{\pi} \int_0^{\infty} e^{-a(y-u)^2} dy \int_{-\infty}^{\infty} e^{-ax^2} dx = \frac{P}{2} \left[1 + \operatorname{erf} \left(\sqrt{2} \frac{u}{r_1} \right) \right] \quad (6.127)$$

where erf is the error function. One can see by inspection of Equation 6.127 that the essential characteristics of this position sensor are determined by u/r_1 . The normalized response, P_d/P , is shown in Figure 6.115 as a function of u/r_1 . When $u = r_1$, then $P_d = 97.7\% P$.

The deflection sensitivity is given by the slope of Equation 6.127,

$$\frac{dP_d}{du} = \sqrt{\frac{2}{\pi}} \frac{P}{r_1} e^{-2\left(\frac{u}{r_1}\right)^2}, \quad \text{with} \quad \left(\frac{dP_d}{du} \right)_{\max} = \left(\frac{dP_d}{du} \right)_{u=0} = \sqrt{\frac{2}{\pi}} \frac{P}{r_1} \quad (6.128)$$

Define the small-signal position sensor sensitivity (units of m^{-1}):

$$\alpha_{\text{KE}} \equiv \frac{1}{P} \left(\frac{dP_d}{du} \right)_{u=0} = \frac{1}{r_1} \sqrt{\frac{2}{\pi}} \quad (6.129)$$

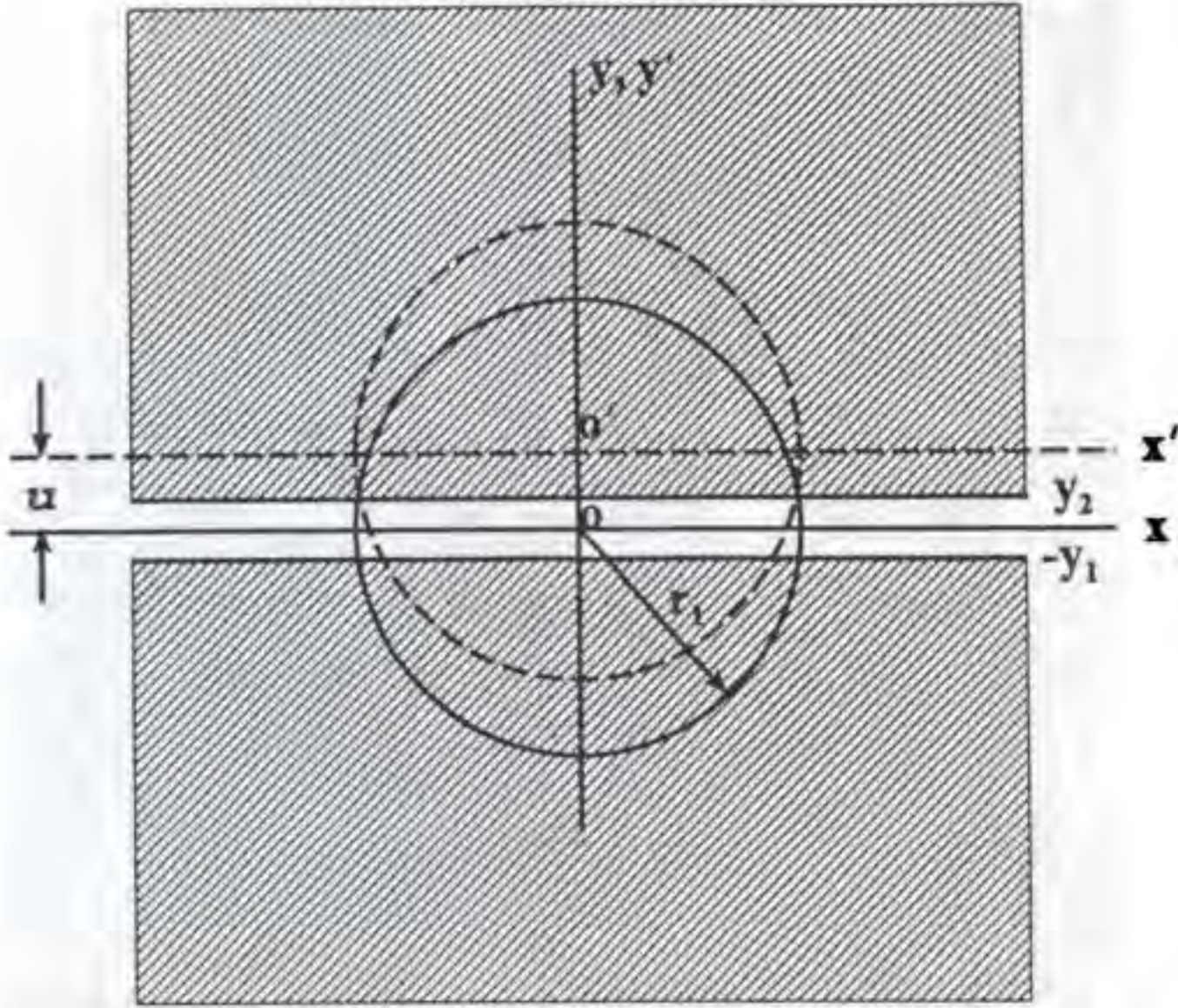


FIGURE 6.116 Deflection of a Gaussian beam initially centered in the insensitive gap of a bicell detector.

The power incident on the lower half of the bicell is given by:

$$P_1 = \frac{aP}{\pi} \int_{-\infty}^{-y_1} e^{-u(y-u)^2} dy \int_{-\infty}^{\infty} e^{-ax^2} dx = \frac{P}{2} \left[1 - \operatorname{erf} \left(\frac{\sqrt{2}}{r_1} (y_1 + u) \right) \right] \quad (6.134)$$

The photocurrent from each detector of the bicell is converted to voltage by identical transimpedance amplifiers: $V_2 = KZP_2$ and $V_1 = KZP_1$; a difference amplifier is then used to obtain the bicell signal voltage:

$$V = V_2 - V_1 = KZ(P_2 - P_1) = \frac{KZP}{2} \left[\operatorname{erf} \left(\frac{\sqrt{2}}{r_1} (y_1 + u) \right) - \operatorname{erf} \left(\frac{\sqrt{2}}{r_1} (y_2 - u) \right) \right] \quad (6.135)$$

The normalized response, $2V/(KZP)$, is shown in Figure 6.117 as a function of u/r_1 for $y_1 = y_2 = r_1/10$.

Suppose that the beam is centered in the gap, $y_1 = y_2$; then, for small displacements, one obtains:

$$V \cong 2 \sqrt{\frac{2}{\pi}} \frac{KZPu}{r_1} e^{-2(y_1/r_1)^2} \quad (6.136)$$

The small-signal sensitivity is:

$$\alpha_{\text{BC}} \equiv \frac{1}{KZP} \left(\frac{dV}{du} \right)_{u=0} = 2 \sqrt{\frac{2}{\pi}} \frac{e^{-2(y_1/r_1)^2}}{r_1} \quad (6.137)$$

This quantity is optimized when $r_1 = 2y_1$, and the optimal sensitivity is $0.484/y_1$.

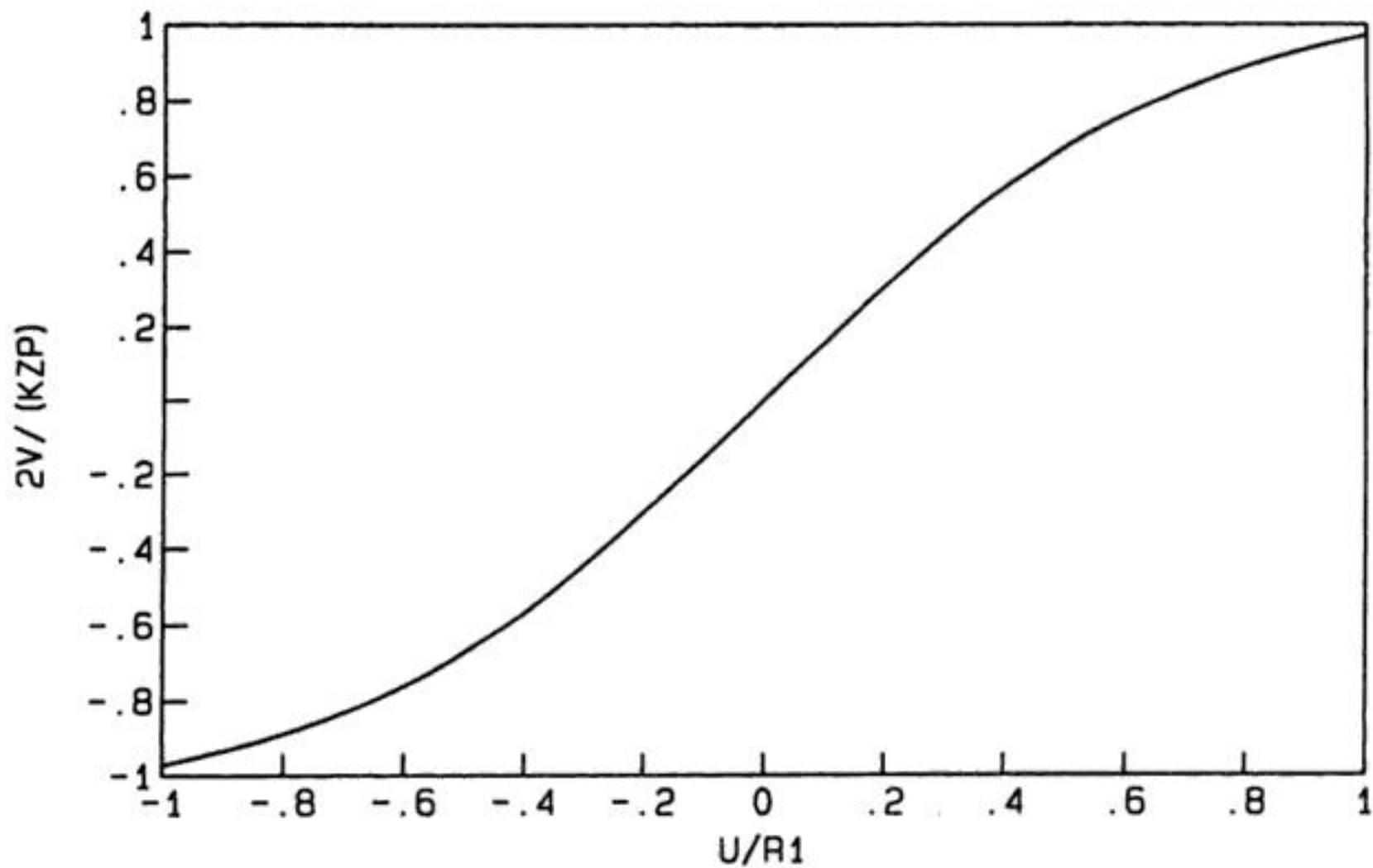


FIGURE 6.117 The normalized response $2V/KZP$ in a bicell detector as a function of u/r_1 for $y_1 = y_2 = r_1/10$.

Continuous Position Sensor

The position information in a continuous position sensor (also known as a lateral-effect photodiode) is derived from the divided path taken by photon-generated electrons to two back electrodes on the device. For a homogeneous device, the current to an electrode depends only on the distance of the centroid of the light beam from that electrode; the currents would be equal in an ideal device when the beam is located at its electrical center.

Consider the analysis of a one-dimensional continuous PSD. The current signal from each electrode is converted to a voltage signal by a transimpedance amplifier with gain, Z . Operational amplifiers are then used to provide the sum signal, $V_s = KPZA_s$, and the difference signal, $V_d = KPZA_d\alpha u$, where A_s and A_d are the sum and difference amplifier gains, respectively. The position sensor sensitivity, α , is then given by:

$$\alpha \equiv \frac{V_d A_s}{u V_s A_d} \quad (6.138)$$

which is determined by measuring V_d and V_s as a function of u .

Characterization of PSDs

The PSD characteristics presented here were measured by mounting the device to a translation stage with an optical encoder driven by a piezoelectric motor [9]; the position accuracy was $\pm 0.1 \mu\text{m}$. An appropriately attenuated He-Ne laser beam was directed at normal incidence to the PSD as it was translated past the beam. The PSD sum and difference voltages were measured with GP-IB digital voltmeters as a function of displacement; the PSD displacement and voltage measurements were computer controlled.

The operation of a knife-edge PSD can be evaluated using the signal from one side of a bicell as well as from a photodetector behind a knife edge. The transimpedance amplifier output voltage corresponding

5. G. Meyer and N. M. Amer, *Appl. Phys. Lett.*, 53, 1045, 1988.
6. J. C. Murphy and G. C. Wetsel, Jr., *Mater. Evaluation*, 44, 1224, 1986.
7. G. C. Wetsel, Jr., S. E. McBride, R. J. Warmack, and B. Van de Sande, *Appl. Phys. Lett.*, 55, 528, 1989.
8. S. E. McBride and G. C. Wetsel, Jr., Surface-displacement imaging using optical beam deflection, *Review of Progress in Quantitative Nondestructive Evaluation*, Vol. 9A, D. O. Thompson and D. E. Chimenti, (Eds.), New York: Plenum, 1990, 909-916.
9. Burleigh Instruments, Inc., Fishers, NY 14453.
10. United Detector Technology, 12525 Chadron Ave., Hawthorne, CA 90250.
11. On-Trak Photonics Inc., 20321 Lake Forest Dr., Lake Forest, CA 92630.

7

Thickness Measurement

John C. Brasunas
NASA/Goddard

G. Mark Cushman
NASA/Goddard

Brook Lakew
NASA/Goddard

- 7.1 Descriptions of the Relatively Mature Measuring Techniques.....7-2
Mechanical • Electronic Gages • Pneumatic Gaging • Optical: Focusing, Shadowing, Comparing • Weighing • Capacitive Gaging • Inductive Gaging (Eddy Current Sensing) • Magnetic Induction • Hall Effect Gage • Far-Field/Time-of-Flight: Ultrasound, Radar, Lidar • Far-Field/Resonance: Ultrasound, Interferometry, Ellipsometry • Far-Field/Absorption, Scattering, Emission: Beta, Gamma, X-Ray, Infrared • Destructive Techniques
- 7.2 Future Directions in Thickness Measurement7-7
Concerning Techniques Mentioned Above • THz Technology • Nanoscale-Scanning Probe Microscopy

One can measure thickness on many scales. The galaxy is a spiral disk about 100 Em (10^{20} m) thick. The solar system is pancake-like, about 1 Tm (10^{12} m) thick. The rings of Saturn are about 10 km thick. Closer to home, Earth's atmosphere is a spherical shell about 40 km thick; the weather occurs in the troposphere, about 12 km thick. The outermost shell of the solid Earth is the crust, about 35 km thick. The ocean has a mean depth of 3.9 km. In the Antarctic, the recently discovered objects believed to be microfossils indicative of ancient Martian life are less than 100 nm thick. In terms of the man-made environment, industry must contend with thickness varying from meters, for construction projects, to millimeters on assembly lines, to micrometers and nanometers for the solid-state, optical, and coatings industries. Perhaps the most familiar way of measuring thickness is by mechanical means, such as by ruler or caliper. Other means are sometimes called for, either because both sides of an object are not accessible, the dimension is either too big or too small for calipers, the object is too fragile, too hot, or too cold for direct contact, or the object is in motion on an assembly line — it may not even be a solid. Thickness may also be a function of position, as either the object may have originally been made with nonuniform thickness, deliberately or not, or the thickness may have become nonuniform with time due either to corrosion, cracking, or some other deterioration. The thickness may also be changing with time due to deliberate growth or etching, as example for thin films. Thus it follows that, in more general terms, measuring thickness might require measuring the topography or height profile of two surfaces and taking the difference. Alternatively, the measurement technique may produce a reading directly related to the difference. Table 7.1 lists some of the many techniques suited to determining thickness, together with the range of usefulness and some comments on accuracy and/or precision.

TABLE 7.1 Thickness Measuring Techniques

Technique	Range	Comments
Mechanical		
Caliper gage, micrometer	1 μm –100 mm	$\pm 3 \mu\text{m}$ accuracy
Electronic gages: LVDT	0–1 m	Precision depends on noise level
Pneumatic gaging	50 nm minimum	
Optical/focusing, shadowing, comparing		
Microscope	5 μm minimum	About 1% accuracy
Comparators/projectors	25–250 nm	
Laser caliper	100 μm –100 mm	Precision of 6 μm or better
Weighing		
Capacitive gaging	Range depends on area	
Inductive gaging (eddy current sensing)	From <1 μm to about 1 cm	
Magnetic induction	0–1.5 mm	Precision of 2.5 μm
Hall effect gage	0–4 mm	10% accuracy
Hall effect gage	0–10 mm	1–3% accuracy
Far-field/time-of-flight		
Sonar/ultrasound	0.5–250 mm	25 μm accuracy
Radar	0.1 to few hundred km	
Lidar	10 m–5 km	
THz technology		
Far-field/resonance		
Resonant ultrasound		
Interferometry: spectral and spatial	1 nm–100 μm	Accuracy about $\lambda/50$
Ellipsometry	0.3 nm–10 μm	0.1 nm accuracy
Far-field/absorption, scattering, emission		
Gamma-ray backscatter	Range to 25 mm	0.5% precision
Beta-transmission	2 μm –1 mm	0.2% precision
Beta-backscatter	100 nm–50 μm	3 to 20% precision
X-ray fluorescence	0–30 μm	
Infrared absorption	Depends on material	
Scanning techniques: scanning probe microscopy		Precision better than 0.1 nm
Destructive techniques: electrolytic	15 nm–50 μm	

7.1 Descriptions of the Relatively Mature Measuring Techniques

The following descriptions will also refer to some of the relevant vendors, whose addresses are found in Table 7.2. Additional vendor information, with specific price or model number identification, is found in Table 7.3. The words “gage” and “gauge” are used interchangeably.

Mechanical

The fundamental tool for measuring thickness is the line-graduated instrument [1, 2]. It is the only mechanical means to make direct measurements. Graduated spacings that represent known distances are used as direct comparisons to the unknown distance. Instruments include bars, rules, and tapes generically called rulers; caliper gages, which employ a positive contact device for improved alignment of the distance boundaries; and micrometers, which typically have greater precision due to a combination of linear and circumferential scales. Caliper precision can be improved with vernier scales or linear transducers. Fixed gages are often used to measure objects on a pass/fail basis. An object of fixed geometry (length, tapered bore, thread, etc.) is compared to a test piece typically for part inspection. Variations include the master gage, an object used to represent the nominal dimension of the part; the limit gage, an object used to represent the limit condition for tolerance dimensioning; and gage blocks or Johansson blocks, an object of fixed length used as a dimensional reference standard. Dial indicators are used to sense displacement

TABLE 7.2 Vendor addresses.

Vendor	Address
Bomem	Quebec, Canada
Brown & Sharpe	North Kingstown, RI
CMI International	Elk Grove Village, IL
Conductus	Sunnyvale, CA
deFelsko	Ogdensburg, NY
Digilab	Cambridge, MA
Digital Instruments	Santa Barbara, CA
Electromatic	Cedarhurst, NY
Fischer	Windsor, CT
Hewlett Packard	Englewood, CO
Kta-Tator	Pittsburgh, PA
Magnetic Analysis Corp.	Mount Vernon, NY
Mattson	Madison, WI
Measurex	Cupertino, CA
Micro Photonics	Allentown, PA
Midac	Costa Mesa, CA
Mitutoyo	Plymouth, MI
Moore Products Co	Spring House, PA
NDC Systems	Irwindale, CA
Nicolet	Madison, WI
Ono Sokki	Addison, IL
Oxford Instruments	Concord, MA
Panametrics	Waltham, MA
Park Scientific Instruments	Sunnyvale, CA
Penny + Giles	Attleboro, MA
Perkin Elmer	Norwalk, CT
Phase-Shift Technology	Tucson, AZ
Rudolf Instruments	Fairfield, NJ
Scantron	Dist. by Micro Photonics
Schaevitz	Pennsauken, NJ
Sentech	Dist. by Micro Photonics
SolveTech	Claymont, DE
Starrett	Athol, MA
Stresstel	Scotts Valley, CA
Transicoil Inc.	Valley Forge, PA
Trans-Tek	Ellington, CT
Willrich Precision Instrument Co.	Cresskill, NJ
J.A. Woolam Co., Inc.	Lincoln, NE
Wyko	Tucson, AZ
Zygo	Middlefield, CT

from a reference plane and display the deviation thereof. The display can be electronically coupled for amplification and/or display purposes. The range of a measuring instrument may be extended if multiple copies of the object to be measured are available. For example, the thickness of a sheet of paper may be measured by a simple ruler if 500 sheets of paper are stacked. (Vendors: Brown & Sharpe, Starrett, Mitutoyo. Also see [3].)

Electronic Gages

A Linear Variable Differential Transformer (LVDT), utilizes multiple toroidal transformers to sense axial displacement of an iron core that is attached to a measuring contact, either directly or by another joint (such as a lever). The displacement has a direct correlation to the distance that other electronics display. Thus, the LVDT serves as a replacement for a lined ruler or micrometer, incorporating an electrical readout. (Vendors: Penny + Giles; Schaevitz; Transicoil Inc.; Trans-Tek.)

TABLE 7.3 Instruments for measuring thickness.

Manufacturer	Model Number	Price	Description
KTa-Tator	TI-12	\$1595.	General-purpose ultrasonic gage, 0.75 mm to 75 mm range
NDC Systems	6100TC	\$49,300.	Backscatter gamma gage for 60 in. web, 25 mm range
NDC Systems		\$66,600.	Transmission beta gage, for continuous web products
Panametrics	25DL	\$2200. to \$3800.	Single-element ultrasonic gage, 50 mm range
Panametrics	26DL Plus	\$1400. to \$2500.	Dual-element ultrasonic gage, 250 mm range
Panametrics	8000	\$6500.	Hall effect magnetic gage, for nonferrous materials, 6 mm range
DeFelsko	Positest 1000-N	\$1995.	Eddy current sensor, Apple Newton read-out, measure out to 1.5 mm nonferrous, nonconducting coating on conducting substrate
Magnetic Analysis	Various	\$1500. to \$100,000.	Ultrasonic, time-of-flight gages
Fischer	Deltascope MP2C	\$1200.	Magnetic induction gage, measure nonmagnetic coating on ferromagnetic substrate
Fischer	IsoScope MP1C	\$1200.	Eddy current gage, measure nonconducting coating on nonferrous conducting substrate
Fischer	Fischerscope MMS	\$6500.	Beta-backscatter system to measure coating thickness
Fischer	Fischerscope X-Ray 1020 video	\$34,000.	X-ray fluorescence system to measure coating thickness
Fischer	Coulescope Sx	\$2500.+ accessories	Electrolytic, destructive system to measure coating thickness
J.A. Woollam Co. Inc.	M-44	Application specific	Variable angle, multiwavelength spectroscopic ellipsometer
Rudolf Instruments	431A31WL633	\$10,100.	Manual, HeNe wavelength ellipsometer
Rudolf Instruments	444A12	\$34,000.	Automatic, HeNe wavelength ellipsometer
Hewlett Packard	HP8712C	\$13,500.	RF vector network analyzer, measure transmission/reflection frequency response to 1.3 GHz, optional to 3 GHz
Stresstel	T-Mike Programmable	\$995.	Dual-element ultrasonic system
Stresstel	TM1D	\$1795.	Single-element ultrasonic system
Measurex	DMC480	Application specific	High-speed X-ray thickness gage
Bomem	MB series	\$20,000. and up	1 cm ⁻¹ resolution Fourier transform spectrometer
Park Scientific Instruments	Autoprobe CP	\$65,000.	Ambient scanning probe microscope
Park Scientific Instruments	Autoprobe VP2	\$130,000.	UHV scanning probe microscope
Digital Instruments	Nanoscope IIIa/D3000	\$90,000.	Small sample scanning probe microscope

Pneumatic Gaging

Pneumatic gages have pressurized air exiting gage orifices. The air velocity differential or backpressure is a function of the separation of the gage and the part. In the direct or open jet method, the pressurized air experiences backpressure due to the impedances posed by the measured part. The typical scenario is that the gage head and the measured part have similar geometry (i.e., a cylindrical gage in a bored hole). By placing two gages on either side of a flat plate, the thickness may be inferred. In the indirect or contact method, the pressurized air pushes on a contact piece that directly contacts the part. Tolerances as small as 50 nm can be measured. (Vendors: Willrich Precision Instrument Co.; Moore Products Co.)

Optical: Focusing, Shadowing, Comparing

This includes microscopes, which can determine thickness either by comparison with a known reference, or by focusing on the front and rear surfaces of a sample, noting the difference in focus position. Comparators project onto a screen what might be noted through a microscope. Laser calipers retrieve dimensions by measuring the shadowing of a laser beam. (Vendors: NDC Systems for laser caliper; Scantron for laser profilometer.)

Weighing

Given a plate of material with known density, first measure the area with some type of calibrated video system. Then, a measurement of weight can be simply converted to an estimate of the thickness. As is common with this technique and most of the following techniques, estimating the thickness requires knowledge of some other property of the material to be measured — in this case, the density.

Capacitive Gaging

Capacitive gaging is realized by inserting a nonmetallic material into a known electric field. Knowing the gage sensor area and the material's dielectric constant, the thickness can be determined. Submicron thickness levels can be achieved. (Vendors: Ono Sokki; SolveTech.)

Inductive Gaging (Eddy Current Sensing)

The principle here is that ac currents in a coil induce eddy currents in a nearby conducting plate [4, 5]. These eddy currents can be sensed by a pickup coil, which may be the exciting coil or a second coil. The presence of the eddy currents manifests itself as a modification of the apparent inductance and/or the loss of the pickup coil. This technique is appropriate for nonferrous metals, and is especially sensitive to thickness variations due to flaws such as cracks or corrosion. There is one particular instance in which it is common to measure thickness rather than variations. That would be the thickness of a nonconducting coating on a nonferrous conducting substrate. The coating thickness creates a gap (lift-off) between the exciting coil and the eddy currents, thereby affecting the eddy current signal. The range of this technique would be about 1 mm. Fischer has an instrument designed for measuring the thickness of a newly laid road surface coating to a depth of 40 cm, by burying a conductive plate below the road. (Vendors: Fischer; deFelsko; CMI International.)

Magnetic Induction

This technique is also used to measure coating thickness, in this case a nonmagnetic coating on a ferromagnetic substrate. The nonmagnetic coating creates a gap (lift-off) between the ferromagnetic substrate and a probe. One way to measure the gap and thereby the thickness is by measuring the force required to pull away a magnetic probe. Another technique would be to magnetically couple the ferromagnetic substrate to a transformer core, with a gap between the substrate and the core. This technique would have a range of about 4 mm. CMI International has an informative brochure describing the relative merits of measuring coating thickness via eddy current, magnetic induction, beta-backscatter, microresistance, and X-ray fluorescence; the choice of technique depends, among other things, on the material to be tested. (Vendors: Fischer; CMI International; Electromatic; deFelsko.)

Hall Effect Gage

This sensor measures the thickness of nonferrous materials with 1% accuracy by sandwiching the material being measured between a magnetic probe on one side and a small target steel ball on the other side [6].

by varying the angle of incidence and by observing at multiple wavelengths. The ability to estimate both thickness and refractive index is an important advantage of this technique, as often the refractive index of a material in thin film form is not the same as the bulk value, and indeed may be a property of the deposition conditions. (FTS vendors: Bomem; Digilab; Mattson; Midac; Nicolet; Perkin-Elmer. FTS system pricing may range from about \$15,000 to over \$100,000, depending on the application. Spatial interferometer vendors: Zygo; Wyko, Phase-Shift Technology. Ellipsometer vendors: J.A. Woollam Co., Inc.; Rudolf Instruments; Sentech. The cost of an ellipsometer may range from \$10,000 for a manual, single-wavelength system to \$200,000 for an automatic, multiwavelength system. Microwave resonance vendor: Hewlett-Packard.)

Far-Field/Absorption, Scattering, Emission: Beta, Gamma, X-Ray, Infrared

These techniques depend on the extinction (scattering or absorption) or emission of photons or massive particles (electrons, protons, neutrons) when transiting the material to be measured. Typically, the extinction or emission shows an exponential dependence on thickness; the dependence becomes linear if the absorption is sufficiently low. These techniques, in particular gamma-ray backscatter and beta-ray transmission, are used to measure continuously moving web materials (paper, metals, fabrics) on assembly lines. Infrared absorption is also suitable if the moisture content is controlled. Beta-backscatter and X-ray fluorescence [10] are used for measuring coatings. In X-ray fluorescence, upon exposure to X-rays, certain elements fluoresce (emit) X-rays at characteristic wavelengths. The strength of this emission is related to thickness. These absorption/emission techniques may sometimes be better suited than time-of-flight ultrasound to the dimensional measurement of objects with complex shapes. (Gamma gage vendor; NDC Systems. X-ray absorption vendor: Measurex. X-ray fluorescence vendors: Fischer; NDC Systems; CMI International. Beta-backscatter vendor: Fischer; Electromatic; CMI International; Measurex. Infrared absorption vendor: NDC Systems. The prices for these systems will depend on the application; a typical system could cost \$500,000.)

Destructive Techniques

Fischer markets a system that removes a coating into an electrolyte and then electrolytically deposits the removed coating. The electrical charge required for deposition is related to the coating thickness.

7.2 Future Directions in Thickness Measurement

Concerning Techniques Mentioned Above

Concerning capacitive sensors, the NASA Langley Research Center is developing sensors based on patterns of conductors sandwiched between insulating layers. The presence of ice over the conductors changes the capacitance, providing a way of sensing ice build-up on aircraft wings. With respect to eddy current sensing, one limitation is that a nonsuperconducting sense coil responds best to high-frequency excitations, and not at all to dc magnetic fields. This limits the technique to fairly high frequencies and thus low penetration depths, since the skin depth becomes shallower with increasing frequency. One possibility is to use a SQUID (superconducting quantum interference detector) as the sensor, since the SQUID is probably the most sensitive sensor of dc and low-frequency magnetic fields. One disadvantage of the SQUID has been the need for liquid helium for cooling for low-temperature superconductors; with the recent availability of high-temperature superconductors (HTS, above 90 K) and now HTS SQUIDS, cooling can be done with liquid nitrogen or single-stage mechanical coolers. In the area of spatial interferometry, work at Lawrence Livermore National Laboratory replaces the reference surface with a single-mode fiber in a process called phase-shifting diffraction interferometry. A measurement accuracy of 1.44 nm rms is quoted, with a goal of 0.1 nm rms. (HTS SQUID vendor: Conductus.)

THz Technology

With the availability of femtosecond pulsed lasers, Bell Labs has been investigating a technique using 100 fs pulses to pulse an antenna in the range of 0.1 THz to 3.0 THz. The terahertz pulses are sent through the material to be tested, detected, and the received pulse shape is analyzed to extract constituent information. This technique may also provide information on thickness.

Nanoscale-Scanning Probe Microscopy

Scanning probe microscopes (SPMs) are used in a wide variety of disciplines, including fundamental surface science, routine surface roughness analysis, and spectacular three-dimensional imaging — from atoms of silicon to micron-sized protrusions on the surface of a living cell [11]. The scanning probe microscope is an imaging tool with a vast dynamic range, spanning the realms of optical and electron microscopes. It is also a profiler with unprecedented 3-D resolution. In some cases, scanning probe microscopes can measure physical properties such as surface conductivity, static charge distribution, localized friction, magnetic fields, and elastic moduli. As a result, applications of SPMs are very diverse. The scanning tunneling microscope (STM), the progenitor of SPMs, utilizes a sharp conductive tip with a bias voltage applied between the tip and the sample. When the tip is within 1 nm of the sample, electrons from the sample begin to tunnel through the 1 nm gap into the tip. If the bias voltage is reversed, the tunneling occurs into the sample. The tunneling current is a function of the separation. Both the tip and the sample must be conductors or semiconductors.

The atomic force microscope (AFM) utilizes a small tip at the end of a cantilever. Forces between the tip and sample cause a deflection in the cantilever, which is translated into a signal. The tip or sample can be scanned covering a large area, producing a topographical map. AFMs can be used on insulators or conductors. AFMs are used in two modes: contact and noncontact. In contact mode, the tip is brought within about 200 pm — about the length of a chemical bond. The electron clouds of the tip and sample atoms interact, netting a repulsive force. For this reason, the contact mode is also called repulsive. Vertical resolution of about 50 pm can be achieved. In noncontact mode, a vibrating cantilever is used in the attractive regime of the van der Waals interactions. The cantilever is typically 2 nm to 20 nm away from the sample surface and has low total force. Noncontact AFM is subsequently less sensitive; thus, sensitive ac detection systems must be employed. The low force does have the advantage of not contaminating the sample surface and is preferred for applications involving silicon wafers and soft or elastic tissues. In noncontact mode, the cantilever is resonated with a small amplitude. As the tip comes near the sample surface, the resultant force changes the spring constant, translating into a deviation of the resonance frequency. This change in resonance (or vibrational amplitude) reflects changes in the sample topology.

Intermittent-contact mode is a combination of noncontact and contact modes and best suited for soft, adhesive, or fragile samples. Contact mode can damage the tip and the sample due to frictional or shear forces and/or create data artifacts from tip/surface adhesion. Noncontact mode produces lower amplitudes and hence lower resolution. Furthermore, surface monolayers of adsorbed gases such as water vapor can produce erroneous results. Intermittent-contact mode avoids these pitfalls by placing the tip in contact with the surface, providing high resolution and then removing the tip to prevent dragging and/or lateral forces. The cantilever is resonated via a piezoelectric crystal (50 kHz to 500 kHz in ambient, 5 kHz to 40 kHz in fluids) overcoming the tip/sample adhesion forces.

In magnetic force microscopy (MFM), the noncontact mode is employed using a tip coated with a ferromagnetic film. Both magnetic and van der Waals interactions are present, but at larger tip/sample separations, the magnetic forces dominate. Multiple scans as a function of tip/sample distance allow differentiation of magnetic forces and topographic information. Magnetic domain structures are resolved to 50 nm via this technique. Current applications of MFM include data storage devices, imaging of micromagnetic structures, IC analysis, imaging of magnetotactic bacteria, and magnetic geophysics. Lateral force microscopy (LFM) is used to generate profiles of changes in surface friction and/or height variations. The probe tip is deflected laterally, indicating some sort of twist. Electronics measure the cantilever deflection. To differentiate between the two effects, LFM and AFM images should be obtained

simultaneously. Phase detection microscopy or phase imaging is an extension of intermittent-contact AFM. It utilizes the phase lag between the driving frequency (cantilever) and the output signal frequency, generating a map of specific mechanical properties such as adhesion, elasticity, and friction. Identification of contaminants, composite materials, and regions of hardness and low surface adhesion can be obtained at the nanometer scale. Additional techniques include force modulation microscopy, where a periodic signal is applied to the cantilever, generating a map of the sample's elastic modulus and/or contaminants; electrostatic force microscopy, where a charged tip is scanned over the sample, revealing the locally charged domains generating a map of the charge carrier density; scanning capacitance microscopy, where a charged tip, kept at a constant tip/sample distance, generates a map of capacitance correlated information such as dielectric material thickness and subsurface charge carrier distributions (i.e., dopant profiles of ion implanted semiconductors); thermal scanning microscopy, where the tip in noncontact mode and a bimetal cantilever are used to map the thermal conductivity of the sample. (Vendors: Park Scientific Instruments; Digital Instruments; Oxford Instruments.)

References

1. R. E. Green (ed.), *Machinery's Handbook, 24th ed.*, New York: Industrial Press, 1992.
2. F. T. Farago and M. A. Curtis, *Handbook of Dimensional Measurement, 3rd ed.*, New York: Industrial Press, 1994.
3. T. Busch, *Fundamentals of Dimensional Metrology, 2nd ed.*, Albany, NY: Delmar Publishers, 1989.
4. R. C. McMaster, P. McIntire, and M. L. Mester (eds.), *Nondestructive Testing Handbook, Vol. 4, 2nd ed.*, American Society for Nondestructive Testing, 1986.
5. D. E. Bray and D. McBride, *Nondestructive Testing Techniques*, New York: John Wiley & Sons, 1992.
6. M. Giannini and A. deChiara, Wall Thickness Gaging in the Blow Mold Industry, distributed by Panametrics.
7. A. S. Birks, R. E. Green, and P. McIntire (eds.), *Nondestructive Testing Handbook, Vol. 7, 2nd ed.*, American Society for Nondestructive Testing, 1991.
8. D. Malacara, *Optical Shop Testing*, New York: John Wiley & Sons, 1978.
9. J. A. Woollam and P. G. Snyder, Variable Angle Spectroscopic Ellipsometry, VASE, distributed by J.A. Woollam Co.
10. H. H. Behncke, Coating thickness measurement by the X-ray fluorescence method, *Metal Finishing*, May, 33-39, 1984.
11. R. Howland and L. Benatar, A Practical Guide to Scanning Probe Microscopy, Park Scientific Instruments, 1993.

8

Proximity Sensing for Robotics

R. E. Saad
University of Toronto

A. Bonen
University of Toronto

K. C. Smith
University of Toronto

B. Benhabib
University of Toronto

8.1	Proximity Definition.....	8-1
8.2	Typical Sensor Characteristics.....	8-2
8.3	Technologies for Proximity Sensing.....	8-2
	Electro-Optical Sensors • Capacitive Sensors • Ultrasonic Sensors • Magnetic Sensors	

The objective of this chapter is to review the state-of-the-art in proximity-sensing technologies for robotics. Special attention is paid to the sensing needs of robotic manipulators for grasping applications, in contrast to the needs of mobile robots for navigation purposes. For a review of the application of proximity sensing to mobile robots, the reader is referred to [1].

Robotic sensors can be categorized into three groups: medium-range (object recognition and gross position/orientation estimation) sensors, short-range (proximity) sensors, and contact sensors. Recent literature [2–6] suggests that robotic end effectors should be equipped with both short-range proximity and contact sensors.

Proximity sensors should be able to measure the position and orientation (pose) of an object's surface. The range must be sufficiently large to compensate for uncertainties in the medium-range pose-estimation process, while maintaining sufficient accuracy to permit effective grasping of the object.

Transducers used by current proximity sensors vary in sophistication. Despite their great variety, however, these transducers and their accompanying electronic interface circuits (together comprising the proximity sensor) cannot presently meet the stringent robustness requirements of most industrial robotic applications. Novel sensing algorithms and techniques still must be developed in order to improve on their current characteristics, and, furthermore, to control both the sensing and grasping processes.

8.1 Proximity Definition

The term "proximity," quantified by "pose" in this chapter, refers to three geometrical parameters x , u , and v as shown in Figure 8.1, where:

x = the translation from the origin of the sensor's reference coordinate frame, F_p , to a target point on the surface of the object measured along X_p . This target point defines the origin of the surface-frame, F_r ,

u = the *vertical* orientation of the object's surface, defined as a rotation around Y_p (of the translated frame), thereby specifying the new Z_r ,

v = the *horizontal* orientation of the object's surface, defined as a rotation around Z_r , thereby specifying Y_r .

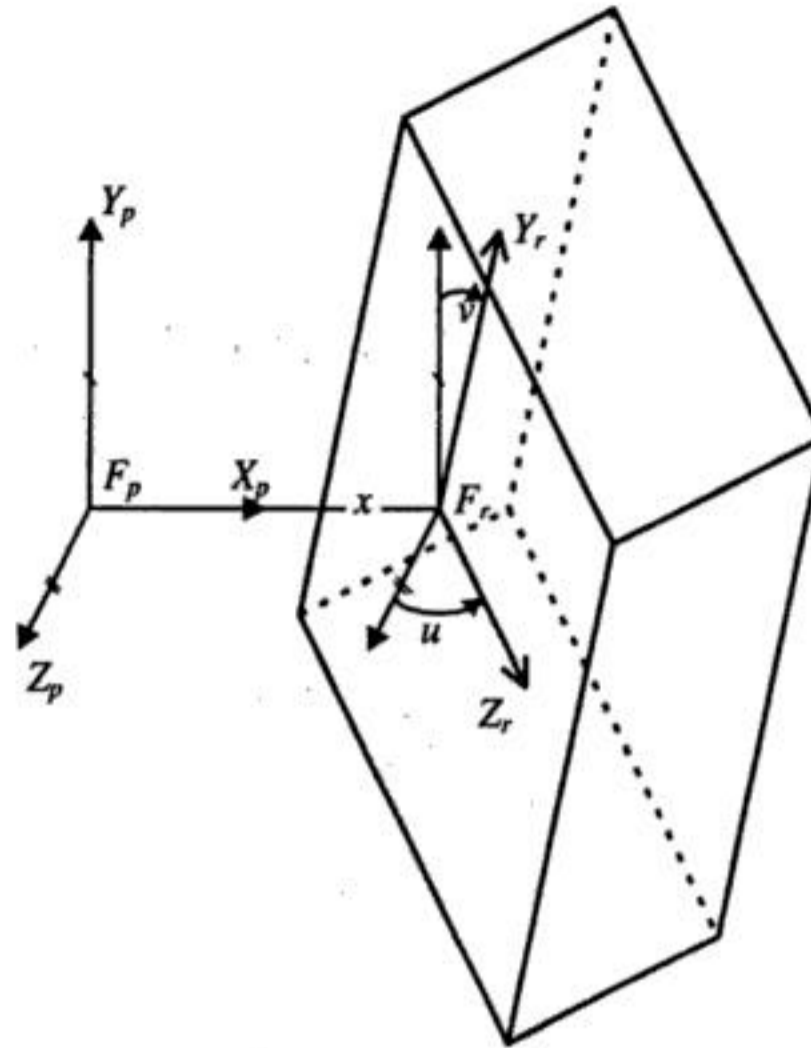


FIGURE 8.1 Proximity parameters.

8.2 Typical Sensor Characteristics

Conventionally, proximity sensors should be capable of measuring distances of up to 50 mm, and two degree-of-freedom orientations equivalent to an overall inclination of up to $\pm 30^\circ$. The intended principal application of the sensor is to act as a guide for the robot. Thus, it would be desirable to have higher sensitivity and accuracy as the gripper approaches the object, namely when both the relative orientation and the distance approach near-zero values.

The signals received by the electronic interface circuit should be processed without limiting the required operating range of the sensor. The interface circuit should also minimize the effect of interference from the surroundings. It should therefore employ solutions to reduce background-noise interference and dynamic-range limitations.

The operation of the robot should not be slowed down by the sensor. Namely, a pose of the object should normally be estimated in 1 ms to 10 ms.

8.3 Technologies for Proximity Sensing

Proximity sensors have employed various transduction media, including sound waves, magnetic fields, electric fields, and light. Presently, electro-optical techniques seem to be the most appropriate for robotic-grasping applications. Such sensors are relatively small in size, have a large range of operation, and impose almost no restrictions on the object's material. However, recently, some new ultrasonic and capacitive proximity sensors have been fabricated directly as ICs, also showing the possibility of very-small-size proximity sensors based on these technologies [7, 8].

Brief descriptions of the principles of the primary technologies used by proximity sensors are given below, with the main emphasis being on optical transducers. A survey of commercial proximity sensors capable of measuring distances can be found in [6].

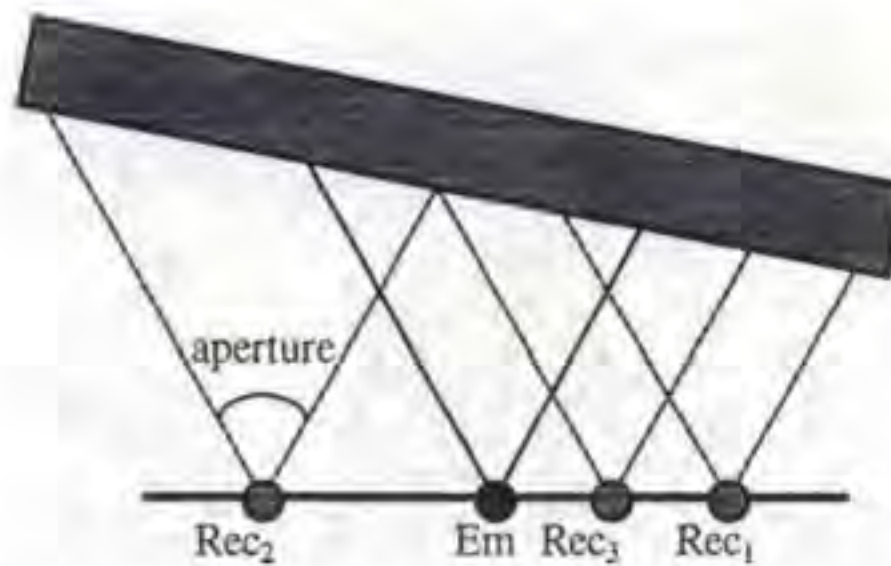


FIGURE 8.5 The basic amplitude-modulated proximity sensor configuration.

Similarly, the orientation angles u and v can be obtained by driving the LEDs as indicated in Table 8.1. For example, v can be determined by modulating LED4 and LED2 by the signals $K_1 \sin \omega t$ and $K_2 \cos \omega t$, respectively. For this case, the phase shift associated with the brightness of the object at point P is given by:

$$\phi_v = \tan^{-1} \left[\frac{K_2 (x - b \tan v)}{K_1 (x + b \tan v)} \right] \quad (8.9)$$

Note that, in order to recover v from Equation 8.9, x must be known. Accordingly, the distance x must be determined first. Correspondingly, the orientation angle (v) can be calculated from the new phase shift ϕ_v . The angle u can be calculated by modulating LED6 and LED5 with $K_1 \sin \omega t$ and $K_2 \cos \omega t$, respectively, and then determining the corresponding phase shift of the associate brightness at point P .

The pose-estimation results using this sensor were quite satisfactory and showed a good agreement between the theory and experiment.

In [10], an experimental setup of a PM distance sensor, similar to the one in [9], was reported for investigating the effect of the geometric and electronic parameters on the performance of the sensor. Optimal parameters were obtained for some targeted sensor-operation characteristics.

Amplitude Modulation

In amplitude-modulated (AM) sensors, the magnitude of the light reflected from a surface is utilized to determine the pose of the object.

AM transducers usually consist of one light source and several photodetectors (Figure 8.5). They were redesigned and optimized several times over the past decade to yield better measurement accuracy [11–14].

Many AM proximity sensors utilize optical fibers to illuminate and collect light from the surfaces of objects. The use of optical fibers, in a Y-guide configuration (Figure 8.6), facilitates the operation of sensitive low-noise circuitry in a shielded environment appropriately remote from the robot's electromagnetic interference sources.

AM transducers primarily use variations of the basic Y-guide transducer. Two important parameters can be varied in the design of Y-guides: the distance, d , between the emitting and receiving fibers (referred to hereafter as the emitter and the receiver, respectively), and the inclination angle, ϑ , of the receiver fiber with respect to the transducer's surface. The emitter is usually placed perpendicular to the transducer's surface, due to symmetry requirements, as will be explained later in this section.

The collection of a sufficient amount of reflected light requires the use of relatively wide-diameter fibers, typically having a 0.3 mm to 2 mm core size. This requirement demands the use of relatively low-grade plastic fibers. Although attenuations of up to 1 dB m^{-1} are common in such plastic fibers, this loss rate is relatively insignificant for Y-guide applications because of the short length of the cables normally used. The numerical aperture (NA) of the plastic fibers, on the other hand, is an important parameter

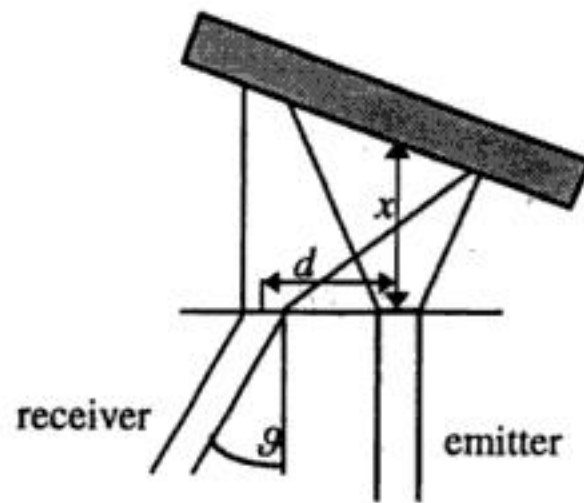


FIGURE 8.6 Y-guide transducer.

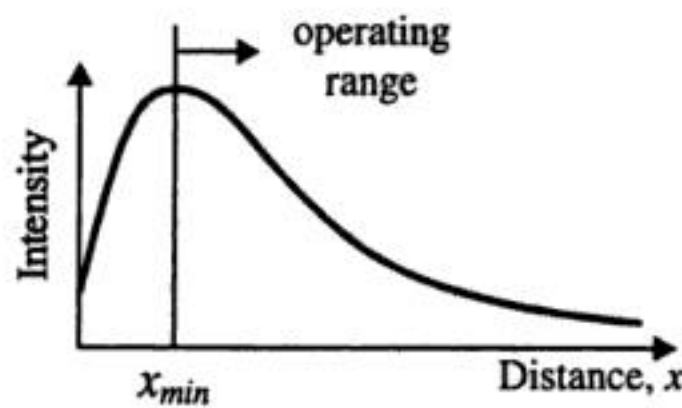


FIGURE 8.7 Y-guide response for distance measurement.

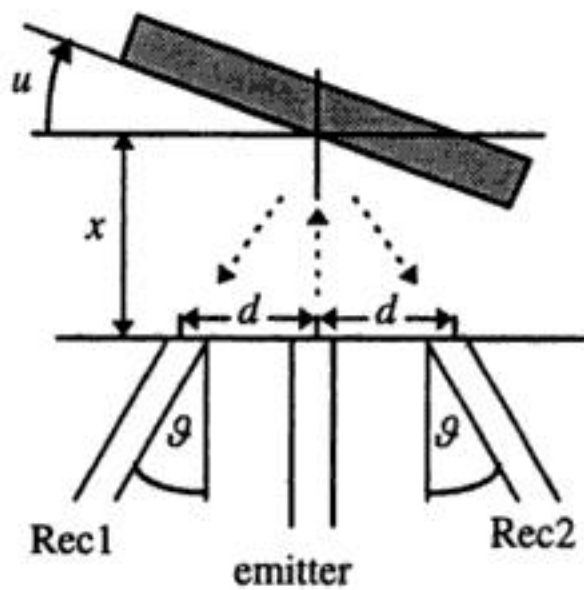


FIGURE 8.8 Typical receiver-pair constellation for orientation measurements.

in the transducer design, since lenses are rarely used in conjunction with AM-type transducers. In particular, the acceptance angle of the fiber is fixed and given by $\alpha = 2\sin^{-1} NA$.

For a Y-guide, the intensity of the light reflected from the surface is not a monotonic function of the distance. Thus, the minimum operating distance of the transducer (x_{min}) is usually limited to a value that will guarantee a monotonic response (Figure 8.7).

For the measurement of surface orientation, a symmetrical three-fiber constellation (Figure 8.8) can be used. In this Y-guide configuration, the emitter is at the center and the two receivers are positioned symmetrically on either side [12]. The light intensities detected by the receivers, for the transducer shown in Figure 8.8, are illustrated in Figure 8.9 as a function of the surface orientation.

In the usual operating range of an AM transducer, the intensity of the light at the receiver is inversely related to the distance squared. As a result, it is conceptually possible to configure a transducer such that

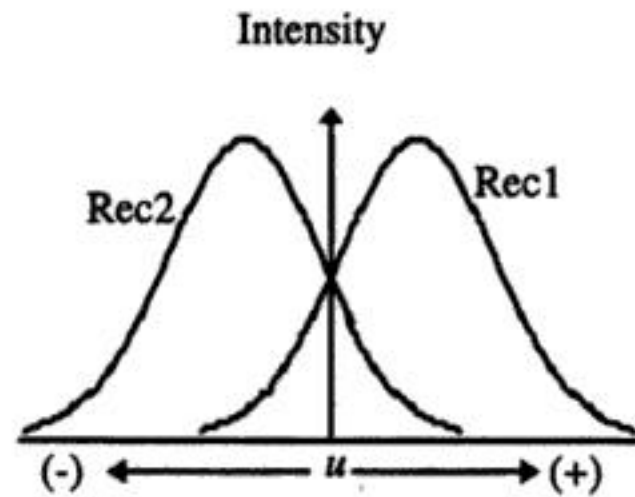


FIGURE 8.9 The light intensity detected by each receiver as a function of the surface orientation (u).

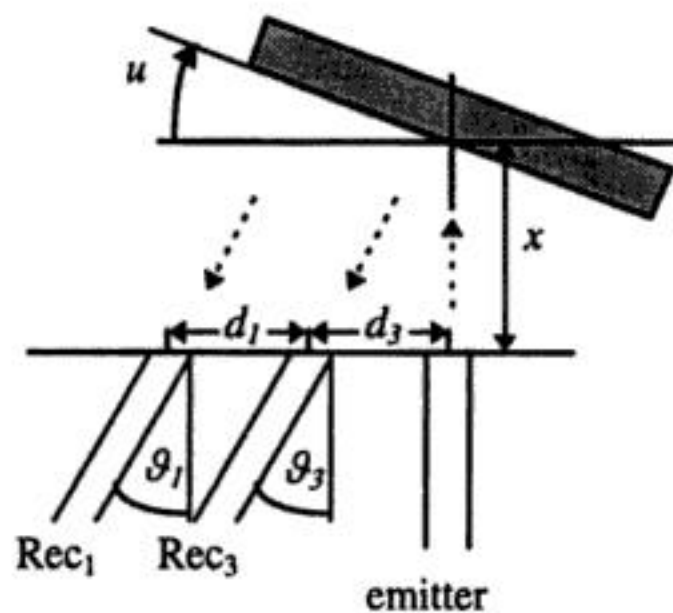


FIGURE 8.10 An asymmetrical receiver-pair constellation for distance measurements.

its sensitivity and accuracy will increase as the gripper nears the contact point, at which both the distance and the orientation of the object's surface are zero [10, 12]. However, in practice, because of the limited dynamic range of the electronic transducer interface, a trade-off exists between the desired maximum accuracy near contact, and the maximum range of operation. These and other considerations must be taken into account when establishing the geometric features of the transducer.

Another important factor to take into account in the design of an AM sensor is the need to reduce, as much as possible, the effect of the variation in the emitting power of the light source, P_0 , on the transducer's measurements. This normally leads to the employment of a pair of receivers. A normalized differential voltage (DV) estimation scheme, such as the following, is then applied to the pair of measurements:

$$DV = \frac{V_{rec1} - V_{rec2}}{V_{rec1} + V_{rec2}} \quad (8.10)$$

where V_{rec1} , V_{rec2} are the voltages measured by receivers 1 and 2. However, in order to eliminate the effect of P_0 on DV, each receiver must linearly convert the light intensity to a corresponding voltage measurement.

In order to use a DV scheme for the measurement of distance, an asymmetrical transducer configuration can be used (Figure 8.10). However, one must note that, although orientation measurements are not affected by variations in distance, distance measurements are significantly affected by the orientation of the surface, e.g., [15].

Accordingly, in using an AM proximity sensor with a DV scheme, the orientation is first approximated, and subsequently the distance is determined. The accuracies of the measured distance and orientation angle can be further improved by an iterative process.

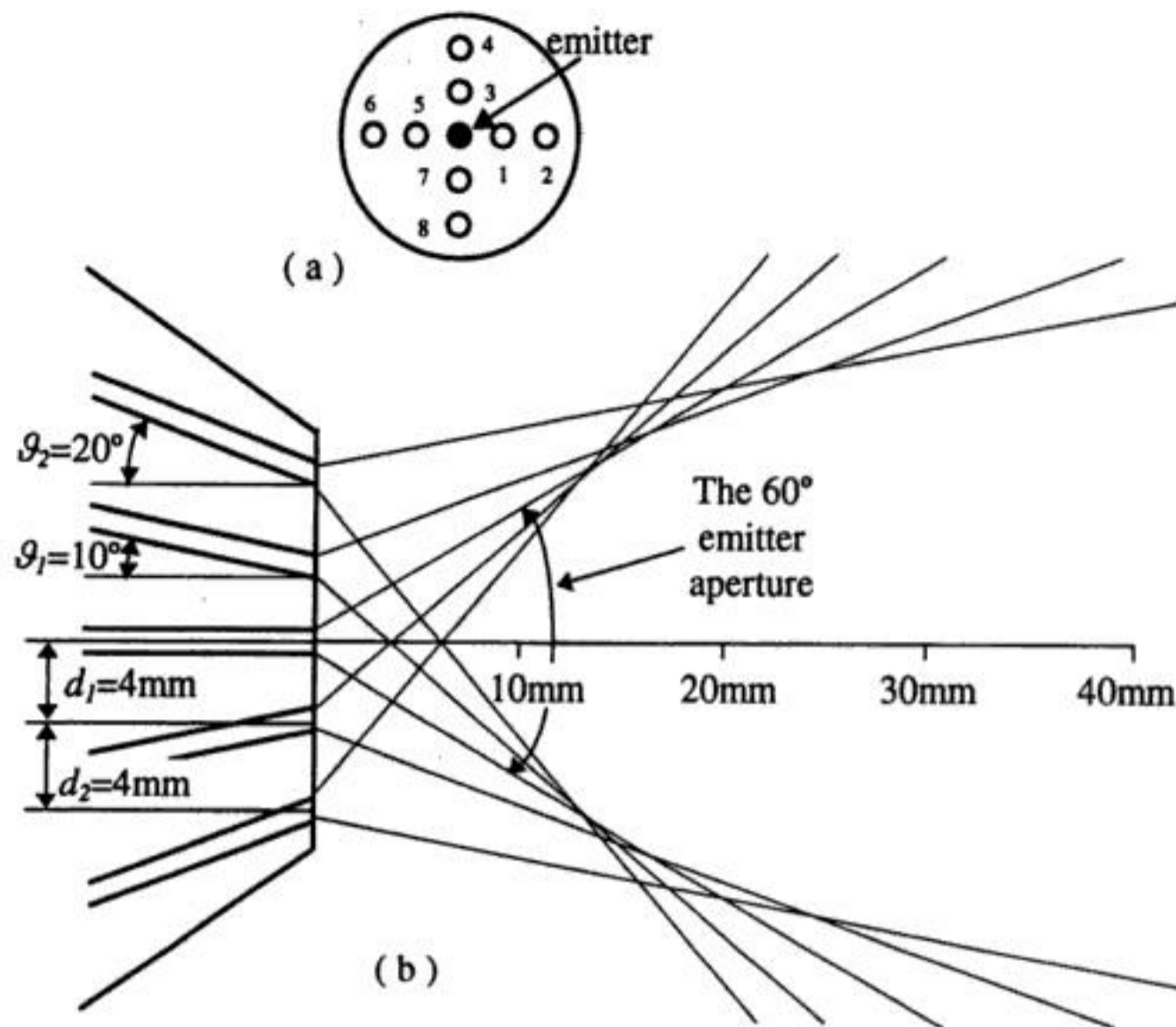


FIGURE 8.11 AM transducer design for the sensor reported in [14]: (a) top view; (b) front view.

Based on the above issues and observations, the outputs of the three receivers of the basic AM proximity sensor (Figure 8.5) can be paired for measuring both distance and orientation: the pair rec_1-rec_2 can be used for orientation measurement, while the pair rec_1-rec_3 can be used for distance measurement.

An experimental AM proximity sensor, capable of estimating the pose of an object with high accuracy, was reported in [14, 16, 17]. The transducer consists of one emitter, placed perpendicularly to the sensor head, and eight inclined receiver elements (Figure 8.11). The receivers of this transducer were paired for the specific measurements of distance, as well as of the vertical or horizontal orientation. However, the pose of the surface was determined with higher accuracy by using a polynomial fit technique (as opposed to the DV scheme described above), that provided relationships between the individual estimated parameters (x , u , and v) and all eight signals received.

The sensor presented in [14], and shown in Figure 8.11, operates in the range of 0 mm to 50 mm and $\pm 20^\circ$. It can achieve an accuracy of $6.25\ \mu\text{m}$ in distance and an accuracy of 0.02° in angular measurements in the near-contact region (0 mm to 6 mm range), using a general calibration-per-group strategy for different material groupings. This implies that the measured object's material belongs to a calibration group, which includes similar object surface characteristics; for example, machined metals. Better accuracies can be achieved using a calibration-per-surface strategy.

A similar configuration to the one shown in Figure 8.11 was reported earlier in [18] for the measurement of distances, where the orientations of each receiver pair relative to the emitter ϑ_1 and ϑ_3 were set at 10° . However, in this case, the apertures of the emitter and receiver were severely restricted by a collimating graded index (GRIN) lens. The emitter diameter was larger than that of the receivers in order to transmit more light. The measurements of the transducer were then processed in two phases: (1) the DVs of all the receiver pairs were processed independently to provide four distance estimations; and, (2) the four distance estimations were then averaged to provide a more accurate estimate, eliminating adverse effects due to variations in surface orientation.

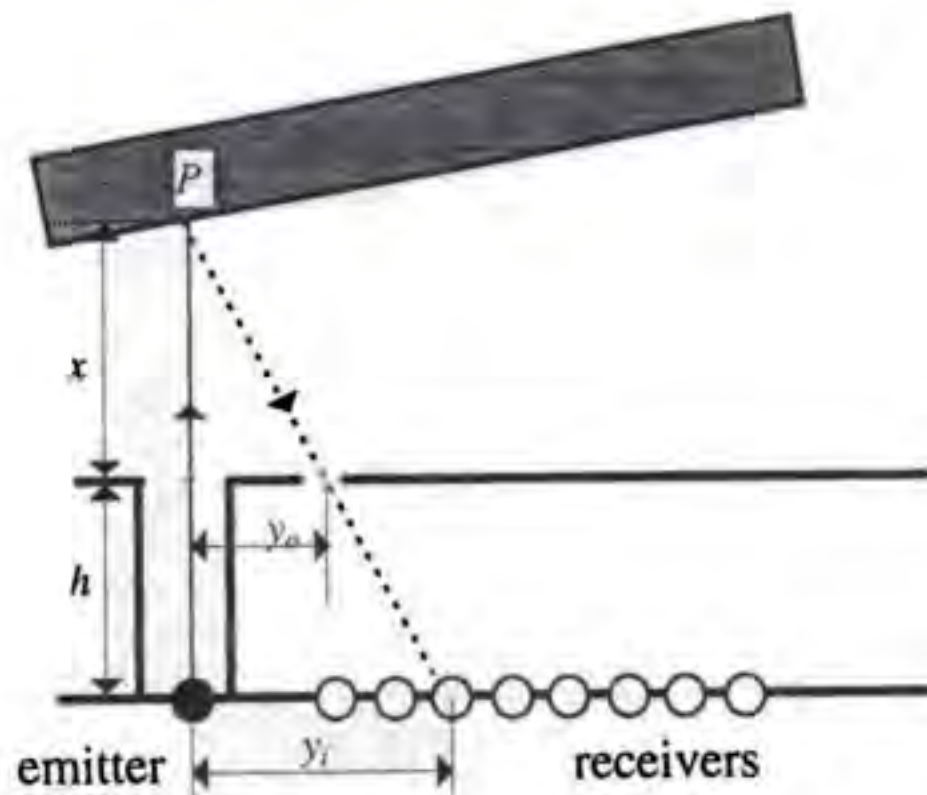


FIGURE 8.12 Basic principle of a proximity sensor for measuring distance based on triangulation.

Geometrical Techniques

Proximity sensors based on geometrical techniques determine the pose of the object by examining the geometrical attributes of the reflected and incident light beams. Two of these techniques, one based on triangulation and the other based on the Gaussian lens law, are presented here.

Figure 8.12 shows the basic configuration of a proximity transducer for measuring distance (x) based on the triangulation technique [19, 20]. The sensor head consists of a laser light source and a linear array of photodetectors (R_i , with $i = 1, 2, \dots, n$). A narrow light beam illuminates point P , and the receivers detect the reflected light from the illuminated point through a transmitting slit. The geometry of the ray trajectory provides the basic information for the estimation of the distance (x). While the light source illuminates the surface of the object, the photodetector array is scanned to detect the light path used for making the output signal maximum. The light path obtained by this scanning is called the effective light path [19]. This light path is the one indicated in Figure 8.12. The distance (x) can be determined by accurately detecting the position (y_i) and precisely measuring the dimensions (h) and (y_o),

$$\frac{x}{x+h} = \frac{y_o}{y_i} \quad (8.11)$$

or

$$x = \frac{y_o h}{y_i - y_o} \quad (8.12)$$

In [26], it is claimed that such a sensor has the following properties: (1) the influence of irregularities, reflectivity, and orientation of the object is negligible; (2) the distance measurement is not affected by illumination from the environment and luminance of the object (their influence is eliminated by comparison of two sensor signals obtained in successive on-and-off states of the light source); and (3) the sensor head is sufficiently small to be used in a robot hand.

An experimental proximity sensor configuration, based on triangulation and capable of measuring both distance and orientation, is shown in Figure 8.13 [21]. The sensor uses six infrared LEDs as light sources, an objective lens, and an area-array detector chip for detecting spot positions. The directions of

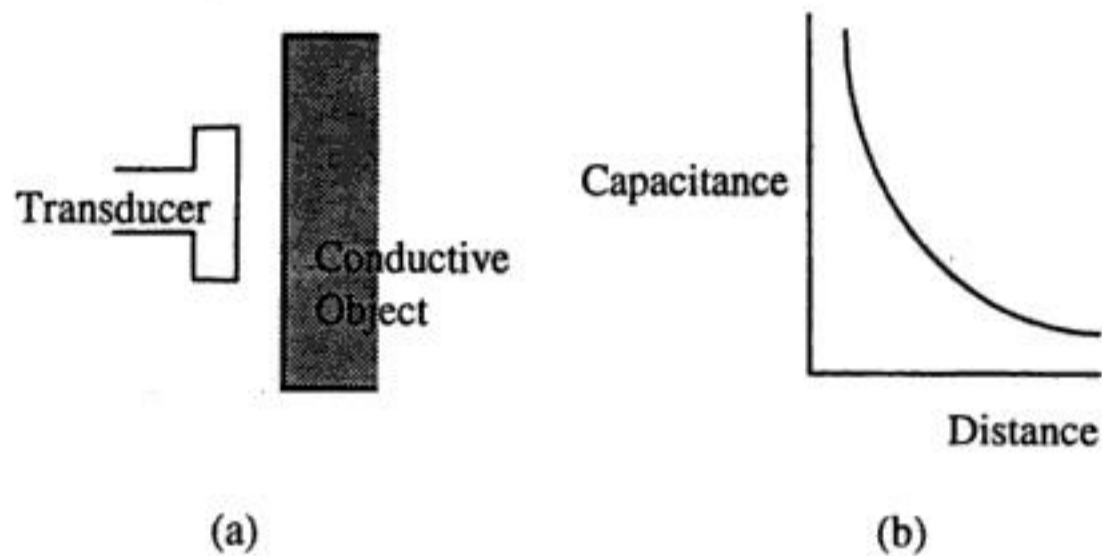


FIGURE 8.16 Capacitive proximity sensor based on the principle of parallel plates, (a) structure and (b) sensor response.

Using the Gaussian lens law has the following advantages over the triangulation principle: (1) the light source can be located at the center of the objective lens, allowing not only the sensor to be compact, but also the amount of light flux input to the lens to be maximized, and (2) the sensitivity can be optimized for a certain range of measurement distance by controlling f .

Time-of-Flight

Time-of-flight-measuring electro-optical sensors are radar-type systems. However, unlike regular radar, which transmit a pulse of radio-frequency energy, these sensors normally use a modulated light beam. The distance to the target is extracted from the measured phase shift of the reflected light. Two problems associated with such sensors are: difficulty in measuring short distances (which requires a very high modulation frequency), and the need for a mechanical scanning/switching system to get additional information (such as orientation) [5, 25].

Photothermal Effect

The photothermal effect transducer uses a strong light beam directed toward the object's surface. The distance to the object is extracted from measurements of the thermal wave generated by the light absorbed by the object. The detection scheme and signal processing are similar to those used in an AM sensor. Since the shape of the thermal wave generated at the surface is surface-texture independent, the photothermal sensor does not suffer from the surface robustness problem associated with AM sensors. However, the photothermal sensor is rather slow, and limited to highly absorbing surfaces [26].

Capacitive Sensors

Capacitive sensors generate and measure changes in an electric field caused by either a dielectric or conducting object in their proximity.

There are basically two types of capacitive proximity sensor. One type uses the principle of a parallel plate capacitor, the other uses the principle of fringing capacitances [8, 27, 28]. For the parallel plate type proximity sensor, the transducer forms one plate and the object measured forms the other plate. The structure of a parallel plate type proximity sensor and its typical response are shown in Figure 8.16 [8].

The parallel plate type proximity sensor is widely applied in industry. However, this type of a sensor has three major limitations: (1) the object being measured must be conductive; (2) the inverse gap-capacitance relationship is highly nonlinear and (3) the sensitivity drops significantly in the case of large gaps.

The second type of capacitive proximity sensor uses the principle of fringing capacitance [8]. The sensor has two "live" electrodes and the object being measured does not need to be part of the sensor system. The target object could be either conductive or nonconductive. However, the measurement of distances is affected by the type of object material. Therefore, separate calibrations must be carried out for different materials.

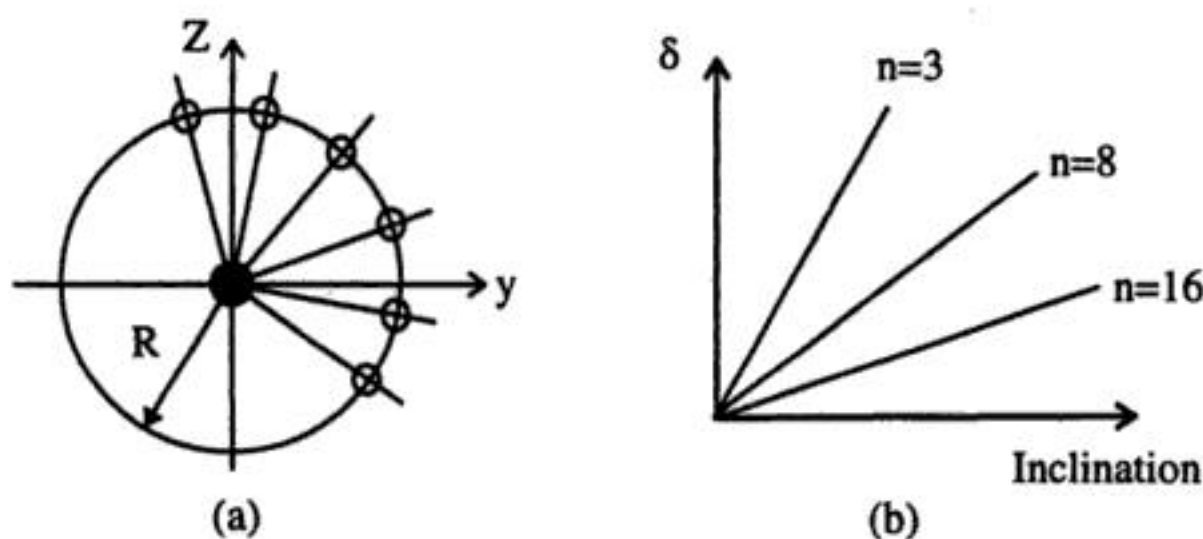


FIGURE 8.17 (a) Configuration of the ultrasonic sensor and (b) surface inclination versus δ .

In [30], an innovative capacitive microsensors was presented. Using micromachining technology, the electrode thickness can be significantly reduced and the fringing effect increased when compared with other capacitive sensors. Consequently, this sensor yields a better sensitivity. An array of such transducers can be implemented to measure the distance and orientation of an object.

Proximity capacitive sensors have the following general advantages: (1) low energy consumption and (2) simple structure. The major disadvantages, however, are that they are influenced by external signals and a calibration-per-surface technique must be carried out, since their operation directly depends on the object's material.

Ultrasonic Sensors

The basic principle underlying ultrasonic ranging sensors is the measurement of the time required for a sound wave to travel from the emitter to the object's surface and return to the detector. By using several such emitters and detectors, one can obtain information about the distance and orientation of the surface.

In [19], a novel method is proposed to measure the orientation angles of an object's surface using the phase differences of reflected echoes. Figure 8.17(a) shows the configuration of a planar sensor head with n receivers, which are equally spaced and located on a circle of radius R around the transmitter T . A linear relationship exists between the difference in lengths, δ , of two reflecting paths for an adjacent pair of receivers and the object inclinations.

In Figure 8.17 (b), the relationship between the inclination of the target surface and δ is shown. It can be observed that the measuring range of the sensor can be enlarged with an increase in the number of receivers. In [29], it is also shown that the measuring range can be enlarged by reducing R . However, it was noted that measurements carried out with a small sensor are potentially less accurate.

Experimental results using a transducer with six ($R = 30$ mm) and eight ($R = 20$ mm) receivers were reported in [29]. With the six-receiver transducer, the measuring range of the orientation angles was $\pm 15^\circ$; while for the eight-receiver transducer, the maximum measuring range was $\pm 30^\circ$. With the six-receiver transducer, the orientation angle could be determined with an accuracy of 0.5° , in the measuring range of $\pm 15^\circ$, and 0.2° when the range was restricted to $\pm 5^\circ$.

One of the major disadvantages of ultrasonic proximity sensors is that they are relatively large in size. However, implementing these sensors using micromachining could solve this problem. In [7], the generation and detection of ultrasound, for proximity sensing, was investigated using micromachined resonant membrane structures.

Magnetic Sensors

A magnetic-type sensor creates an alternating magnetic field, whose variation provides information about the object's position.

The simplest magnetic sensors are reed microswitches or Hall effect switches. However, the most commonly used sensors in robotics are based on the electromagnetic inductive principle, emphasizing

eddy current generation. The basic principle consists of creating a magnetic field using appropriate coils around a core with high permeability and an oscillator with a frequency excitation high enough to minimize the penetration of the field inside a conductive material. The main problems with magnetic sensors are their high size/range ratio and difficulty in providing reliable distance measurements in varying magnetic environments.

Acknowledgments

The authors would like to thank Martin Bonert for careful review and critique of this chapter. We also acknowledge the financial support of Natural Sciences and Engineering Research Council of Canada.

References

1. H. R. Everett, *Sensors for Mobile Robots: Theory and Application*, Natick, MA: A. K. Peters, Ltd., 1995.
2. B. Espiau, An overview of local environment sensing in robotics applications, *Sensors and Sensory Systems for Advanced Robots, NATO ASI Series, F43*, 125-151, 1988.
3. W. D. Koenigsberg, Noncontact distance sensor technology, *SPIE, Intelligent Robots*, 449, 519-531, 1988.
4. Å. Wernersson, B. Boberg, B. Nilsson, J. Nygårds, and T. Rydberg, On sensor feedback for gripping an object within prescribed posture tolerances, *IEEE, Int. Conf. on Robotics and Automation*, Nice, France, 1992, 1654-1660.
5. A. Bradshaw, Sensors for mobile robots, *Measurement and Control*, 23(2), 48-52, 1990.
6. R. Volpe and R. Ivlev, A survey and experimental evaluation of proximity sensors for space robotics, *IEEE Int. Conf. on Robotics and Automation*, 4, 3466-3473, 1994.
7. O. Brand, H. Baltes, and U. Baldenweg, Ultrasound-transducer using membrane resonators realized with bipolar IC technology, *IEEE Conf. on Micro Electro Mechanical Systems*, Oiso, Japan, 1994, 33-38.
8. R. C. Luo and Z. Chen, Modeling and implementation of an innovative micro proximity sensor using micromachining technology, *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, Yokohama, Japan, 1993, 1709-1716.
9. R. Masuda, Multifunctional optical proximity sensor using phase modulation, *J. Robotic Systems*, 3(2), 137-147, 1986.
10. O. Partaatmadja, B. Benhabib, and A.A. Goldenberg, Analysis and design of a robotic distance sensor, *J. Robotic Systems*, 10, 427-445, 1993.
11. O. Partaatmadja, B. Benhabib, A. Sun, and A. A. Goldenberg, An electrooptical orientation sensor for robotics, *IEEE Trans. on Robotics and Automation*, 8, 111-119, 1992.
12. O. Partaatmadja, B. Benhabib, E. Kaizerman, and M.Q. Dai, A two-dimensional orientation sensor, *J. Robotic Systems*, 9, 365-383, 1992.
13. P. P. L. Regtien, Accurate optical proximity detector, *IEEE Conf. on Instrumentation and Measurement Technology*, San Jose, CA, 1990, 141-143.
14. A. Bonen, R. E. Saad, K. C. Smith, and B. Benhabib, Active-sensing via a novel robotic proximity sensor, *Int. Conf. on Recent Advances in Mechatronics (ICRAM'95)*, Istanbul, 1995, 1053-1058.
15. Y. F. Li, Characteristics and signal processing of a proximity sensor, *Robotica*, 12, 335-341, 1994.
16. A. Bonen, R. E. Saad, K. C. Smith, and B. Benhabib, A novel calibration technique for electro-optical proximity sensors, *Int. Conf. on Industrial Electronics, Control and Instrumentation (IECON'95)*, Orlando, FL, 1995, 1190-1195.
17. A. Bonen, R. E. Saad, K. C. Smith, and B. Benhabib, A novel optoelectronic interface-circuit design for sensing applications, *IEEE Trans. Instrum. Meas.*, 45, 580-584, 1996.
18. H. Bukow, Fiber optic distance sensor for robotic applications, *SME Conf., Sensors*, MS86-938, Detroit, MI, 1986.
19. T. Okada, Development of an Optical Distance Sensor for Robots, *Int. J. Robotics Res.*, 1, 3-14, 1982.

20. M. A. Kujoory, Real-Time Range and Elevation Finder, *Proc. IEEE*, 72(12), 1821-1822, 1984.
21. M. Fuhrman and T. Kanade, Optical proximity sensor using multiple cones of light for measuring surface shape, *Optical Eng.*, 23, 546-553, 1984.
22. T. Okada and U. Rembold, Proximity sensor using a spiral-shaped light-emitting mechanism, *IEEE Trans. on Robotics and Automation*, 7, 798-805, 1991.
23. S. Lee, Distributed optical proximity sensor system: HexEYE, *IEEE Int. Conf. Robotics and Automation*, 2, 1567-1572, Nice, France, 1992.
24. S. Lee and J. Desai, Implementation and evaluation of HexEye: a distributed optical proximity sensor system, *Proc. IEEE Int. Conf. Robotics and Automation*, 3, 2353-2360, Nagoya, Aichi, Japan, 1995.
25. S. Shinohara et al., Compact and high precision range finder with wide dynamic range using one sensor head, *IEEE Conf. Instrumentation and Measurement Technology*, Atlanta, GA, 1991, 126-130.
26. M. Ito, K. Hane, F. Matsuda, and T. Goto, Proximity sensing technique using the photothermal effect, *J. Japan Soc. Precision Eng.*, 58, 139-144, 1992.
27. B. E. Noltingk, A novel proximity gauge, *J. Scientific Instruments, Series 2*, 2, 356-360, 1969.
28. B. E. Noltingk, A. E. T. Nye, and H. J. Turner, Theory and application of a proximity gauge using fringing capacitance, *Proc. ACTA IMEKO*, 1976, 537-549.
29. S. Nakajima and Y. Takahashi, An ultrasonic orientation sensor with distributed receivers, *Advanced Robotics*, 4, 151-168, 1990.
30. A. Moldoveanu, Inductive proximity sensors, fundamentals and standards, *Sensors*, 10(6), 11-14, 1993.

9

Distance

9.1	Basic Distinctions Between Range Measurement Techniques.....	9-1
	Contact or Noncontact • Active or Passive • Time-of-Flight, Triangulation, or Field Based • Form of Energy • Coherent or Noncoherent Detection • Ranging, Range Imaging, or Position Tracking	
9.2	Performance Limits of Ranging Systems	9-7
	Range Accuracy • Depth of Field • Maximum Range • Lateral Resolution • Rate of Acquisition	
9.3	Selected Examples of Ranging, Range Imaging, and Motion Tracking Systems	9-10
	Laser-Based Active Triangulation Ranging and Range Imaging Sensors • Laser-Based Lidar Range Imaging Sensors • Position Tracking with Active Targets	
9.4	A Sampling of Commercial Ranging, Range Imaging, and Motion Tracking Products.....	9-15

W. John Ballantyne
Spar Aerospace Ltd.

The tools and techniques of distance measurement are possibly one of humankind's longest-running inventive pursuits. The scale shown in Figure 9.1 illustrates the enormous range of distances that science and engineering have an interest in measuring [1]. This chapter concerns itself with methods to measure a relatively small segment of this range — from centimeters to kilometers. Even within this limited segment, it would hardly be possible to list, much less describe, all of the distance measurement approaches that have been devised. Nevertheless, the small sampling of technologies that are covered here should be of help to a broad range of readers.

Distance measurement, at its most basic, is concerned with determining the length of a unidimensional line joining two points in three-dimensional space. Oftentimes, a collection of distance measurements is called for, so that the shape, the orientation, or the changes in position of an object can be resolved. Therefore, one must consider not only the measurement of distances, but also their spatial and temporal distributions. The terminology “ranging” will be used in reference to systems that perform single sensor-to-target measurements, “range-imaging” for systems that collect a dense map or grid of spatially distributed range measurements, and “position tracking” for systems that record the time history of distance measurement to one or several targets.

9.1 Basic Distinctions Between Range Measurement Techniques

Range measurement devices may be classified according to some basic distinctions. Generalizations can be made based on these broad classes, thereby facilitating the process of comparison and selection. The following subsections identify the fundamental bases for classification.

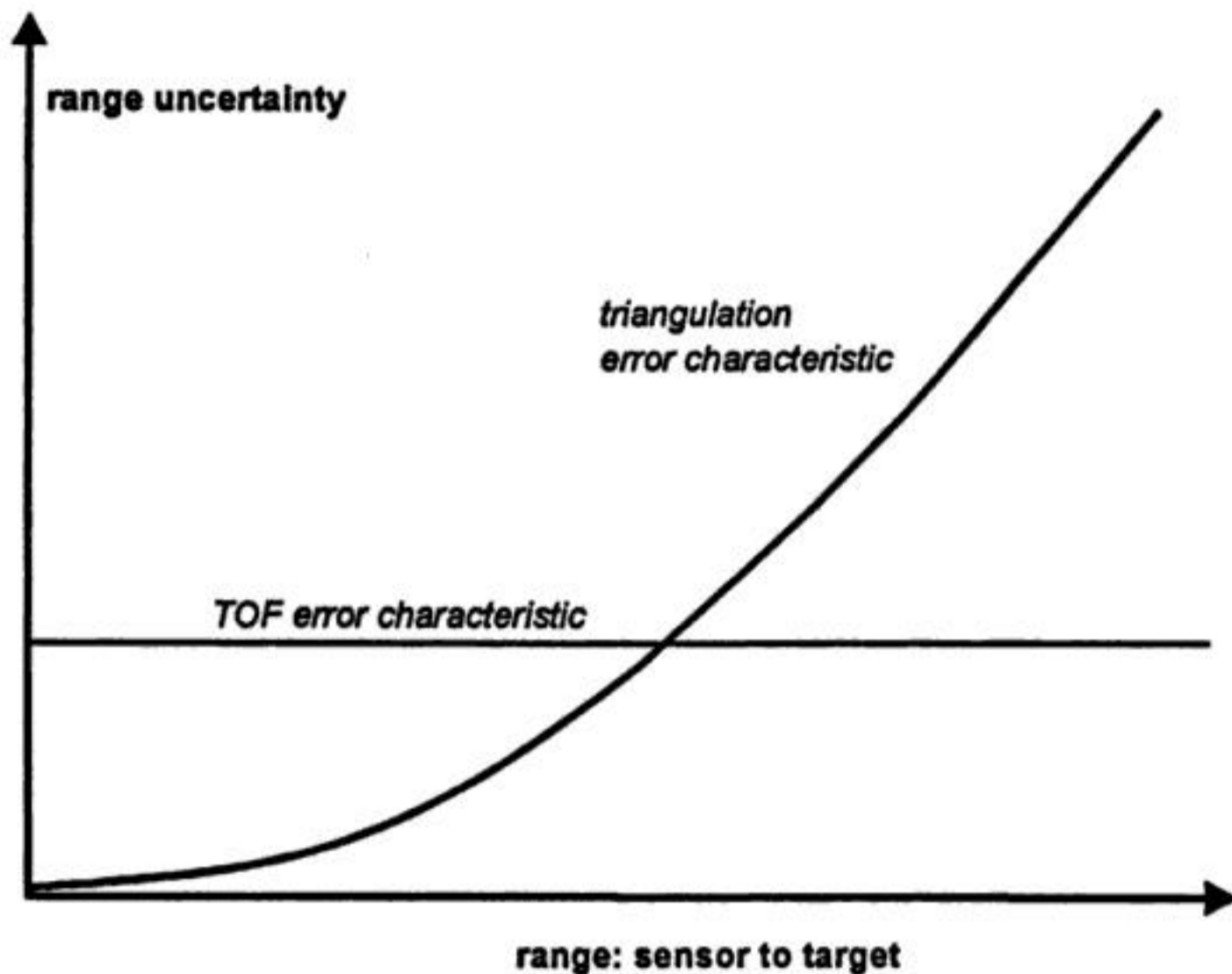


FIGURE 9.3 Time-of-flight (TOF) and active triangulation techniques tend to exhibit error characteristics related to their fundamental principles of operation. The dominant error source in TOF systems is usually the shortest measurable time interval, but this is a detection issue and is essentially independent of distance. Active triangulation systems are typically more accurate at close distances, but geometry considerations dictate that the effects of their error sources will increase with the square of distance.

An interesting distinction between field-based approaches and wave-based approaches is that the former, although they employ energy fields, do not rely on the propagation and conversion (and concomitant losses) of energy. That is, they may employ stationary fields, like those generated by a magnet or static charge. Such fields encode position information by their very shape. Sound and light, although having a wave nature, can be exploited in the same manner as stationary fields because of their distance-dependent intensity.

Field-based techniques must confront some basic issues that limit their range of application. First, the characteristics of most practically exploitable fields are typically influenced by objects or materials in the vicinity, and it is not always possible to ensure that these influences will remain constant. Second, the variation of fields through space is highly nonlinear (typically inverse square or inverse cube), implying that the sensitivity of a measurement is strongly affected by proximity to the source. Notwithstanding these concerns, devices have been developed and are available that perform very well in the situations for which they are intended [7].

Form of Energy

As discussed above, all noncontact, active ranging devices employ some form of energy. This is true whether time-of-flight, triangulation, or field-based principles apply. The following subsections describe the various forms of energy employed and some generalizations about the effectiveness of each in various situations.

Sound

Ranging systems based on sound energy are usually of the pulsed-echo TOF type and employ carrier frequencies in the so-called “ultrasonic” (beyond audible) range of frequencies. Besides being inaudible

(an obvious benefit), ultrasonic frequencies are more readily focused into directed beams and are practical to generate and detect using piezoelectric transducers. Ultrasonic signals propagate through air, but long-distance transmission is much more effective in liquids, like water, where higher density-to-viscosity ratios result in higher wave velocity and lower attenuation per unit distance. Ultrasonic ranging techniques (or SONAR, for SOund NAVigation and Ranging) were first developed for subsea applications, where sound is vastly superior to electromagnetic energy (including light) in terms of achievable underwater transmission distances [5]. Low-cost, portable sonar systems are widely used by sport fishermen as “fish finders” [6].

The frequencies typically used in sonic ranging applications are at a few tens of kilohertz to a few hundred kilohertz. A basic trade-off in the choice of ultrasonic frequency is that while high frequencies can be shaped into narrower beams, and therefore achieve higher lateral resolution, they tend to fade more quickly with distance. It may be noted that beam widths narrow enough for range imaging applications (less than 10°) are effective in a fluid medium, but attenuate too quickly to be practical in air. Interestingly, although sound energy attenuates more rapidly in air than in water, useful short-range signals can be generated in air with relatively low power levels because the much lower density of air requires smaller dynamic forces in the transducer for a given wave amplitude.

When comparing sound energy to electromagnetic energy for TOF-based techniques, one needs to remember that sound, unlike light, propagates at not only much lower speeds, but with considerably more speed variation, depending on the type and state of the carrying media. Therefore, factors like air humidity and pressure will affect the accuracy of a TOF ranging device. For underwater applications, salinity and depth influence the measurement. The lower speed of sound has a detrimental impact on the rate at which range samples can be collected. For example, a target 10 m away takes at least 60 ms to measure through an air medium. This may not seem like a long time to wait for a single sample, but it becomes an issue if the application involves multiple sampling, as in motion tracking or collision avoidance sensing.

Stationary Magnetic Fields

Stationary or pseudostationary (i.e., low frequency) magnetic fields are only used in field-based approaches. An advantage of such fields is that they are easily and cheaply produced by either a permanent magnet or electrical coil. Since stationary fields do not transmit energy, the targets cannot be passive — they must actively sense the properties of the field at their particular location. A variety of sensing technologies may be used to make measurements of the direction and intensity of a magnetic field, including flux gate, Hall effect, and magnetostrictive type magnetometers. A comprehensive list of such technologies is given in [7].

Radio Frequencies

Echo-type TOF ranging systems based on the band of the electromagnetic spectrum between approximately 1 m and 1 mm wavelength are known as RADAR (RADio Detection And Ranging). Radio waves can be used for long-distance detection in a variety of atmospheric conditions. As in the case of sound waves, there are trade-offs to be addressed in the choice of frequency. Long waves tend to propagate better over long distances, but short waves can be focused into narrow beams capable of better lateral discrimination. An interesting application of short-range radar is ground-penetrating radar, which can be used to locate and image subsurface objects [8]. Here, the frequency vs. range trade-off is particularly acute because of the need to balance reasonable imaging capability (narrow beam) with good depth penetration (long wave).

An example of a TOF one-way (active receiver) system that uses radio frequencies is the global positioning system (GPS). The distance between a receiver on land is determined by each of several orbiting satellites equipped with a transmitter and a very precise Cesium clock for synchronization. A good description of GPS and its use in vehicle navigation is available in [9].

Light Frequencies

Beyond the radio portion of the electromagnetic spectrum are the infrared, visible, and ultraviolet frequencies. These frequencies can be produced by lasers and detected by solid-state photosensitive devices and are useful for both TOF and active triangulation ranging. Echo-type TOF techniques are known as LIDAR (Light Detection And Ranging), in keeping with the terminology introduced earlier.

While light frequencies attenuate more than radio frequencies through cloud and fog, they can have very narrow beam widths, allowing superior lateral resolution and target selectivity.

Coherent or Noncoherent Detection

Echo-type TOF devices, whether sonar, radar, or lidar, can be further classified according to whether the detection approach measures time-of-flight directly (noncoherent) or exploits an inherent periodicity in the emitted energy to ascertain the flight distance (coherent).

Noncoherent techniques face the problem of timing short intervals. This is not a serious challenge in the case of sound waves, where a meter round trip corresponds to 6 ms, but is somewhat more problematic for light and radio waves, where that distance equates to only 6 ns. Accuracy of noncoherent detection typically relies on the averaging of repeated measurements.

Coherent detection is achieved by combining a portion of the emitted signal with the reflected signal to produce a third signal indicating the amount of phase delay. The signals are continuous wave (CW) as opposed to pulsed. Coherent detection techniques are classified as amplitude modulated (AMCW) or frequency modulated (FMCW).

A basic issue with coherent detection techniques is the inability to distinguish between integral multiples of the basic modulation wavelength. Any coherent detection system must employ techniques to resolve the so-called "ambiguity interval." Noncoherent techniques do not face this problem.

Ranging, Range Imaging, or Position Tracking

Ranging devices are typically pointed toward a target to produce a single range reading. A common example of simple ranging is the feedback sensor used in auto-focus cameras. There are many active ranging devices currently available based on TOF (i.e., radar, sonar, lidar) and active triangulation principles.

Range imaging devices use the same principles as ranging devices, except that they include some form of scanning that is employed to generate an array of spatially distributed range samples. Sometimes, the scanning action is accomplished by means intrinsic to the sensor (e.g., spinning and nodding mirrors, or phased-array antenna) so that the reference location remains fixed. In this case, the data are recorded in the polar form (range, elevation, azimuth) as shown in Figure 9.4. In other cases, the sensor might scan on only one axis internally while the second scan dimension is realized by moving the sensor location through some set pattern. It is not uncommon to record the "intensity" or return energy associated with a range sample as well. The intensity map may be presented as a "gray scale" image and, like a black and white photograph, often contains additional information useful in interpreting a scene. Range images can be used to produce three-dimensional graphic representations of scenes and objects. A common use of range imaging is aerial terrain mapping.

Position tracking devices are used to measure the change in an object's position and orientation over time. Basic issues in position tracking are the acquiring of, and locking on to, specific target points. These issues can be avoided by employing active targets, and most systems available today are of this type.

9.2 Performance Limits of Ranging Systems

The performance characteristics of available ranging systems vary widely, as do the requirements of the applications for which they are designed. The following subsections review the most basic performance categories and the technical issues of performance limits.

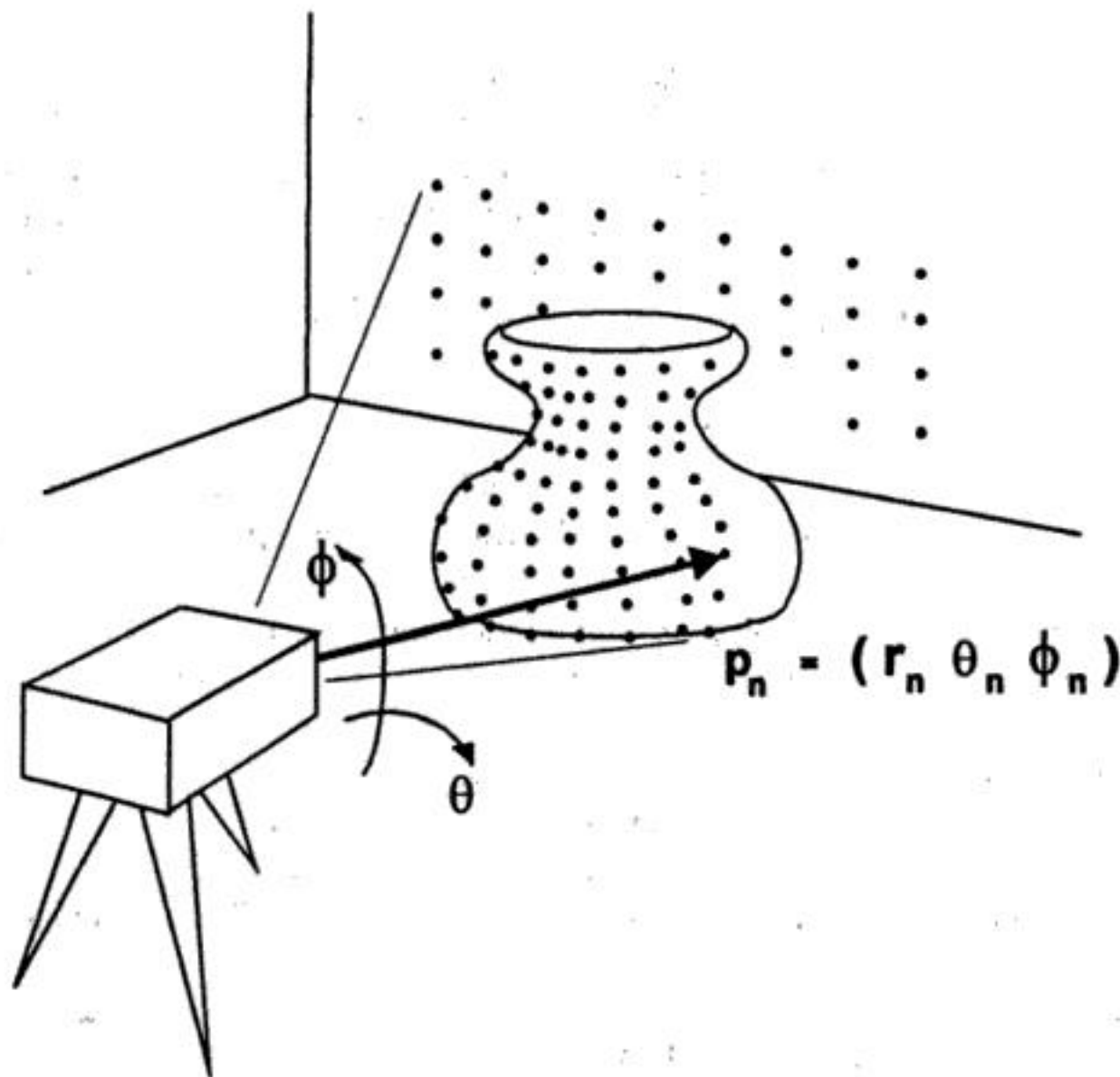


FIGURE 9.4 Range images are typically an array of individual range values sampled while changing the pointing direction (e.g., azimuth and elevation angles) of a ranging device. A digital range image of the polar form shown can be readily transformed into rectangular coordinates if required.

Range Accuracy

As illustrated in Figure 9.3, TOF and active triangulation techniques differ fundamentally in their error vs. distance characteristics. Currently available systems based on active triangulation achieve better repeatability and accuracy in the less than 1 m range than do TOF systems, but are seldom used at distances of several meters. Hymarc Ltd. and Perceptron Inc. each offer laser triangulation systems with 3σ accuracy of 25 mm and 50 mm, respectively [10, 11].

In principle, TOF systems could achieve accuracy rivaling active triangulation, but the most promising detection technique — a variation of laser interferometry, which solves the ambiguity interval problem [12] — has yet to make its commercial debut.

Depth of Field

Depth of field refers to the interval of distance through which a stationary reference ranging system can measure without resorting to a change in configuration. Large depth of field is often an important characteristic in practical applications. For example, if the distance to the target is poorly known a priori, then a large depth of field is desirable.

Passive optical triangulation approaches like stereography and photogrammetry tend to have restricted depth of field because they rely on camera-type imaging, which is inherently limited by depth of focus. Timed-interval TOF systems have excellent depth of field because they do not rely strongly on optical imaging except to concentrate the collected return energy on the detector. Some active triangulation systems do rely on optical imaging of the projected laser spot, but the design employed by Hymarc Ltd. regains a large depth of field by tilting the detector array with respect to the lens plane [13].

Maximum Range

Any active ranging, range imaging, or position tracking system has a practical maximum distance that it can measure. This is because the controlled energy, whether propagated as a wave or established as a field, must spread before reaching the detector. The spreading inevitably increases with distance and all detectors, no matter what form of energy they measure, require a certain minimum amount to exceed their inherent "noise floor."

The "classical radar range equation" is introduced in many texts on radar (e.g., [14]). Jelalian [15] points out that the equation is equally applicable to lidar, which, after all, just employs a higher frequency version of electromagnetic wave. In fact, the same idea applies to sonar and to active triangulation systems as well. The equation computes the power of the received signal as:

$$P_R = P_T G_T / 4\pi R^2 \times \rho A / 4\pi R^2 \times \pi D^2 / 4 \times \eta_{\text{atm}} \eta_{\text{sys}} \quad (9.2)$$

where P_R = power at the receiver
 P_T = power transmitted
 G_T = transmitter gain
 R = range to target
 ρ = reflectivity of target
 A = effective area of target
 D = diameter of collecting aperture
 η_{atm} = atmospheric transmission coefficient
 η_{sys} = system transmission coefficient

Equation 9.2 applies when the target area is smaller than the footprint of the incident beam, which is often the case for radar and sonar ranging. However, in the case of laser-based systems, the relatively narrow beam usually means that the laser spot is small compared to the target. For a transmitted beam that spreads with a solid angle θ_T , the illuminated patch area is:

$$\sigma_{\text{spot}} = \pi R^2 \theta_T^2 \quad (9.3)$$

The definition of transmitter gain is based on the notion of the solid angle beam width as compared to an omnidirectional transmitter

$$G_T = 4\pi / \theta_T^2 \quad (9.4)$$

One can substitute for Equation 9.4 for G_T and Equation 9.3 for the variable s in Equation 9.2 to produce the range equation for a small spot size.

$$P_R = P_T / R^2 \times \rho / 4 \times \pi D^2 / 4 \times \eta_{\text{atm}} \eta_{\text{sys}} \quad (9.5)$$

The importance of this equation is primarily in the $1/R^2$ dependence. Any ranging system that works by bouncing energy off a diffuse reflective target encounters severe signal attenuation with increasing distance. Given a detector with a fixed noise floor, the only ways to improve maximum range are to increase the transmitted power or the collecting area. In practice, there are design constraints that limit both of these measures. For example, laser power must sometimes be limited for eye-safety considerations, and increased collecting area can imply a proportional increase in sensor packaging volume.

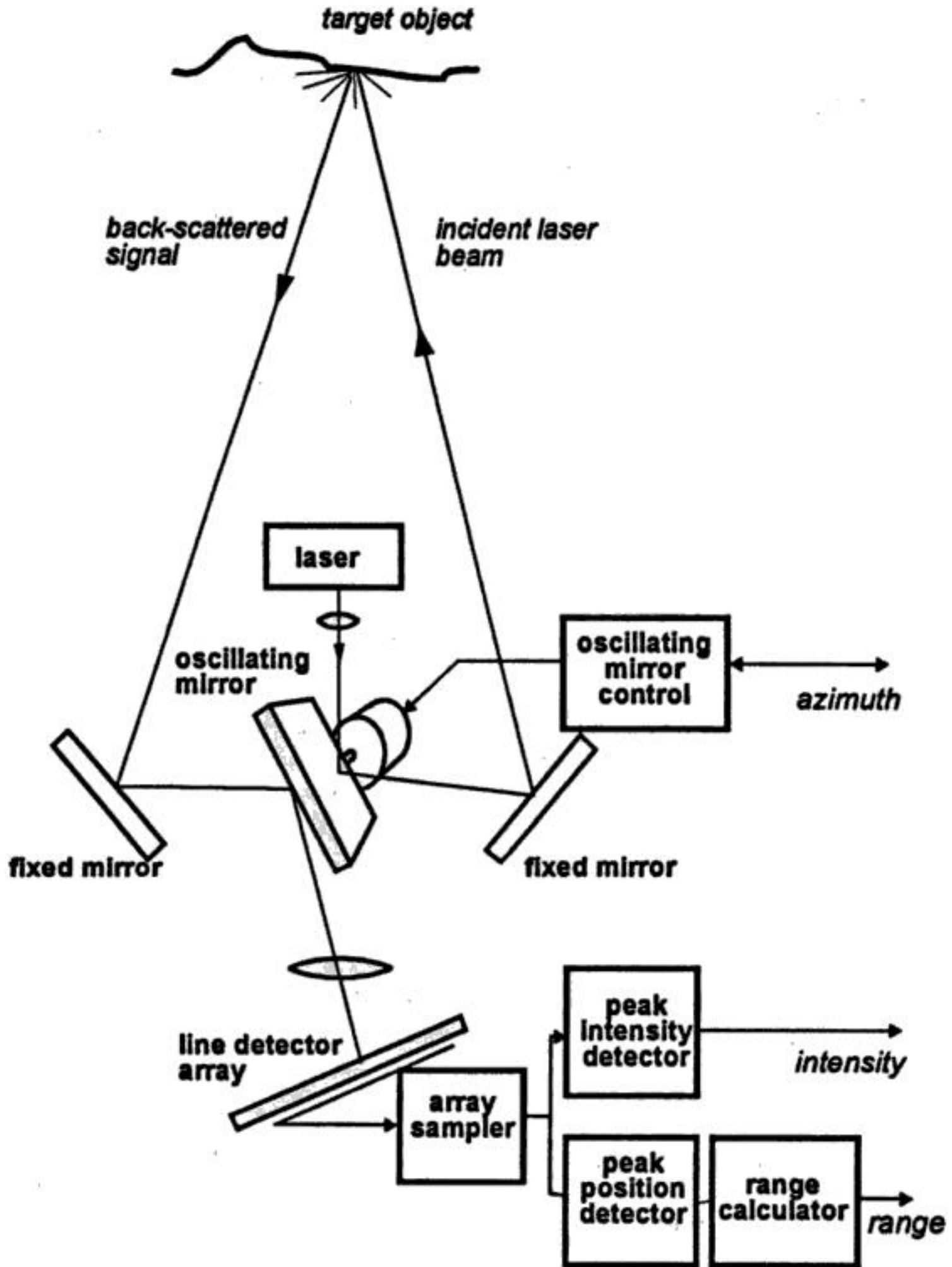


FIGURE 9.6 The Hymarc laser triangulation line scanner uses the synchronized scanning principle. Both sides of an oscillating mirror are used to sweep both the projected beam and the axis of detection over the target. The detector array is tilted to the lens plane to maximize the depth of focus.

Laser-Based Lidar Range Imaging Sensors

AM Lidar (Phase-Based Detection)

Perceptron Inc. also offers a scanning lidar under the name "LASAR" that can produce high-resolution range images through a large measurement volume. The device uses a near-infrared laser that is projected through a collimating telescope to form a spot on the first surface encountered. The spot is swept over

Position Tracking with Active Targets

Active target approaches are not convenient in some applications, but they are an excellent way to track the changing positions of several target points simultaneously. Active targets are a way of getting around the "correspondence problem" mentioned earlier. The two systems introduced here are interesting to compare. One employs light energy and triangulation; the other uses a magnetic field-based approach. They are both used for real-time tracking and recording of human kinetics, robotics, and other moving objects.

Active Target Triangulation

The "OPTOTRAK" system offered by Northern Digital Ltd. [19] uses infrared light emitting diodes (LEDs) as targets. The LEDs are multiplexed so that only one at a time can be seen by the camera system, avoiding the correspondence problem. The unique form of stereo ranging is based on three line detectors with lenses that transform the point source LED illumination into a focused line. The simplified triangulating geometry is shown in Figure 9.8. It may be shown from this geometry that the target position (x_p, y_p, z_p) can be determined from the detector outputs u_{left} , u_{right} , and v as follows:

$$x_p = b(u_{\text{right}} + u_{\text{left}}) / 2(u_{\text{right}} - u_{\text{left}}) \quad (9.8)$$

$$y_p = bv / (u_{\text{right}} - u_{\text{left}}) \quad (9.9)$$

$$z_p = fb / (u_{\text{right}} - u_{\text{left}}) \quad (9.10)$$

where f and b are the lens-to-detector distance and the baseline separation respectively. In practice, the image space to object space mapping is much more complicated than Equations 9.8 to 9.10, and involves a camera model with more than 60 parameters that are determined through a calibration process.

OPTOTRACK offers high sampling rate, large measurement volume, and high accuracy compared to many other position tracking systems.

Magnetic Position Tracking

A position/orientation tracking sensor based on a three-axis magnetic dipole transmitter and a three-axis magnetic loop detector has been developed by Polhemus Inc. [20]. The transmitted fields are alternating current for ease of detection (i.e., transformer coupled) and time-multiplexed so that the field due to each axis can be distinguished from the others. Distance between transmitter and detector is determined by exploiting the $1/R^3$ relationship between field strength and distance from the source. Orientation of the detector is determined by exploiting the directionality of magnetic fields and the direction sensitivity of loop detectors.

An issue with respect to the use of ac fields is the distortions in field shape that occur if metal objects are present, and the consequent effect on sensor accuracy. These distortions result from eddy currents in the conducting metal. Ascension Technology Corp. has developed a variation on the Polhemus sensor based on dc magnetic fields. The switching transient due to time-multiplexing does produce an eddy current effect, but it is allowed to die out before measurement is made. Details of the dc technique are available in [21].

An important difference between optical and magnetic tracking technologies is that the former require an unbroken line of sight to the targets while the latter do not. This gives magnetic trackers an advantage in some applications. On the other hand, the $1/R^3$ field distribution characteristic of magnetic tracking

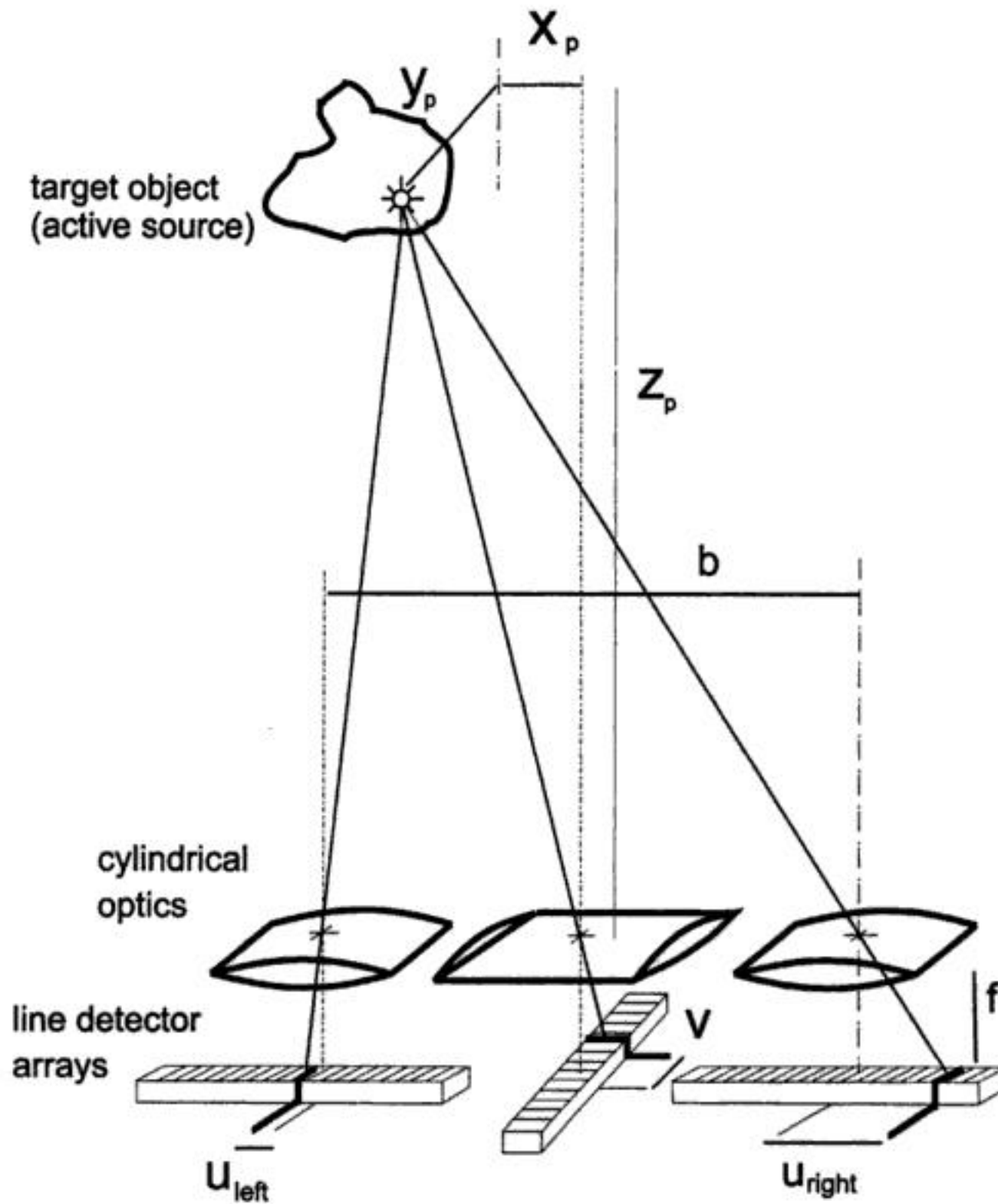


FIGURE 9.8 The OPTOTRAK position tracking system employs a novel arrangement of cylindrical optics and one-dimensional detectors to triangulate the 3-D position of an infrared LED target. Up to 255 individual multiplexed targets can be tracked by the system.

implies an extreme sensitivity loss with distance, whereas optical triangulation has a more benign $1/R$ characteristic. This, to some extent, explains why the volume of measurement and accuracy of optical triangulation systems is generally much better than for magnetic systems.

9.4 A Sampling of Commercial Ranging, Range Imaging, and Motion Tracking Products

Table 9.1 contains information collected from vendor literature. Be advised when comparing specifications that test conditions, standards, and interpretations can vary significantly. The specifications, therefore, should serve only as a rough guide.

TABLE 9.1 Ranging, Range Imaging, and Position Tracking Products and Vendors

Class	Trade Name	Principle	Features	Contact
Ranging (contact)	MicroScribe-3DX	Instrumented arm	50 in. spherical work volume, 0.3 mm accuracy	Immersion Corp. (408) 467-1900, info@immerse.com
Ranging (noncontact)	LASERVISION	TOF, laser	50 m range, 4.9 mm accuracy @ 15 m, integrated electronic level	ZIRCON Corp., (408) 866-8600
Range-Imaging (line scan)	HYSCAN	Active triangulation laser	40 mm depth of field, 70 mm swath, 0.025 mm accuracy, 10,000 points/s	Hymarc Ltd., (613) 727-1584, info@hymarc.com
Range-Imaging (line scan)	TriCam	Active triangulation laser	120 mm depth of field, 60 mm swath, 0.05 mm accuracy	Perceptron Inc., (810) 478-7710, inquiry@perceptron.com
Range-Imaging (line scan)	ALTM 1020	TOF laser time-interval	330-1000 m range, 15 cm accuracy, 20° swath	Optech Inc., (416) 661-5904
Range-Imaging (area scan)	Rangecam 7000	Laser or strobe triangulation	uses standard CCD camera and light plane projector	Range Vision Inc. (604) 473-9411
Range-Imaging (area scan)	LASAR	TOF, AM Lidar	2-40 m range, 60 × 70° max field of view, 360,000 samples/s	Perceptron Inc., (810) 478-7710
Position Tracking	OPTOTRAK	Active target triangulation	up to 255 targets, submillimeter accuracy, 5000 3 DoF samples/s	Northern Digital Inc., (519) 884-5142
Position Tracking	Flock of Birds	Magnetic field based	up to 30 position/orientation targets, approx. 10 mm accuracy, 144 6-DoF samples/s	Ascension Technology Corp. (802) 860-6440

References

1. R. Resnick and D. Halliday, *Physics (Part 1)*. New York: John Wiley & Sons, 1966. 4.
2. P. J. Besl, Range imaging sensors. General Motors Research Publication, GMR-6090, General Motors Research Laboratories, Warren, MI, March, 1988.
3. R. Resnick and D. Halliday, *Physics (Part 1)*. New York: John Wiley & Sons, 1966. 3.
4. D. F. McAllister (ed.), *Stereo Computer Graphics and Other True 3D Technologies*, Princeton, NJ: Princeton University Press, 1993. Ch. 4.
5. L. E. Kinsler and A. R. Frey, *Fundamentals of Acoustics, 2nd. ed.*, New York: John Wiley & Sons, 1962, Chs. 9, 15.
6. W. Diedrich, Foundations of reading sonar, *The In-Fisherman*, April-May, 42-56, 1996.
7. E. B. Blood, Device for quantitatively measuring the relative position and orientation of two bodies in the presence of metals utilizing direct current magnetic fields, U.S. Patent 4,945,305, Jul. 31, 1990.
8. W. J. Steinway and C. R. Barrett, Development status of a stepped-frequency ground penetrating radar, in *Underground and Obscured Object Imaging and Detection, SPIE Proceedings*, Vol. 1942, Orlando, FL, April 1993, 34-43.
9. J. Borenstein, H. R. Everett, and L. Feng, Where am I? Sensors and Methods for Autonomous Mobile Robot Positioning, 1995 Edition. University of Michigan report for the United States Dept. of Energy Robotics Technology Development Program, Ann Arbor, MI, 1995. Ch. 3.
10. Hymarc Ltd., 1995. Product Information, Hyscan 3D Laser Digitizing Systems. Ottawa, Ontario, Canada.
11. Perceptron Inc., 1995. Product Information, TriCam Non-Contact Measurement Solutions. Farmington Hills, MI.
12. F. E. Goodwin, Frequency Modulated Laser Radar, U.S. Patent 4,830,486, May 16, 1989.
13. F. Blais, M. Rioux, and J.-A. Beraldin, Practical considerations for a design of a high precision 3D laser scanner system, *SPIE Vol. 959, Optomechanical and Electro-Optical Design of Industrial Systems*, 1988.

14. D. K. Barton, *Radar System Analysis*, Englewood Cliffs, NJ: Prentice-Hall, 1964. Ch. 4.
15. A. V. Jelalian, *Laser Radar Systems*, Artech House, 1992. Ch. 1.
16. E. S. Cameron, R. P. Srumski, and J. K. West, Lidar Scanning System, U.S. Patent 5,006,721, Apr. 9, 1991.
17. Acuity Research Inc., 1995. Product Information, Accurange 4000. Menlo Park, CA.
18. R. R. Clark, Scanning rangefinder with range to frequency conversion, U.S. Patent 5,309,212, May 3, 1994.
19. Northern Digital Inc., 1990. Product Literature, OPTOTRACK 3D Motion Measurement System, Waterloo, Ontario, Canada.
20. F. H. Raab, E. B. Blood, T. O. Steiner, and H. R. Jones, Magnetic position and orientation tracking system, *IEEE Trans. Aerospace Electronic Systems*, Vol. AES-15, No. 5, September 1979.
21. E. B. Blood, Device for quantitatively measuring the relative position and orientation of two bodies in the presence of metals utilizing direct current magnetic fields, U.S. Patent 4,945,305, July 31, 1990.

10

Position, Location, Altitude Measurement

Dimitris E. Manolakis
Technological Education Institute

Mark Stedham
University of Alabama in Huntsville

Partha P. Banerjee
University of Alabama in Huntsville

Seiji Nishifuji
Yamaguchi University

Shogo Tanaka
Yamaguchi University

Halit Eren
Curtin University of Technology

C.C. Fung
Curtin University of Technology

Jacob Fraden
Advanced Monitors Corporation

10.1	Altitude Measurement	10-1
	Ground-Based Height Estimation • Onboard Derived Height Estimation • Estimation of Vertical Position with the Global Positioning System (GPS) • Special Topics	
10.2	Attitude Measurement	10-17
	Attitude Sensors for Ships, Aircraft, and Crane Lifters • Attitude Sensors for Spacecraft Applications • Automatic On-Line Attitude Measurement for Ships and Crane Lifters • Aircraft Attitude Determination • Spacecraft Attitude Determination • PALADS	
10.3	Inertial Navigation	10-34
	The Principles • Errors and Stabilization • Vehicular Inertial Navigation • Aircraft • Underwater • Robotics	
10.4	Satellite Navigation and Radiolocation	10-48
	Accuracy of Electronic Fix • Radionavigation Systems • Satellite Relay Systems • Transponders • Global Satellite Navigation Systems	
10.5	Occupancy Detection	10-62
	Ultrasonic Sensors • Microwave Motion Detectors • Micropower Impulse Radar • Capacitive Occupancy Detectors • Triboelectric Detectors • Optoelectric Motion Detectors	

10.1 Altitude Measurement

Dimitris E. Manolakis

Accurate monitoring of aircraft cruising height is required in order to reduce vertical separation to a minimum standard. Interest here focuses on the measurement of the distance between aircraft level and the sea surface level. This distance can be estimated onboard via barometric altimeters or it can be measured — either onboard or in ground stations — via electronic radio wave systems. The indication of the first equipment is referred to as pressure altitude, or simply altitude, whereas that of the second category is referred to as geometric height or simply height.

The altitude information at air traffic control (ATC) centers is based on pressure altitude measurement that the aircraft transponder system sends after it receives an appropriate interrogation — known as mode C interrogation — transmitted by a secondary surveillance radar. Actually, the altitude information is an atmospheric pressure measurement transformed to altitude indication through a formula expressing the pressure/altitude relationship. When a flight level is cleared for an aircraft, it actually means that the pilot must keep flying on an isobaric surface. However, the altimetry system may present systematic errors (biases) that are different for each airplane, and that significantly affect safety. Thus, the altimetry

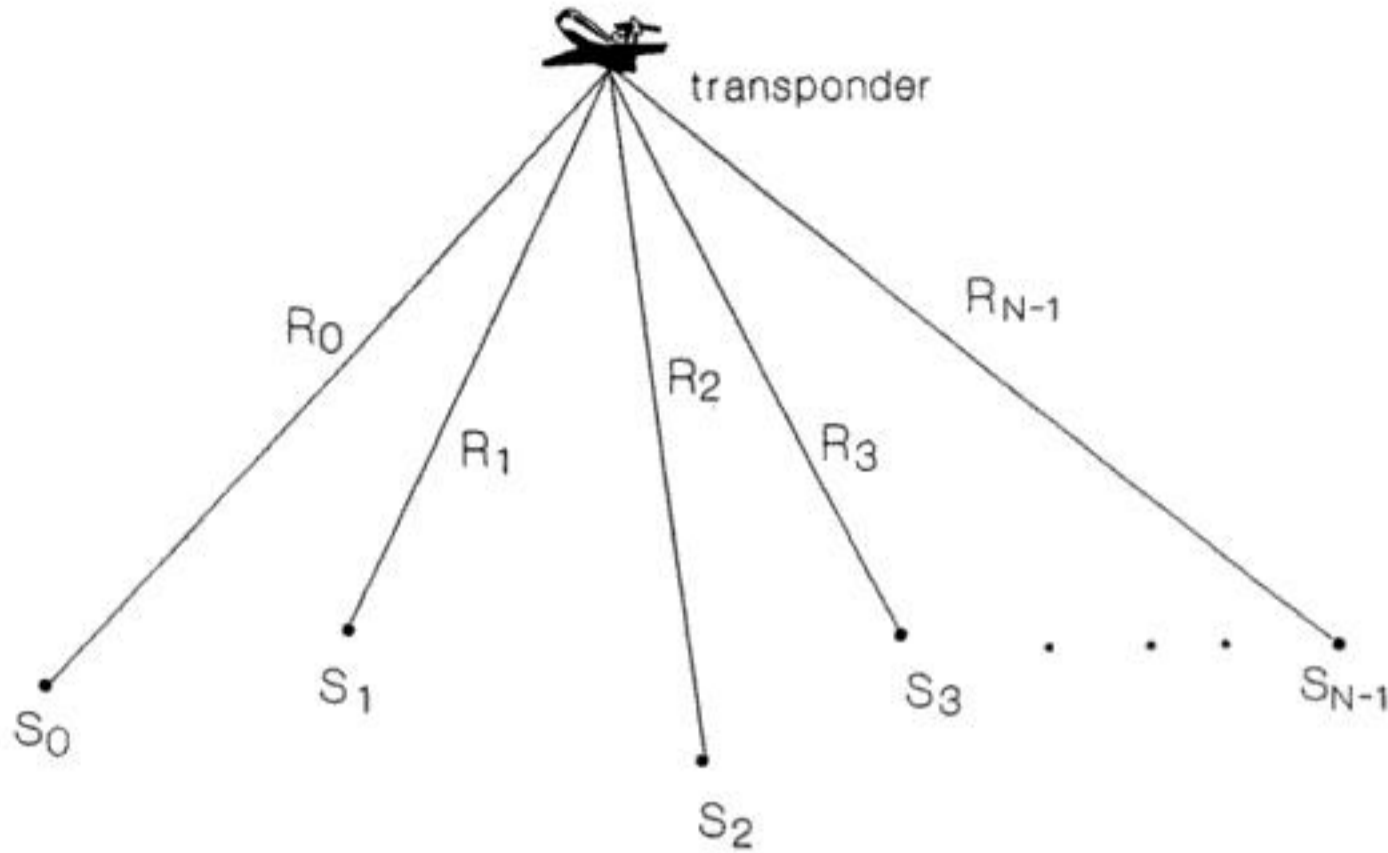


FIGURE 10.5 Typical configuration of the height estimation systems that are based on N SSR stations. One of the stations is active, i.e., it both transmits the interrogations and receives the replies, whereas the other stations are receivers only.

$$T_i = T_s + \frac{R_i}{c} \quad i = 0, 1, 2, 3 \quad (10.23)$$

$$R_i^2 = (x - x_i)^2 + (y - y_i)^2 + (z - z_i)^2, \quad i = 0, 1, 2, 3 \quad (10.24)$$

The above system of eight equations can be solved for the unknown quantities. The unknown quantities used by Rice are $(R_0, R_1, R_2, R_3, x, y, z, T_s)^T$. However, an equivalent approach is to substitute for R_i in Equation 10.23, which becomes:

$$T_i = T_s + \frac{1}{c} \sqrt{(x - x_i)^2 + (y - y_i)^2 + (z - z_i)^2} = f_i(\mathbf{q}) \quad i = 0, 1, 2, 3 \quad (10.25)$$

where $\mathbf{q} = [x, y, z, T_s]^T$ is the unknown vector. Thus, there are four nonlinear equations to be solved for \mathbf{q} . One suitable method, for example, is the Newton-Raphson method, which iteratively approximates the solution via the following formula:

$$\mathbf{q}_{k+1} = \mathbf{q}_k + \mathbf{F}(\mathbf{q}_k)^{-1} (\mathbf{T} - \mathbf{f}(\mathbf{q}_k)) \quad (10.26)$$

where $\mathbf{T} = [T_0, T_1, T_2, T_3]^T$ is the measurement vector and \mathbf{F} is the Jacobian matrix:

$$\mathbf{F} = \frac{\partial \mathbf{f}}{\partial \mathbf{q}} = \begin{bmatrix} \frac{\partial f_0}{\partial x} & \frac{\partial f_0}{\partial y} & \frac{\partial f_0}{\partial z} & \frac{\partial f_0}{\partial T_s} \\ \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} & \frac{\partial f_1}{\partial z} & \frac{\partial f_1}{\partial T_s} \\ \frac{\partial f_2}{\partial x} & \frac{\partial f_2}{\partial y} & \frac{\partial f_2}{\partial z} & \frac{\partial f_2}{\partial T_s} \\ \frac{\partial f_3}{\partial x} & \frac{\partial f_3}{\partial y} & \frac{\partial f_3}{\partial z} & \frac{\partial f_3}{\partial T_s} \end{bmatrix} \quad (10.27)$$

Notice that the time of interrogation transmission, as well as the transponder nominal delay time, are not involved in the measurements. The measured quantities are only the TOAs at the station sites. Thus, the height estimate is not affected by any transponder bias.

The theoretical and experimental research has been conducted at the GEC Marconi Research Center, U.K. The optimum station arrangement is to locate the three of them equispaced on a circle and the fourth in the middle. The typical circle radius is 35 km. The same magnitude holds for the measurement range. The method may be implemented in on-line or off-line mode. In the first case, there must be transmitters at the stations to transmit their measurements of TOA to the height monitoring center. The Vertical Dilution Of Precision (VDOP) is a performance index defined by the ratio:

$$VDOP = \frac{\sigma_z}{\sigma_{rte}} \quad (10.28)$$

where σ_{rte} is the SD of the relative timing errors. The VDOP expresses the effect of the relative geometry to system performance. The VDOP of this system achieves a typical value of 3.

Height Estimation with SSRs and Pseudoranges

This approach has been investigated by Nagaoka at Electronic Navigation Institute of Tokyo [9]. The system consists of N SSR receiving stations; see Figure 10.5. One of them, say S_0 , must be active to transmit interrogations to aircraft. The time of interrogation transmission, T_0 , and the times of signal arrival at the receiving stations T_i , $i = 0, 1, \dots, N-1$, are measured. Thus, N pseudorange measurements r_i are obtained where $r_i = c(T_i - T_0)$. Let T_D denote the transponder delay and D denote the distance corresponding to this delay, $D = cT_D$. Then, for each pseudorange measurement r_i , the following relation holds:

$$\begin{aligned} r_i &= D + R_i + R_0 = D + \sqrt{(x-x_i)^2 + (y-y_i)^2 + (z-z_i)^2} + \sqrt{(x-x_0)^2 + (y-y_0)^2 + (z-z_0)^2} \\ &= f_i(\mathbf{q}) \quad i = 0, 1, \dots, N-1 \end{aligned} \quad (10.29)$$

where $\mathbf{q} = [x, y, z, D]^T$ is the unknown vector. The set of N measurements $\rho = [\rho_0, \rho_1, \dots, \rho_{N-1}]^T$, $N \geq 4$, and the unknown vector are related through Equation 10.30.

$$\mathbf{r} = \mathbf{f}(\mathbf{q}) \quad (10.30)$$

The unknown vector \mathbf{q} can be obtained from the solution of Equation 10.30 with a nonlinear weighted least squares method. Thus, the best estimate of \mathbf{q} is iteratively calculated as:

$$\hat{\mathbf{q}}_{k+1} = \hat{\mathbf{q}}_k + (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T (\mathbf{r} - \mathbf{f}(\hat{\mathbf{q}}_k)) \quad (10.31)$$

where \mathbf{F} is the Jacobian matrix

$$\mathbf{F} = \frac{\partial \mathbf{f}}{\partial \mathbf{q}} = \begin{bmatrix} \frac{\partial f_0}{\partial x} & \frac{\partial f_0}{\partial y} & \frac{\partial f_0}{\partial z} & \frac{\partial f_0}{\partial D} \\ \cdot & \cdot & \cdot & \cdot \\ \frac{\partial f_{N-1}}{\partial x} & \frac{\partial f_{N-1}}{\partial y} & \frac{\partial f_{N-1}}{\partial z} & \frac{\partial f_{N-1}}{\partial D} \end{bmatrix} \quad (10.32)$$

The estimate of \mathbf{q} is free of the transponder delay systematic error because the estimation is based not on the nominal delay, but on the actual delay time, which is one of the parameters to be estimated, whereas the rest of the parameters are the aircraft 3-D position coordinates x, y, z .

The station arrangement proposed by Nagaoka, when there are four stations, is an equilateral triangle formed by the three stations, whereas the fourth station is located in the center. The VDOP, defined as the ratio σ_z/σ_R (where σ_R is the observation error) has a typical value of 4 when the aircraft is above the center at a height equal to the baseline radius. The VDOP increases as the aircraft flies higher and longer and as the baseline radius becomes smaller.

Height Measurement with SSRs and Range Differences

This approach has been proposed by Manolakis and Lefas [10, 11]. The system consists of $N-1$ receiving SSR stations S_i , $i = 1, N-1$, and one station, say S_0 , which is both receiver and interrogator, see Figure 10.5. The stations receive the reply and the time difference of arrival (TDOA) between a reference station, say S_0 , and station S_i is measured. A set of $N-1$ TDOA or equivalently range difference (RD) measurements is collected at each time the transponder sends a reply signal. The height estimation derived from this set of measurements is not affected by any transponder delay systematic error since this error is inherently subtracted from the measurements used. This system could be referred to as RD height monitoring unit (RDHMU). The systems that derive the position fix based on this kind of measurement are known as TDOA or RD or hyperbolic systems.

Let τ_i denote the TDOA between stations S_i and S_0 , and d_i denote the corresponding RD measurement, $d_i = c \tau_i$. The following relation holds:

$$\begin{aligned} d_i &= R_i - R_0 = \sqrt{(x-x_i)^2 + (y-y_i)^2 + (z-z_i)^2} - \sqrt{(x-x_0)^2 + (y-y_0)^2 + (z-z_0)^2} \\ &= f_i(\mathbf{q}) \quad i = 1, 2, \dots, N-1 \end{aligned} \quad (10.33)$$

where $\mathbf{q} = [x, y, z]^T$ is the unknown aircraft position vector. The vector of RD measurements $\mathbf{d} = [d_1, d_2, \dots, d_{N-1}]$ is expressed as:

$$\mathbf{d} = \mathbf{f}(\mathbf{q}) \quad (10.34)$$

A commonly employed method to solve for \mathbf{q} in this nonlinear equation is the Taylor series method or equivalently the Gauss-Newton iterative method. The best estimate of \mathbf{q} is iteratively approximated as:

$$\hat{\mathbf{q}}_{k+1} = \hat{\mathbf{q}}_k + (\mathbf{F}^T \mathbf{F})^{-1} \mathbf{F}^T (\mathbf{d} - \mathbf{f}(\hat{\mathbf{q}}_k)) \quad (10.35)$$

where \mathbf{F} is the Jacobian matrix:

$$\mathbf{F} = \frac{\partial \mathbf{f}}{\partial \mathbf{q}} = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} & \frac{\partial f_1}{\partial z} \\ \frac{\partial f_{N-1}}{\partial x} & \frac{\partial f_{N-1}}{\partial y} & \frac{\partial f_{N-1}}{\partial z} \end{bmatrix} \quad (10.36)$$

In the case of four stations the best arrangement is an equilateral triangle with the fourth station in the center. The SD of height estimation error σ_z will be 15 m when the baseline radius is 6 km, the flying height is 9 km, and σ_{TDOA} is 10 ns.

Work on proof of principles and system development of a HMU based on the concept of TDOA measurement of SSR signals has been conducted by Roke Manor Research Ltd., U.K. [12].

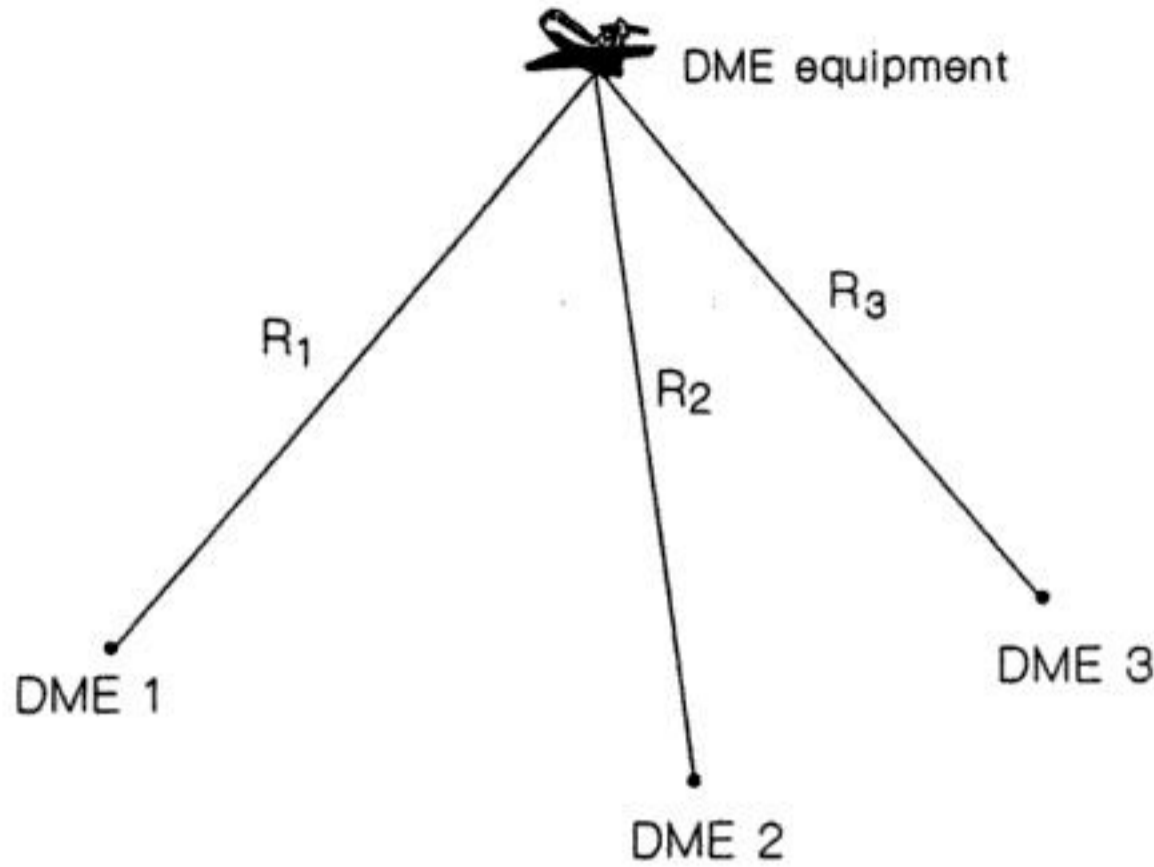


FIGURE 10.6 Configuration of the on-board height estimation system that utilizes the distance measurements derived from the DME equipment.

Onboard Derived Height Estimation

Height Measurement with Distance Measuring Equipment (DME)

This approach for deriving the geometric height onboard the aircraft using DME equipment was first reported by Rekkas et al. [13], whereas more efficient and general techniques have been proposed by Manolakis [14, 15]. Using the DME interrogation equipment, the distance from three DME ground stations is measured onboard (see Figure 10.6). The three stations are located under an airway. The height is then computed from the range measurement vector $\mathbf{R} = [R_1, R_2, R_3]^T$ by trilateration. An exact and efficient solution of the nonlinear measurement equation was derived in [15]. Specifically, the height is computed from the closed form:

$$z = g(\mathbf{R}) = \frac{-b(\mathbf{R}) + \sqrt{d(\mathbf{R})}}{2a} \quad (10.37)$$

where $b(\mathbf{R})$ and $d(\mathbf{R})$ are the following simple polynomial-type functions

$$b(\mathbf{R}) = b_0 + b_1 R_1^2 + b_2 R_2^2 + b_3 R_3^2 \quad (10.38)$$

$$d(\mathbf{R}) = d_{00} + d_{01} R_1^2 + d_{02} R_2^2 + d_{03} R_3^2 + d_{11} R_1^4 + d_{22} R_2^4 + d_{33} R_3^4 + d_{12} R_1^2 R_2^2 + d_{13} R_1^2 R_3^2 + d_{23} R_2^2 R_3^2 \quad (10.39)$$

The coefficients a , b_i , d_{ij} are analytically defined in the Appendix of [15]. An important aspect of these coefficients is that they are completely defined by the ground stations' coordinates (x_i, y_i, z_i) , which are fixed. Thus, the coefficients are calculated only once at the moment the aircraft enters the data acquisition area. Then, every time a new set of range measurements is available, the height is computed from the above equations using the range measurements and the stored coefficients. Define the ratio σ_z/σ_R as the VDOP of this technique, where σ_R is the SD of the ranging error. The VDOP is 1 in the case where the stations form an equilateral triangle inscribed in a circle with 10 km radius and the aircraft is above the triangle center at a height of 8 km.

Estimation of Vertical Position with the Global Positioning System (GPS)

The research and development of the GPS has been coordinated by the U.S. Department of Defense. Another similar system is the Global Navigation Satellite System (GLONASS) developed by the former Soviet Union. The GPS is a satellite system providing users with accurate timing and ranging information. The system is available with reduced accuracy to civilian users. Many companies, mainly from the U.S., produce GPS receivers. Let (x, y, z) and (x_i, y_i, z_i) be the coordinates of the user and satellite s_i . The GPS receiver of the user derives the pseudorange measurement D_i , and the corresponding measurement equation is:

$$D_i = R_i + cT_b = \sqrt{(x-x_i)^2 + (y-y_i)^2 + (z-z_i)^2} + b = f_i(\mathbf{q}) \quad (10.40)$$

where T_b is the user clock bias, and $\mathbf{q} = [x, y, z, b]^T$ is the unknown vector that incorporates the bias term b . Thus, in order to estimate the 3-D position of the aircraft, four pseudorange measurements are required at least; consequently, four satellites must be visible from the receiver. The set of N pseudorange measurements $\mathbf{D} = [D_1, D_2, \dots, D_N]^T$ defines the following matrix measurement equation:

$$\mathbf{D} = \mathbf{f}(\mathbf{q}) \quad (10.41)$$

which is solved for \mathbf{q} with the Gauss-Newton least squares iterative method, that is:

$$\hat{\mathbf{q}}_{k+1} = \hat{\mathbf{q}}_k + (\mathbf{A}^T \mathbf{W} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{W} (\mathbf{D} - \mathbf{f}(\hat{\mathbf{q}}_k)) \quad (10.42)$$

where \mathbf{A} is the partial derivatives matrix:

$$\mathbf{A} = \frac{\partial \mathbf{f}}{\partial \mathbf{q}} = \begin{bmatrix} \frac{\partial f_1}{\partial x} & \frac{\partial f_1}{\partial y} & \frac{\partial f_1}{\partial z} & \frac{\partial f_1}{\partial b} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_N}{\partial x} & \frac{\partial f_N}{\partial y} & \frac{\partial f_N}{\partial z} & \frac{\partial f_N}{\partial b} \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{x1} & \mathbf{a}_{y1} & \mathbf{a}_{z1} & 1 \\ \vdots & \vdots & \vdots & \vdots \\ \mathbf{a}_{xN} & \mathbf{a}_{yN} & \mathbf{a}_{zN} & 1 \end{bmatrix} \quad (10.43)$$

The elements a_{xi} , a_{yi} , a_{zi} , of the partial derivatives matrix \mathbf{A} are the direction cosines from the receiver to the satellite s_i . The weighting matrix is the inverse of the covariance matrix of the pseudorange measurement errors, $\mathbf{W}^{-1} = \mathbf{E}(\delta \mathbf{D} \delta \mathbf{D}^T)$. The weighting is generally used to take into account the possible different performances of each satellite, although usually the same performance is assumed for all satellites; that is, $\mathbf{W} = \mathbf{I}$. The VDOP, defined as σ_z / σ_D , depends on the geometry which varies continuously, even in the case of a fixed receiver, because the satellites are not geostationary but move in such orbits as to complete a rotation in 12 h. The world mean value of VDOP is about 2 [16]. Typical VDOP values range from 1.5 to 7, depending on the area of the receiver and on the time of day. The ranging error for the precision positioning service (available only to U.S. military users) has been specified to be less than 6 m (SD). For the standard positioning service, normally available to civilian users, the specified ranging error is double (12 m, SD), whereas it will be about 40 m when selective availability is activated by the Department of Defense. The corresponding measured ranging errors found to be smaller than the specified ones. Namely, for the three operating conditions mentioned, the corresponding values for ranging errors were found to be 2.3 m, 6 m, and 20 m, respectively [17]. The multiplication of the ranging SD error by the VDOP yields the standard deviation of the height estimation error.

TABLE 10.3 Specification of a Servo-Type Accelerometer

Measurement range	$\pm 5 g$
Resolution	Less than $5 \mu g$ (dc)
Sensitivity	$2 V g^{-1}$
Output resistance	560Ω
Torquer current	$3.5 mA g^{-1}$
Case alignment	Less than $\pm 1^\circ$
Frequency response	$450 Hz (\pm 3 dB)$
Temperature range	-25 to $+70^\circ C$
Power source	$\pm 15 V$ (dc)
Consumption current	Less than $15 mA$
Size	$28.4 mm \times 24.5 mm$
Mass	$46 g$ (including the cable $10 g$)

Note: g : gravitational acceleration (according to the type TA-25D-05 by TOKIMEC).

sensing system that overcomes such difficulty will be introduced later. Although application is limited to static inclined surfaces with minute tilt angles, a dielectric-type inclinometer employing electrodes and a bubble kept in an electrolyte can achieve high accuracies on the order of 10^{-4}° .

Attitude Sensors for Spacecraft Applications

Attitude measurement for spacecraft usually requires two or more sensors for detecting the reference sources needed to satisfy attitude requirements. The choice of which sensors to employ is primarily influenced by the direction the spacecraft is usually pointing as well as the accuracy requirements for attitude determination [2]. Table 10.4 summarizes some performance parameters for these sensors as well as typical manufacturers.

Inertial measurement units generally consist of gyroscopes coupled with accelerometers, which together measure both rotational and translational motion. These IMUs may be either gimbal mounted (movement about a gimbal point, independent of the spacecraft) or a strapdown system (rigidly mounted to the spacecraft body), where expansive software is used to convert sensor outputs into reference frame measurements. IMUs tend to suffer gyro drift and other bias errors and, when used for spacecraft attitude measurements, are often used with one or more of the sensors discussed below.

Sun sensors detect the visible light from the sun, measuring the angle between the sun's radiation and the detector's photocell. The sun is a commonly chosen attitude reference source since it is by far the

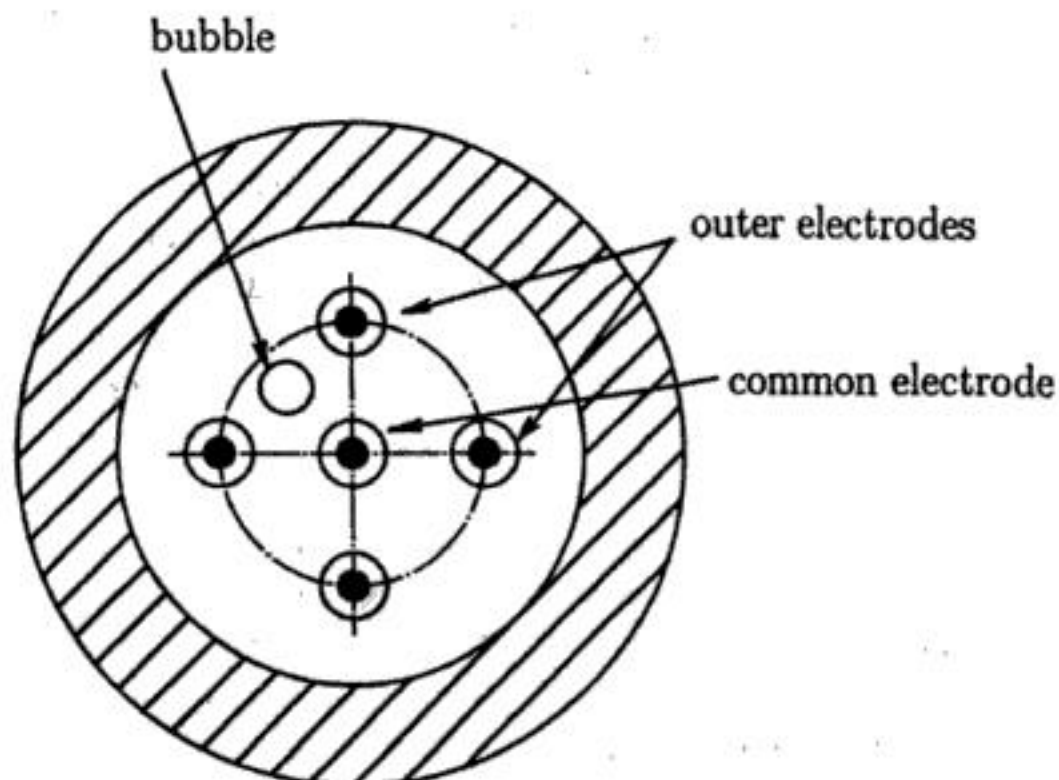
**FIGURE 10.9** Dielectric-type inclinometer (front view).

TABLE 10.4 Spacecraft Attitude Determination Sensors

Sensor	Accuracy	Mass (kg)	Typical vendors
IMU	1 to 5×10^{-6} g	3 to 25	Northrop Grumman, Bendix, Kearfott, Honeywell, Hamilton, Standard, Litton, Teledyne
Sun Sensor	10^{-2} to 3°	0.5 to 2	Adcole, TRW, Ball Aerospace
Horizon Sensor	10^{-1} to 1°	2 to 5	Barnes, Ithaco, Lockheed Martin, Lockheed Barnes
Star Sensor	10^{-3} to 10^{-2}	3 to 7	Ball Aerospace, Bendix, Honeywell, Hughes
Magnetometer	0.5 to 3°	~1	Schonstedt, Develco

Adapted from Larsen, W. J. and Wertz, J. R., Eds., *Space Mission Analysis and Design*, Torrance, CA: Microcosm Inc. and Dordrecht, The Netherlands: Kluwer Academic Publishers, 1992, p.360.

visually brightest object in the sky, having a total radiation per unit area of 1353 W m^{-2} at Earth distances [3]. Also, it is generally accepted as a valid point source for most attitude applications, having an angular radius of 0.25° at Earth distances. Increased measurement accuracy can be obtained by determining its centroid. Even though sun sensors are quite accurate (0.01° to 3.0°), they do require clear fields of view, and sometimes suffer periods of eclipse from both the Earth and the moon [4]. Also, sensitive equipment (such as imaging devices) must be protected from the powerful radiation of direct sunlight. When the sun is available, the angle between it and the sensor's primary axis is referred to as the *sun angle*.

For spacecraft in near-Earth orbits, the Earth is the second brightest object in the sky and covers as much as 40% of the sky. Earth *horizon sensors* detect the interface between the Earth's edge (or limb) and the space background. Horizon sensors can detect either of the Earth's visible limb (albedo sensor), infrared limb, or air glow. The infrared limb is the edge between the warm Earth and the cold space background. The air glow is a region of the atmosphere around the Earth that is visible to the spacecraft when it is on the night side of the Earth. Accuracies for horizon sensors are in the 0.1° to 1.0° range. Increased accuracy requires Earth oblate spheroid modeling [4]. Some problems associated with albedo detection include the distortion effects of the Earth's atmosphere, falsely identifying the day/night terminator crossing as the true Earth limb, and the considerable variability of the Earth's albedo in the visible spectrum (varies from land, sea, ice).

Most sensors used to detect the Earth's horizon are scanning sensors with narrow fields of view that measure the time between horizon crossings. In general, two horizon crossings occur per sensor scan period: one crossing when the sensor scans from the space background onto the Earth, followed by a second crossing when the sensor scans from the Earth back to space. The combination of horizon crossing times, scan rate, and spacecraft altitude allows for the computation of the Earth's apparent *angular radius*. The apparent angular radius will be smaller than the real (or physical) angular radius if the spacecraft is tilted away from the Earth nadir vector. The nadir vector is defined as the vector connecting the center of the spacecraft to the center of the Earth. To see this effect, one needs to compute the Earth's physical radius ρ , which for a given spacecraft altitude h (in kilometers), is given by $\rho = \sin^{-1}[(6371)/(6371 + h)]$.

If the spacecraft horizon sensor is pointing exactly *nadir*, then the apparent angular radius as measured by the sensor will agree with the physical radius given by the above relation for ρ . However, if the horizon sensor is pointed away from nadir, the horizon crossing times will be smaller than when pointing exactly nadir. This results in an apparent angular radius that is smaller than the physical radius by an amount proportional to the angle between the sensor axis and the nadir vector. This angle is referred to as the *nadir angle*.

Star sensors are used when extreme accuracy requirements are necessary. This high degree of sensor accuracy (0.003° to 0.01°) can be attributed mainly to the point source nature and precise fixed location of stars in space. Star sensors may be categorized as either star trackers or star mappers. A star tracker utilizes a wide field of view in order to search for a given star of specific brightness. A star mapper is similar to a tracker, except that it scans over many stars, recording their relative positions and angular separations. By comparing the recorded data with that from a *star catalog* (database), exact spacecraft

orientation can be obtained. The angle between the star line-of-sight and the sensor's primary axis is referred to as the *star angle*.

The accuracy of star sensors is obtained with higher costs, however. Star sensors are generally heavier and consume more power than other types of attitude sensors. In addition, star sensors are quite sensitive to stray light sources such as sunlight reflected from the spacecraft or the Earth and sunlight scattered from dust particles and jet exhausts [4]. Most rely on optical shielding to reduce the effects of stray light.

Magnetic sensors (called *magnetometers*) measure both the magnitude and direction of the Earth's magnetic field. The difference in orientation between the measured field and the true field translates into attitude determination. Magnetometer accuracies (0.5° to 3.0°) are usually less than the other sensor types because of the uncertainty in the Earth's true field, which tends to change or shift over time. In addition, the Earth's magnetic field decreases with increasing altitude, and magnetometers are generally limited to altitudes of about 6000 km. For this reason, magnetometers are often used with one of the other sensor types already discussed for improved measurement accuracy [2].

Automatic On-Line Attitude Measurement for Ships and Crane Lifters

For on-line attitude measurement for ships and crane lifters, the first thing that comes to mind is to use gyros. However, because they often suffer from drifts, accurate attitude measurements might not be achieved using the gyros. Accordingly, one uses attitude on-line measurement systems that do not utilize gyros but servo-type accelerometers and inclinometers. The philosophy of the measurement systems introduced here is to make the best use of the system dynamics of the object and the sensors and to apply Kalman filters or adaptive filters to achieve high measurement accuracy.

Attitude Measurement for Ships

On-line accurate measurement of a ship's attitude is extremely important in exact search of the seabed patterns with sonars [5, 6]. It is also required by high-performance ships like hovercrafts from the viewpoint of suppressing swings by the waves. The measurement of a ship's attitude can usually be reduced to that of the heaving, rolling, and pitching of the ship. For such a measurement, a heave sensor has been used, whose output is given by double integration of the output of an accelerometer vertically directed with a gyroscope. However, since the initial values of heaving displacement and its velocity are unknown, the output will contain a bias that increases with time, and the accuracy of the sensor deteriorates considerably. From this viewpoint, one introduces a strapdown-type on-line measurement system that adequately processes the outputs of the two servo-type inclinometers and one accelerometer mounted on the ship [7].

Location of Sensors and Outputs

The two servo-type inclinometers and one servo-type accelerometer are located on the deck (at the point A) of vertical distance L from O, the intersection of rolling and pitching axes (see Figure 10.10). The two inclinometers are set in such a way that the rolling and the pitching angles are measured respectively. The accelerometer is set upward to the deck to obtain the information on the heaving. Because inclinometers were originally developed for the measurement of the tilt angles of static inclined surfaces, the rigid pendulum inside the sensor is considerably affected by the ship's acceleration other than the gravitational one. Applying Lagrange's equations of motion [8, 9] to rigid pendulums and calculating the torques to keep their deflections from the principal axes almost zero yields the sensor outputs [7]:

$$z_1(t) = \theta(t) - \frac{L}{g} \ddot{\theta}(t) + v_1(t) \quad (10.45)$$

$$z_2(t) = p(t) - \frac{L}{g} \ddot{p}(t) + v_2(t) \quad (10.46)$$

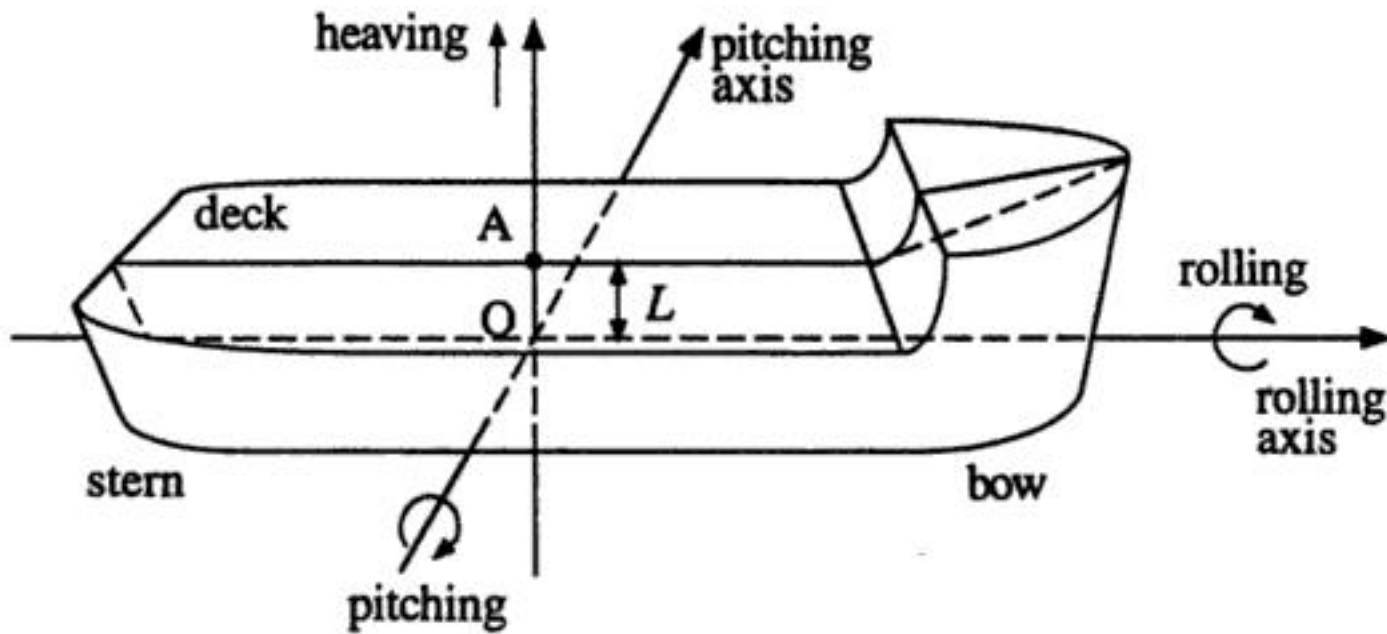


FIGURE 10.10 Location of sensors.

where $z_1(t)$, $z_2(t)$, $\theta(t)$, $p(t)$, and g denote, respectively, the outputs of the two inclinometers, the rolling and the pitching angles, and the gravitational acceleration ($v_1(t)$ and $v_2(t)$: noises of the outputs, including the approximation errors in deriving the outputs).

On the other hand, the accelerometer output is expressed as:

$$z_3(t) = (g + \alpha(t)) \cos\theta(t) \cos p(t) + v_3(t) \tag{10.47}$$

where $\alpha(t)$ and $v_3(t)$ represent, respectively, the heaving acceleration and the accelerometer noise.

Dynamics of Attitude Signals

It is well known that each of the heaving, rolling, and pitching in inshore seas has two dominant waves in a short interval. That is, a sinusoidal wave of long periodic length (in the range of 6 s to 10 s) and a sinusoidal wave of short periodic length (in the range of 2 s to 3 s) [10–12]. Thus, one model each of the signals in a short interval by a composite wave of the two dominant sinusoidal waves. For the heaving (in a short interval), the displacement is modeled by:

$$x(t) = a_1 \sin(\omega_1 t + \varphi_1) + a_2 \sin(\omega_2 t + \varphi_2) \tag{10.48}$$

with the parameters $\{a_i\}$, $\{\varphi_i\}$, and $\{\omega_i\}$ unknown. From the 4th-order differential equation satisfied by the $x(t)$, we obtain the linear dynamic equation [7]:

$$\dot{\mathbf{x}}(t) = A\mathbf{x}(t), \quad A \equiv \begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ -\omega_1^2 & \omega_2^2 & 0 & -(\omega_1^2 + \omega_2^2) \end{bmatrix} \tag{10.49}$$

where $\mathbf{x}(t) \equiv (x_1, x_2, x_3, x_4)^T$ ($x_n = d^{n-1}x/dt^{n-1}$ ($n = 1, \dots, 4$)). On the other hand, the rolling and pitching angles can be modeled by:

$$x(t) = a_1 \sin(\omega_1 t + \varphi_1) + a_2 \sin(\omega_2 t + \varphi_2) + b \tag{10.50}$$

because there are usually some biases associated with them. From the 5th-order differential equation which Equation 10.50 satisfies, we get the similar state variable representation of the model as

Equation 10.49. In practice, the heaving, rolling, and pitching signals have many nondominant sinusoidal waves in addition to the dominant ones. Therefore, Equation 10.49 is modified by introducing a white Gaussian noise $w(t)$ with zero mean and adequate variance σ^2 as follows:

$$\dot{\mathbf{x}}(t) = A\mathbf{x}(t) + \Gamma w(t) \quad (10.51)$$

where $\Gamma = (0,1,0,0)^T$ for the heaving and $\Gamma = (0,1,0,0,0)^T$ for the rolling and pitching. The higher the order of the models, the better the measurement accuracy will be. If we consider the on-line measurement of the signals, Equation 10.51 will be sufficient.

On-Line Attitude Measurement

The observation Equations 10.45 and 10.46 are expressed using their own state vector $\mathbf{x}(t)$. The observation equations in a discretized form are:

$$y_k = H\mathbf{x}_k + v_k \quad (10.52)$$

where $H = [1,0,-L/g,0,0]$ and y_k , \mathbf{x}_k , and v_k , respectively, denote $y(t)$, $\mathbf{x}(t)$, and $v(t)$ of the corresponding signals at the k -th sampling instant [7, 9]. The discretized form of the dynamic Equation 10.51 is:

$$\mathbf{x}_{k+1} = F\mathbf{x}_k + \mathbf{w}_k \quad (10.53)$$

where

$$F \equiv \Phi(t) \Big|_{t=\Delta T}, \quad \Phi(t) \equiv L^{-1} \left\{ (sI - A)^{-1} \right\}. \quad (10.54)$$

Here, L^{-1} and ΔT , respectively, denote the inverse Laplace transformation and the sampling period. The discretized transition noise \mathbf{w}_k becomes a white Gaussian noise with zero mean and covariance:

$$W = \sigma^2 \int_0^{\Delta T} \Phi(\Delta T - \tau) \Gamma \Gamma^T \Phi^T(\Delta T - \tau) d\tau \quad (10.55)$$

The measurement of the rolling and pitching can thus be reduced to the state estimation of the linear discrete dynamic systems (Equations 10.52 and 10.53), if the angular frequencies ω_1 and ω_2 are given and v_k is assumed to have a white Gaussian property. The state estimation is achieved by a Kalman filter [7, 13]. However, difficulties in implementing the filter are that the exact values of the two angular frequencies are a priori unknown and also time variant. To overcome the difficulty, adequate candidates $\{(\omega_1^i, \omega_2^i); 1 \leq i \leq M\}$ for the parameters $\{\omega_1, \omega_2\}$ are set and a bank of Kalman filters is used. Then, the final estimate is obtained as the conditional expectation of the state estimate as follows:

$$\hat{\mathbf{x}}_{k/k}^0 \equiv \sum_{i=1}^M p_k^i \hat{\mathbf{x}}_{k/k}^i \quad (10.56)$$

where $\hat{\mathbf{x}}_{k/k}^i$ represents the state estimate $\hat{\mathbf{x}}_{k/k}$ for the i -th candidate $\Omega_i = (\omega_1^i, \omega_2^i)$, and p_k^i denotes the conditional posteriori probability of the i -th candidate calculated based on the Bayesian theorem:

$$p_k^i = \frac{P(y_k / \Omega_i, Y^{k-1}) p_{k-1}^i}{\sum_{j=1}^M p_{k-1}^j P(y_k / \Omega_j, Y^{k-1})} \quad (10.57)$$

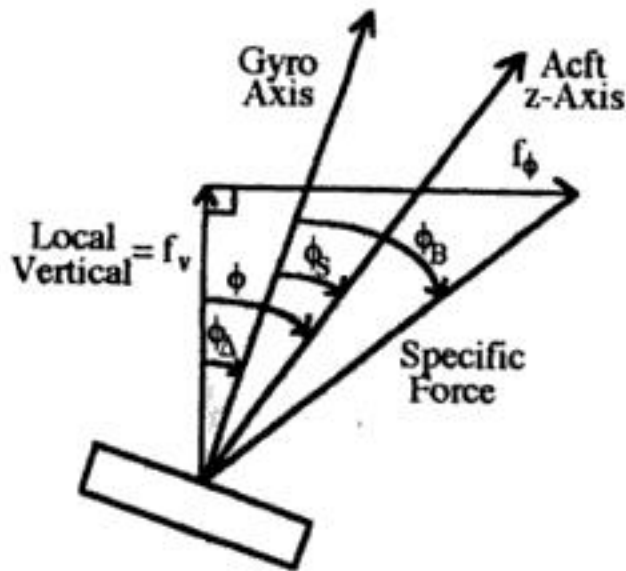


FIGURE 10.13 Vertical gyro analysis.

Figure 10.13, depicting an aircraft with a roll angle of ϕ with respect to the local vertical. The sensed roll angle ϕ_s is given by the difference in the actual roll angle and the gyro roll drift error ϕ_Δ :

$$\phi_s = \phi - \phi_\Delta \quad (10.64)$$

In order to compensate for this drift error, gyros employ a specific force sensor such as an electrolytic bubble device, which senses drift error. This correction device senses the angular difference between the specific force vector \mathbf{f} acting on the aircraft roll axis and the gyro axis, as shown in Figure 10.13. Thus,

$$\phi_B = \tan^{-1} \left[\left(\frac{f_\phi}{f_v} \right) - \phi_\Delta \right] \cong \left(\frac{f_\phi}{f_v} \right) - \phi_\Delta \quad (10.65)$$

where f_ϕ is the side horizontal component of \mathbf{f} and $f_v =$ force of gravity is the vertical component. A similar analysis for the pitch angle θ yields:

$$\theta_s = \theta - \theta_\Delta \quad (10.66)$$

$$\theta_B = \tan^{-1} \left[\frac{f_\theta}{f_v} - \theta_\Delta \right] \cong \frac{f_\theta}{f_v} - \theta_\Delta \quad (10.67)$$

where f_θ is the back horizontal component of \mathbf{f} . Next, define the gyro angular momentum vector by:

$$\mathbf{H}_{VG} = [J_x \dot{\phi}_\Delta, J_y \dot{\theta}_\Delta, -h] \quad (10.68)$$

where J_x and J_y are the sensor moments of inertia and h is the gyro spin angular momentum. In addition, define the inner gimbal axes angular velocity vector as:

$$\boldsymbol{\omega}_{VG} = [\dot{\phi}_\Delta, \dot{\theta}_\Delta, 0] \quad (10.69)$$

Finally, define the gimbal torque vector by:

$$\mathbf{Q}_{VG} = [Q_{cx} + Q_{dx}, Q_{cy} + Q_{dy}, 0] \quad (10.70)$$

where

$$Q_{cx} = \text{gimbal roll control torque} = -k_c \theta_B \quad (10.71a)$$

$$Q_{cy} = \text{gimbal pitch control torque} = k_c \phi_B \quad (10.71b)$$

$$Q_{dx} = \text{gimbal roll disturbance torque} = -k_d(\dot{\phi}_\Delta - \dot{\phi}) + \text{random torques} \quad (10.71c)$$

$$Q_{dy} = \text{gimbal pitch disturbance torque} = -k_d(\dot{\theta}_\Delta - \dot{\theta}) + \text{random torques} \quad (10.71d)$$

and the k_c and k_d are constant scaling factors related to each torque component.

Using the vectors defined in Equations 10.68 through 10.70, the gyro equations of motion are given by:

$$\frac{\partial}{\partial t}(\mathbf{H}_{VG}) + (\boldsymbol{\omega}_{VG} \times \mathbf{H}_{VG}) = \mathbf{Q}_{VG} \quad (10.72)$$

Taking the Laplace transform of the expansion of Equation 10.72, with the assumption that $J_x \cong J_y = J$, yields the following gyro equations of motion in the Laplace domain:

$$\begin{bmatrix} J_x s^2 + k_d s & -(hs + k_c) \\ hs + k_c & J_y s^2 + k_d s \end{bmatrix} \begin{bmatrix} \phi_\Delta(s) \\ \theta_\Delta(s) \end{bmatrix} \equiv \begin{bmatrix} -k_c \theta_B + k_d s \phi(s) + \text{random torques} \\ (k_c/g) f_\phi(s) + k_d s \theta(s) + \text{random torques} \end{bmatrix} \quad (10.73)$$

For normal gyro operation, $J_x \cong J_y \cong 0$ and $k_d/h \ll 1$; so these factors may be ignored in Equation 10.73. Thus, solving for the desired roll and pitch angles under these assumptions gives [1]:

$$\phi_s = \begin{cases} \phi & \omega \gg k_c/h \\ \phi - f_\phi/g & \omega \ll k_c/h \end{cases} \quad (10.74)$$

$$\theta_s = \begin{cases} \theta & \omega \gg k_c/h \\ \theta - f_\theta/g & \omega \ll k_c/h \end{cases} \quad (10.75)$$

A DG is a two degree-of-freedom gyro with its spin axis mounted nominally horizontal and pointing in the direction of magnetic north. It employs a single specific force sensor mounted on the inner gimbal [1]. The DG measures the third required aircraft angle, yaw, generally denoted by ψ . The sensed yaw angle ψ_s is given by the difference in the actual yaw angle ψ (angle between the aircraft z-axis and true north) and the gyro heading angle drift error ψ_Δ (angle between the gyro axis and true north):

$$\psi_s = \psi - \psi_\Delta \quad (10.76)$$

Define the gyro angular momentum vector by:

$$\mathbf{H}_{DG} = [J_y \dot{\theta}_\Delta, J_z \dot{\psi}_\Delta, -h] \quad (10.77)$$

and the inner gimbal axes angular velocity vector as:

$$\boldsymbol{\omega}_{DG} = [\dot{\theta}_\Delta, \dot{\psi}_\Delta, 0] \quad (10.78)$$

and the gimbal torque vector as

$$Q_{DG} = [Q_{cy} + Q_{dy}, Q_{cz} + Q_{dz}, 0] \quad (10.79)$$

Here, the torque vector components are given by:

$$Q_{cy} = k_c(M_\Delta - \psi_\Delta) \quad (10.80a)$$

$$Q_{cz} = -k_c\theta_B \quad (10.80b)$$

$$Q_{dy} = -k_d(\dot{\theta}_\Delta - \dot{\theta}) + \text{random torques} \quad (10.80c)$$

$$Q_{dz} = -k_d(\dot{\psi}_\Delta - \dot{\psi}) + \text{random torques} \quad (10.80d)$$

where M_Δ = magnetic compass heading error (from true north). Therefore, the DG equations of motion are given in Laplace domain as:

$$\begin{bmatrix} J_y s^2 + k_d s & -(hs + k_c) \\ hs + k_c & J_z s^2 + k_d s \end{bmatrix} \begin{bmatrix} \theta_\Delta(s) \\ \psi_\Delta(s) \end{bmatrix} \equiv \begin{bmatrix} k_c M_\Delta(s) + k_d s \theta(s) + \text{random torques} \\ -k_c \theta_B + k_d s \psi(s) + \text{random torques} \end{bmatrix} \quad (10.81)$$

The desired yaw angle measurement for the DG is thus given as [1]:

$$\psi_s = \begin{cases} \psi & \omega \gg k_c/h \\ \psi - M_\Delta & \omega \ll k_c/h \end{cases} \quad (10.82)$$

As indicated in Table 10.2, the accuracies of both VGs and DGs are approximately 1° . An improvement of over 2 orders of magnitude can be obtained through the use of inertial measurement units, which are described next.

Inertial Measurement Units (IMUs)

Inertial measurement units consist of gyroscopes and accelerometers that together provide full three-axis attitude measurements. Most are mounted on stable gimballed platforms that remain locally horizontal via torquing devices. An IMU aboard an aircraft cannot measure exactly the local vertical due to the fact that the specific force acting on the aircraft has a horizontal component due to vehicle motion. In addition, since the vehicle is moving with respect to the inertial reference frame, the Earth's magnetic pole cannot be determined precisely [1].

These problems (errors) are minimized by aligning the IMU to be exactly horizontal and north pointing while the aircraft is stationary. Once platform motion begins, the IMU may be constantly realigned by sensing changes in the direction of vertical and north, and then applying appropriate torques to the platform to keep it properly aligned. This realignment is accomplished by integrating the two orthogonal accelerometer outputs to determine the components of horizontal velocity. This data, combined with the Earth's rotation rate, yields the desired rates of change in local vertical and true north at the vehicle's current latitude and longitude. Performing a second integration of the sensor outputs yields an estimate of relative position.

Analysis in Bryson et al. [1], has shown that the pitch angle (variation in platform horizontal position) is given by the IMU sensor output as:

$$\theta(t) = \frac{-\varepsilon}{\omega_s} \sin(\omega_s t) - \frac{b}{g} \quad (10.83)$$

where ε = gyro drift rate error, b = specific force sensor error, and $\omega_s \equiv$ Schuler frequency = $\sqrt{g/R}$, [g = force of gravity, R = Earth's radius]. Thus, the platform root-mean-square pitch angle becomes:

$$\theta_{\text{rms}} = \left[\frac{1}{2} \left(\frac{\varepsilon}{\omega_s} \right)^2 + \left(\frac{b}{g} \right)^2 \right]^{\frac{1}{2}} \quad (10.84)$$

Using typical values for ε ($\cong 0.015^\circ \text{ h}^{-1}$), ω_s ($\cong 0.71^\circ \text{ h}^{-1}$), and b ($\cong 0.01$) yields an rms pitch angle error of $\theta_{\text{rms}} = 0.01^\circ$. Thus, it is apparent that under normal operating conditions the IMU provides a two orders-of-magnitude improvement in sensor accuracy when compared to the VG and DG.

Spacecraft Attitude Determination

Most spacecraft attitude determination techniques rely upon finding the orientation of a single axis in space (e.g., the spacecraft z -axis) plus the spacecraft rotation about this axis. This provides a full three-axis attitude solution. In order to achieve this, reference sources that are external to the spacecraft must be used. Specifically, full three-axis spacecraft attitude determination requires at least two external vector measurements. Commonly used reference sources for these external vector measurements include the sun, Earth, moon, stars, planets, and the Earth's magnetic field. In addition, IMUs are also used to provide the necessary attitude measurements.

Attitude Determination Methodology

The first step in attitude determination is to determine the angles between the spacecraft's primary axis and the two (or more) attitude reference sources. For example, suppose a particular spacecraft is using the sun and the Earth for attitude reference. The two angles in this case are referred to as the sun angle β_s and the nadir angle Γ_N . Since the orientation of even a single spacecraft axis is unknown at this point, these angles establish two *cones* along which the attitude vector \mathbf{A} must lie. Since the attitude vector must lie on both cones, it must lie along the intersection between the two cones [4] (See Figure 10.14). The two vectors, notably \mathbf{A}_1 and \mathbf{A}_2 , resulting from the intersection of these two cones may be determined by the following method derived by Grubin [15]. Let \mathbf{S} represent the sun vector, \mathbf{E} the spacecraft nadir vector, and \mathbf{A} the desired attitude vector, each defined in Cartesian space as follows:

$$\mathbf{S} = (S_x, S_y, S_z) \quad (10.85)$$

$$\mathbf{E} = (E_x, E_y, E_z) \quad (10.86)$$

$$\mathbf{A} = (A_x, A_y, A_z) \quad (10.87)$$

Let the vectors \mathbf{S} , \mathbf{E} , and \mathbf{N} define a set of base unit vectors with:

$$\mathbf{N} = \frac{\mathbf{S} \times \mathbf{E}}{|\mathbf{S} \times \mathbf{E}|} = (N_x, N_y, N_z) \quad (10.88)$$

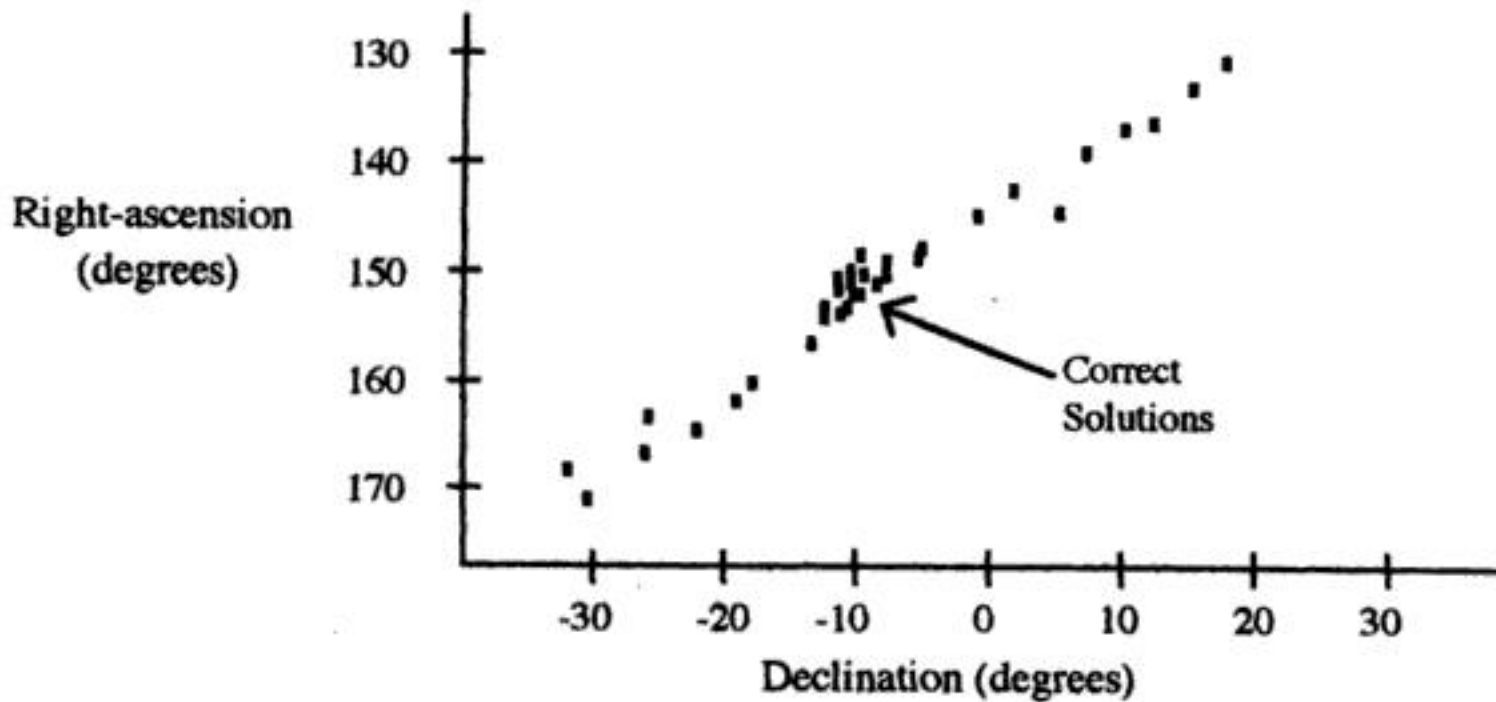


FIGURE 10.15 Method of trace averaging for resolving ambiguous attitude solutions.

values for I_z whenever the two cones do not intersect. Such occurrences are usually attributed to sensor error or random noise fluctuations. In this case, one can add a predetermined sensor bias to both sensors in order to “broaden” the cone angles, thus forcing the cones to intersect.

It should be noted that for most applications involving spacecraft attitude determination, the principle coordinate system used is the *celestial sphere* coordinate system. This coordinate system has the z -axis aligned with the Earth’s polar axis, and the x -axis aligned with the intersection of the Earth’s equatorial plane and the Earth’s orbital plane around the sun (i.e., aligned with the *vernal equinox*). In this coordinate system, all vectors are considered unit vectors and the two principle measurements describing a vector’s position are the *right-ascension* and *declination angles*, denoted Ω and Δ , respectively. Thus, the sun vector \mathbf{S} and the Earth nadir vector \mathbf{E} used in Equations 10.85 and 10.86 will, in general, be given as right-ascension and declination angles that can be converted to Cartesian coordinates via the following set of transformations:

$$x = \cos(\Omega)\cos(\Delta); \quad y = \sin(\Omega)\cos(\Delta); \quad z = \sin(\Delta) \quad (10.91a)$$

$$\Omega = \tan^{-1}(y/x); \quad \Delta = \sin^{-1}(z) \quad (10.91b)$$

The final step in measuring three-axis attitude is to determine which attitude solution is correct, \mathbf{A}_1 or \mathbf{A}_2 , and then measure the rotation about this axis. The two ambiguous attitude solutions may be resolved by comparison with a priori attitude information, if available, or through the use of *trace averaging* [4]. Trace averaging is a method of plotting each attitude solution on a right-ascension versus declination plot and choosing the area of greatest concentration as the correct solution, as demonstrated in Figure 10.15. Since the attitude is assumed to change more slowly than the attitude sensor’s sample rate, over short time intervals the data for the correct solution usually form a “cluster” near the correct attitude; the data for the incorrect solution are usually much more scattered.

Once the correct attitude vector has been obtained, the orientation of the remaining two orthogonal axes may be found by measuring the rotation, or phase angle, of the spacecraft about the preferred axis. Any sensor measurement that provides this phase angle may be used. An example of this technique is provided by the panoramic annular lens attitude determination system (PALADS), described in the next section. This imaging system uses a unique “three-dimensional” lens that provides simultaneous detection of two (or more) reference sources [16]. This information, combined with the orientation of the single axis, uniquely determines three-axis attitude.

4. J. R. Wertz (ed.), *Spacecraft Attitude Determination and Control*, Chapters 11 and 12, The Netherlands: Reidel Publishing Company, 1980.
5. R. D. Angelari, A deterministic and random error model for a multibeam hydrographic sonar system, *Proc. OCEANS'78. The Ocean Challenge*, 1978, 48-53.
6. C. de Moustier, T. Hylas, and J. C. Phillips, Modifications and improvements to the Sea Beam system on board R/V Thomas Washington, *Proc. OCEANS'88 — A Partnership of Marine Interests*, 1988, 372-378.
7. S. Tanaka and S. Nishifuji, Automatic on-line measurement of ship's attitude by use of a servo-type accelerometer and inclinometers, *IEEE Trans. Instrum. Meas.*, 45, 209-217, 1996.
8. D. G. Shultz and J. L. Melsa, *State Functions and Linear Control Systems*, New York: McGraw-Hill, 1967.
9. Y. Takahashi, M. J. Rabins, and D. M. Auslander, *Control and Dynamic Systems*, Reading, MA: Addison-Wesley, 1971.
10. S. Tanaka and S. Nishifuji, On-line sensing system of dynamic ship's attitude by use of servo-type accelerometers, *IEEE J. Oceanic Eng.*, 20, 339-346, 1995.
11. S. Tanaka, On automatic attitude measurement system for ships using servo-type accelerometers (in Japanese), *Trans. SICE*, 27, 861-869, 1991.
12. D. E. Cartwright and M. S. Longuet-Higgins, The statistical distribution of the maxima of a random function, *Proc. Roy. Soc. London, Ser. A*, 237, 212-232, 1956.
13. R. E. Kalman, A new approach to linear filtering and prediction problems, *Trans. ASME, J. Basic Eng.*, 82, 35-45, 1960.
14. S. Tanaka, S. Kouno, and H. Hayashi, Automatic measurement and control of attitude for crane lifters (in Japanese), *Trans. SICE*, 32(1), 97-105, 1996.
15. C. Grubin, Simple algorithm for intersecting two conical surfaces, *J. Spacecraft Rockets*, 14(4), 251-252, 1977.
16. M. A. Stedham and P. P. Banerjee, The panoramic annular lens attitude determination system, *SPIE Proceedings, Space Guidance, Control, and Tracking II*, Orlando, FL, 17-18 April, 1995.
17. J. A. Gilbert, D. R. Matthys, and P. Greguss, Optical measurements through panoramic imaging systems, *Proc. Int. Conf. Hologram Interferometry and Speckle Metrology*, Baltimore, MD, November 4-7, 1990.

10.3 Inertial Navigation

Halit Eren and C. C. Fung

The Principles

The original meaning of the word *navigation* is "ship driving." In ancient times when sailing boats were used, navigation was a process of steering the ship in accordance with some means of directional information, and adjusting the sails to control the speed of the boat. The objective was to bring the vessel from location A to location B safely. At present, navigation is a combination of science and technology. No longer is the term limited to the control of a ship on the sea surface; it is applied to land, air, sea surface, underwater, and space.

The concept of inertial-navigator mechanization was first suggested by Schuler in Germany in 1923. His suggested navigation system was based on an Earth-radius pendulum. However, the first inertial guidance system based on acceleration was suggested by Boykow in 1938. The German A-4 rocket, toward the end of World War II, used an inertial guidance system based on flight-instrument type gyroscopes for attitude control and stabilization. In this system, body-mounted gyro-pendulum-integrating accelerometers were used to determine the velocity along the trajectory. The first fully operational inertial auto-navigator system in the U.S. was the XN-1 developed in 1950 to guide C-47 rocket. Presently, inertial navigation systems are well developed theoretically and technologically. They find diverse applications,

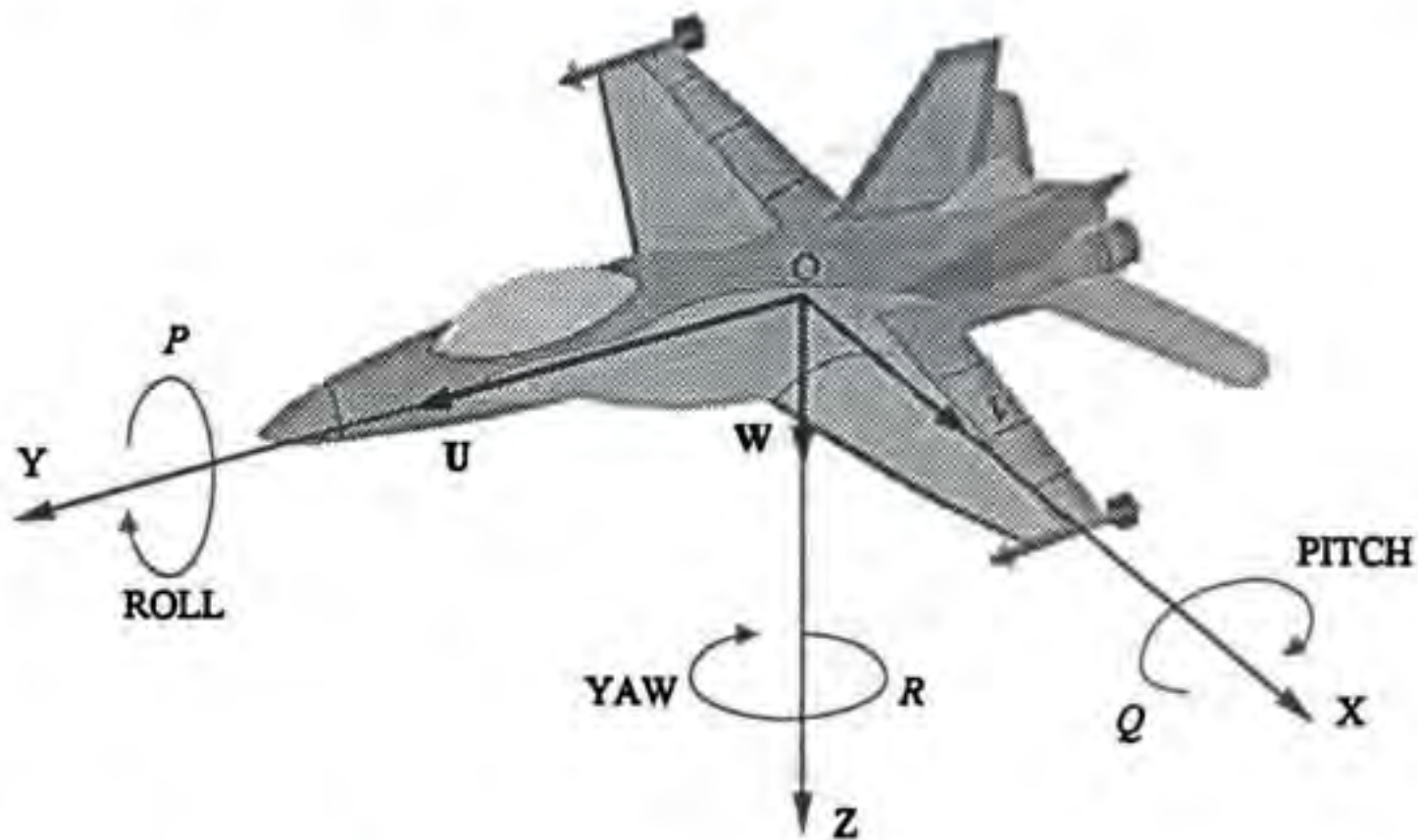


FIGURE 10.17 In inertial navigation, the movement of a vehicle, rocket, ship, aircraft, robot, etc. with respect to a reference axis is monitored. On the Earth's surface, the conventional reference is the Earth's fixed axes — North, East, and Down. A vehicle such as an aircraft or a marine vessel will have its own local axes, known as roll, pitch, and yaw.

allowing the choice of appropriate navigation devices, depending on cost, accuracy, human interface, global coverage, time delay, autonomy, etc.

Inertial navigation is a technique using a self-contained system to measure a vehicle's movement and determine how far it has moved from its starting point. Acceleration is a vector quantity involving magnitude and direction. A single accelerometer measures magnitude but not direction. Typically, it measures the component of acceleration along a predetermined line or direction. The direction information is usually supplied by gyroscopes that provide a reference frame for the accelerometers. Unlike other positional methods that rely on external references, an *inertial navigation system (INS)* is compact and self-contained, as it is not required to communicate to any other stations or other references. This property enables the craft to navigate in an unknown territory.

Inertial navigation can be described as a process of directing the movement of a vehicle, rocket, ship, aircraft, robot, etc., from one point to another with respect to a reference axis. The vehicle's current position can be determined from "dead reckoning" with respect to a known initial starting reference position. On the Earth's surface, the conventional reference will be North, East, and Down. This is referred to as the *Earth's fixed axes*. A vehicle such as an aircraft or a marine vessel will have its own *local axes*: roll, pitch, and yaw, as shown in Figure 10.17.

The inertial sensors of the INS can be mounted in such a way that they stay leveled and pointing in a fixed direction. This system relies on a set of gimbals and sensors attached on three axes to monitor the angles at all times. This type of INS is based on a *navigational platform*. A sketch of a three-axis platform is shown in Figure 10.18. Another type of INS is the *strapdown system* that eliminates the use of gimbals. In this case, the gyros and accelerometers are mounted to the structure of the vehicle. The measurements received are made in reference to the local axes of roll, pitch, and yaw. The gyros measure the movement of angles in the three axes in a short time interval (e.g., 100 samples per second). The computer then uses this information to resolve the accelerometer outputs into the navigation axes. A schematic block diagram of the strapdown system is shown in Figure 10.19.

The controlling action is based on the sensing components of acceleration of the vehicle in known spatial directions, by instruments which mechanize Newtonian laws of motion. The first and second integration of the sensed acceleration determine velocity and position, respectively. A typical INS includes

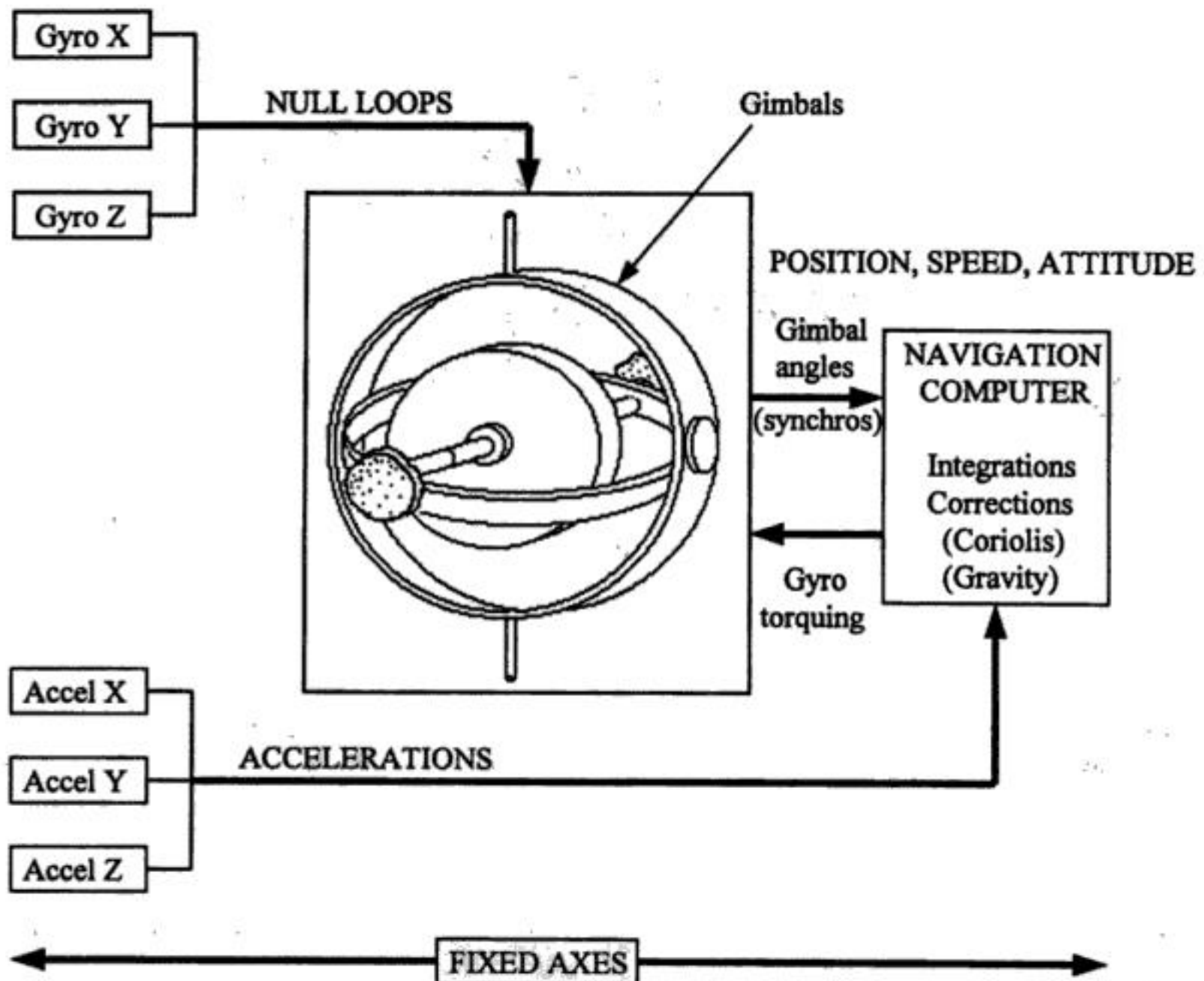


FIGURE 10.18 Some Inertial Navigation Systems, INS, are based on a navigational platform. The inertial sensors are mounted in such a way they can stay leveled at all times, pointing in a fixed direction. This system uses a set of gimbals and sensors attached on three axis in the x , y , and z directions to monitor the angles and accelerations constantly. The navigation computer makes corrections for coriolis, gravity, and other effects.

a set of gyros, a set of accelerometers, and appropriate signal processing units. Although the principle of the systems may be simple, the fabrication of a practical system demands a sophisticated technological base. The system accuracy is independent of altitude, terrain, and other physical variables, but is limited almost purely by the accuracy of its own components. Traditional INSs mainly relied on mechanical gyros and accelerometers, but today there are many different types available, such as optical gyroscopes, piezoelectric vibrating gyroscopes, active and passive resonating gyroscopes, etc. Also, micromachined gyroscopes and accelerometers are making an important impact on modern inertia navigation systems. A brief description and operational principles of gyroscopes and accelerometers suitable for inertial navigation are given below.

Major advances in INS over the years include the development of the *electrostatic gyro* (ESG) and the laser gyro. In ESG, the rotor spins at a speed above 200×10^3 rpm in a near-vacuum environment. The rotor is suspended by an electrostatic field; thus, it is free from bearing friction and other random torques due to mechanical supports. Hence, its operation results in a superior performance compared to others, closely resembling the performance of a theoretical gyro. Although no system can claim to reach perfection, an ESG requires less frequent updates as compared to other mechanical gyros.

Gyroscopes

There are two broad categories: (1) mechanical gyroscopes and (2) optical gyroscopes. Within both of these categories, there are many different types available. Only the few basic types will be described to

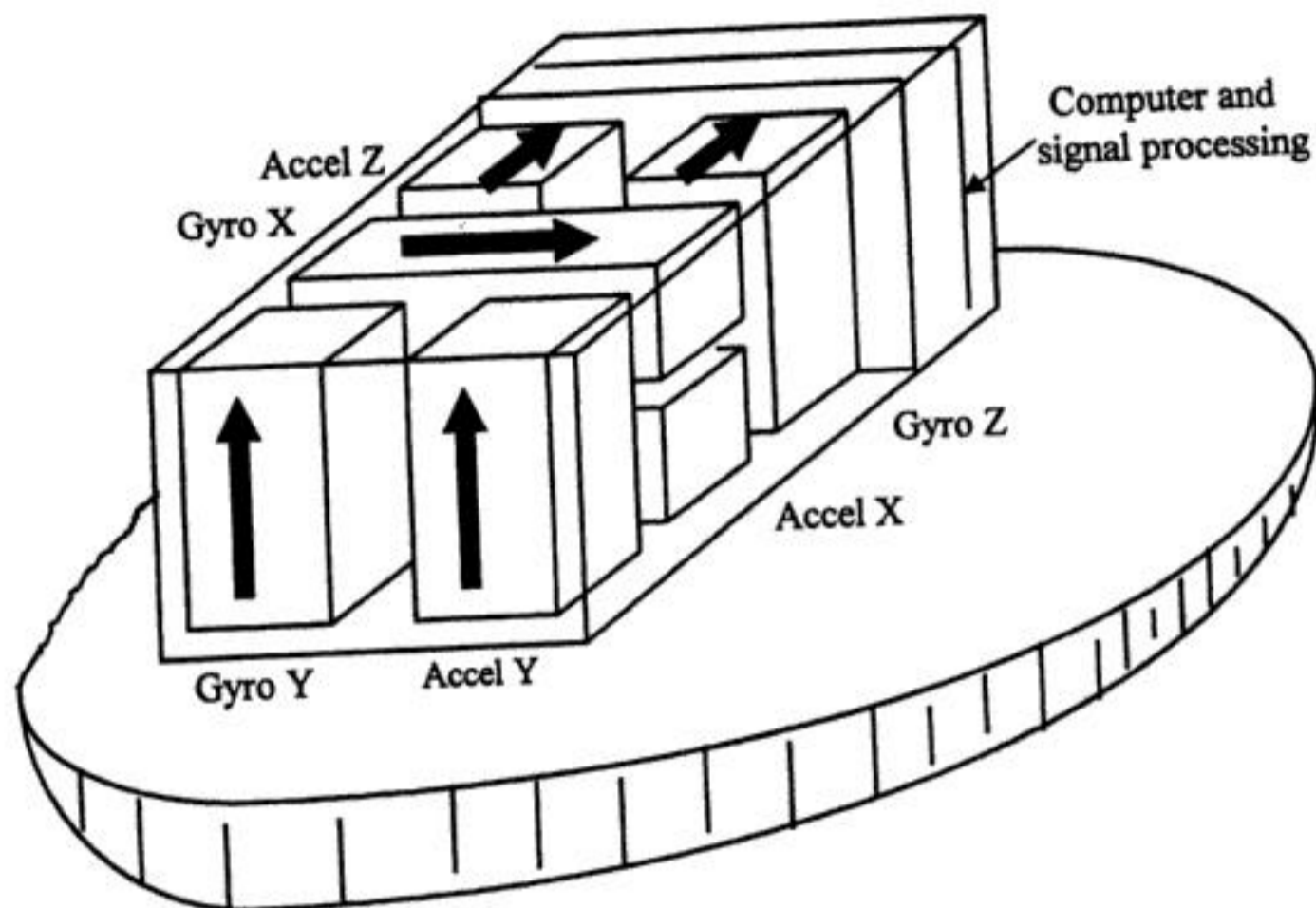


FIGURE 10.19 The use of a strapdown system eliminates the need for gimbals. The gyros and accelerometers are mounted rigidly on the structure of the vehicle, and the measurements are referenced to the local axes of roll, pitch, and yaw. The gyros measure the movement of angles in the three axes in short time intervals to be processed by the computer. This information is used, together with the accelerometer outputs, for predicting navigation axes.

illustrate the operating principles; detailed information may be found in the references listed at the end of this chapter.

Mechanical gyroscopes: The first mechanical gyroscope was built by Foucault in 1852, as a gimballed wheel that stayed fixed in space due to angular momentum while the platform rotated around it. They operate on the basis of *conservation of angular momentum* by sensing the change in direction of an angular momentum. There are many different types, which are:

1. *Single degree of freedom gyroscopes:* include the rate, rate integrating, spinning rotor flywheel, electron, and particle gyros.
2. *Two degree of freedom gyroscopes:* incorporate the external gimbal types, two-axis floated, spherical free-rotor, electrically suspended, gas-bearing free-rotor gyros.
3. *Vibrating gyroscopes:* include the tuning fork, vibrating string, vibrating shell, hemispherical resonating, and vibrating cylinder gyros.
4. *Continuous linear momentum gyroscopes:* incorporate a steady stream of fluid, plasma, or electrons, which tends to maintain its established velocity vector as the platform turns. One typical example is based on a differential pair of hot-wire anemometers to detect the apparent lateral displacement of the flowing air column.

The operating principle of all mechanical gyroscopes is based on the conservation of angular momentum, as shown in Figure 10.20. The angular momentum is important since it provides an axis of reference. From Newton's second law, the angular momentum of a body will remain unchanged unless it is acted upon by a torque. The rate of change of angular momentum is equal to the magnitude of the torque, in vectorial form as:

$$T = dH/dt \quad (10.92)$$

where H = angular momentum (= inertia \times angular velocity, $I\omega$).

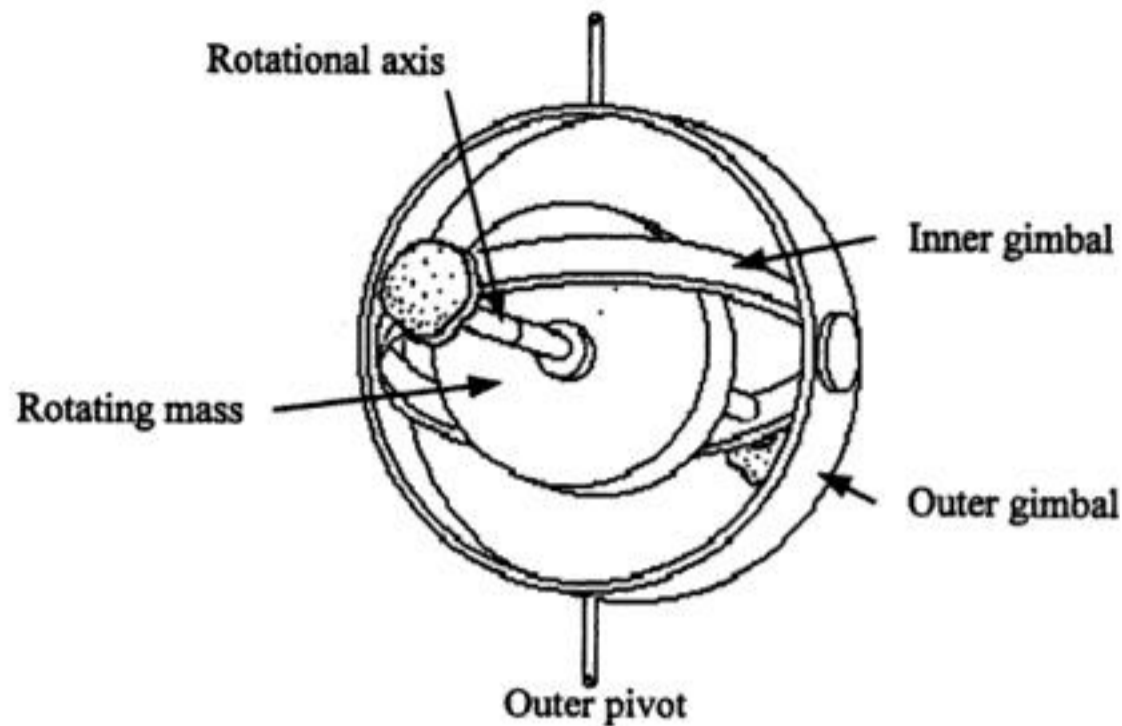


FIGURE 10.21 In a double-axis flywheel gyro, an electrically driven rotor is suspended by a pair of precision low-friction bearings at the rotor axle. The rotor bearings are supported by a circular inner gimbal ring. The inner gimbal ring in turn pivots on a second set of bearings attached to an outer gimbal ring. The pivoting action of the inner gimbal defines the horizontal axis of the gyro, which is perpendicular to the spin axis of the rotor. The outer gimbal ring is attached to the instrument frame by a third set of bearings. This arrangement always preserves the predetermined spin-axis direction in inertial space.

suspension has the property of always preserving the predetermined spin-axis direction in inertial space. Equations governing the two degrees of freedom gyroscope can be written using Equations 10.92 to 10.95. The torque with respect to an inertial reference frame can be expressed as:

$$T = \dot{H}_I \quad (10.96)$$

If the Earth is taken as a moving reference frame, then:

$$\dot{H}_I = \dot{H}_E + \omega_{IE} H \quad (10.97)$$

If the gyroscope itself is mounted on a vehicle (e.g., aircraft) that is moving with respect to the Earth, then:

$$\dot{H}_E = \dot{H}_B + \omega_{EB} H \quad (10.98)$$

The case of the gyroscope can be mounted on a platform so that it can rotate relative to the platform; then:

$$\dot{H}_B = \dot{H}_C + \omega_{BC} H \quad (10.99)$$

Finally, the inner gimbal can rotate relative to the case, hence:

$$\dot{H}_C = \dot{H}_G + \omega_{GC} H \quad (10.100)$$

Substituting Equations 10.97 to 10.100 into Equation 10.96 yields:

$$T = \dot{H}_G (\omega_{GC} + \omega_{BC} + \omega_{EB} + \omega_{IE}) H \quad (10.101)$$

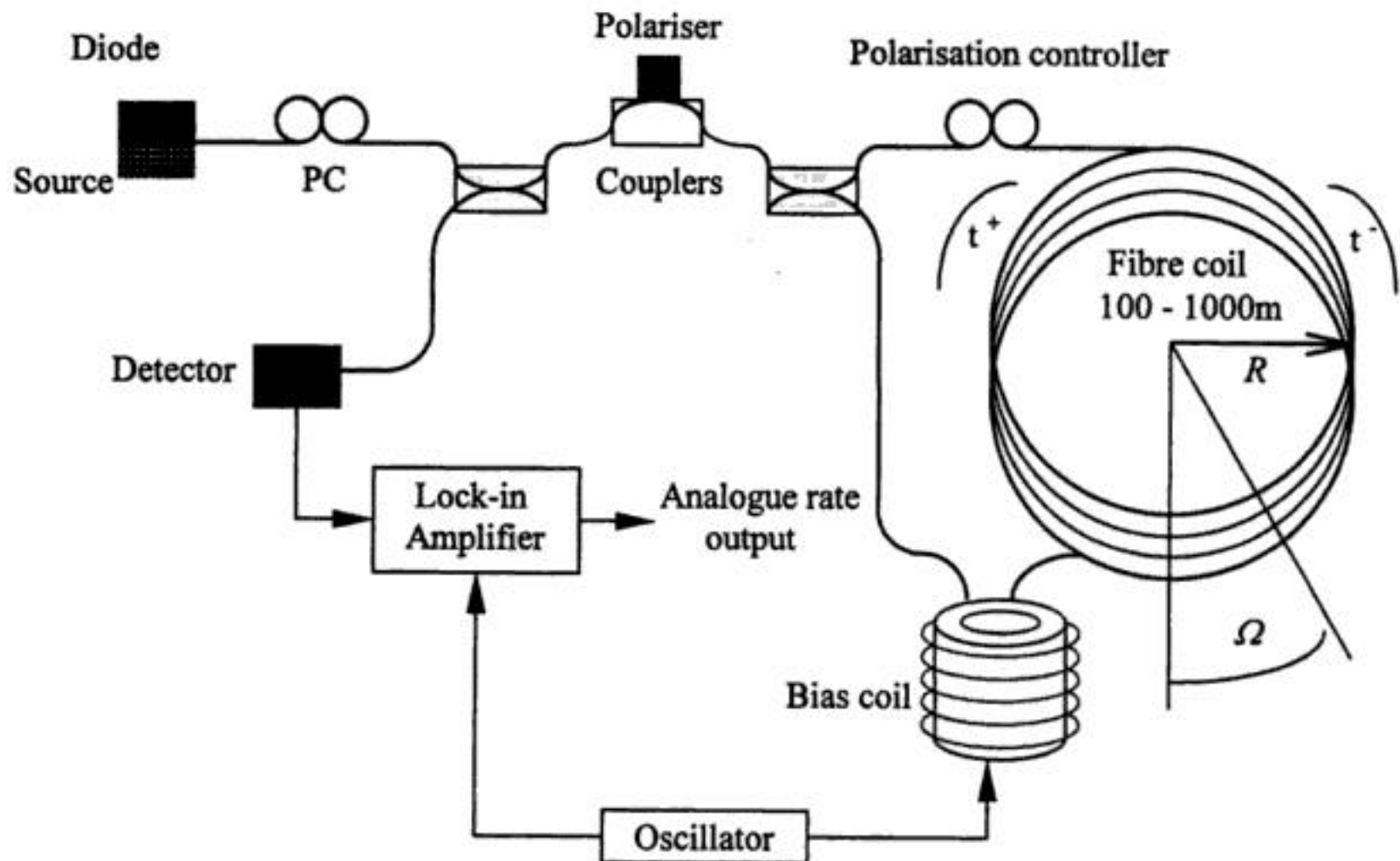


FIGURE 10.22 A typical fiber-optic gyroscope. This gyroscope is based on the inertial properties of light, making use of the Sagnac effect. The Sagnac effect describes interferometer fringe shift against rotation rate. Two light waves circulate in opposite directions around a path of radius R , beginning at source S . When the gyro is stationary, the two beams arrive at the detector at the same time and no phase difference is recorded. If the optical path is rotating with a velocity, the light traveling in the opposite direction to rotation returns to the source sooner than that traveling in the same direction. The two beams interfere to form a fringe pattern and the fringe position may be recorded, or the phase differences of the two beams may be sensed by photodetectors.

The most commonly used accelerometer in navigation systems is based on pendulous types. These accelerometers can be classified as:

1. Generic pendulous accelerometer
2. Q-flex type accelerometers
3. Micromachined accelerometers (A typical example of a modern micromachined accelerometer is given in Figure 10.23.)

Accelerations in the three axes are measured by suitably positioned accelerometers. Since accelerometers contain errors, the readings must be compensated by removing fixed biases or by applying scaling factors. The errors may be functions of operating temperature, vibration, or shock. Measurement of time must be precise as it is squared within the integration process for position determination. The Earth's rotation must also be considered and gravitational effects must be compensated appropriately.

Errors and Stabilization

Errors

In general, inertial navigation is an initial value process in which the location of the navigating object is determined by adding distances moved in known directions. Any errors in the system cause misrepresentation of the desired location by being off-target. The major disadvantage of an inertial guidance system is that its errors tend to grow with time. These errors in the deduced location are due to a number of reasons, including: imperfect knowledge of the starting conditions, errors in computation, and mainly errors generated by gyros and accelerometers.

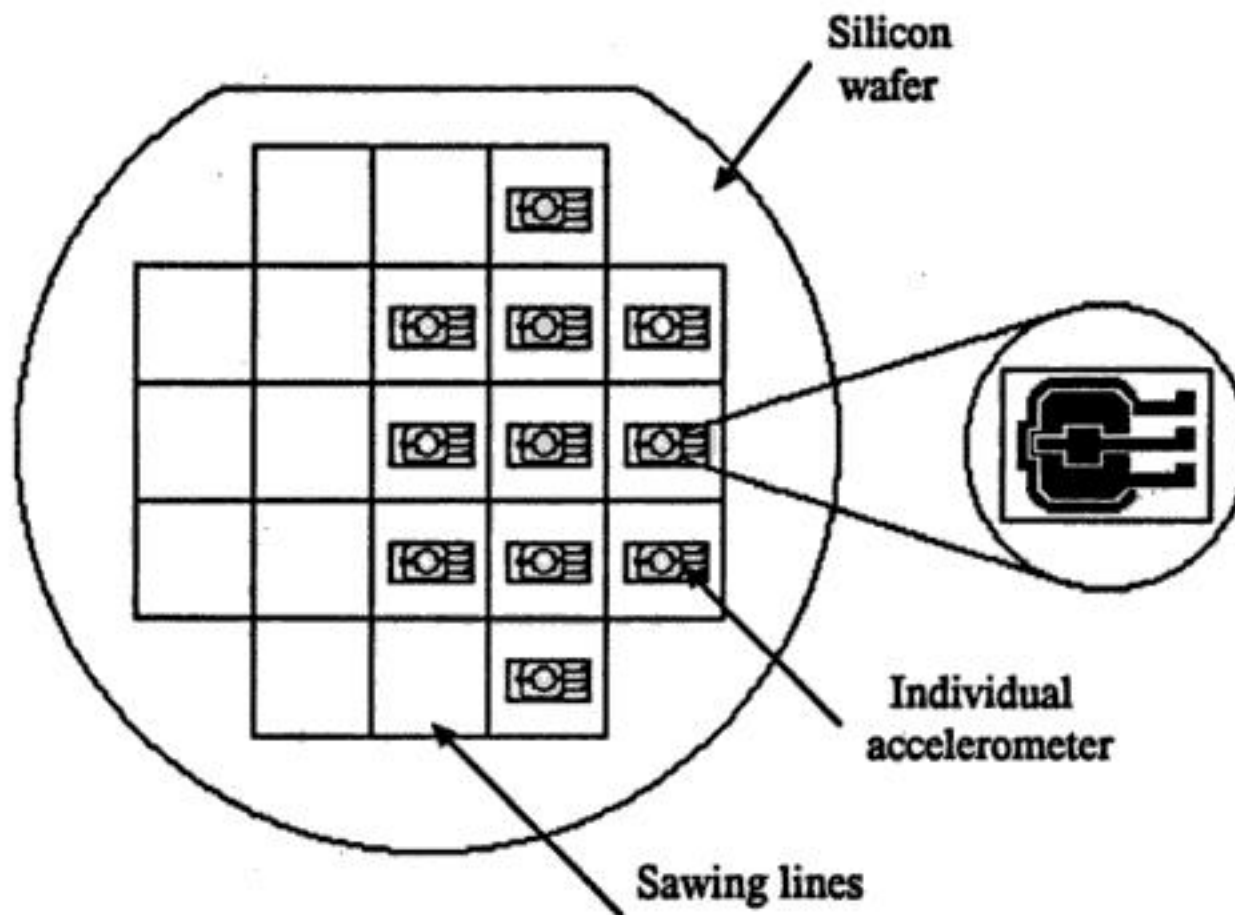


FIGURE 10.23 A typical example of a modern micromachined accelerometer. Multiple accelerometers can be mounted on a single chip, sensing accelerations in the x , y , and z directions. The primary signal conditioning is also provided in the same chip. The output from the chip is usually read in digital form.

If the error build-up with time becomes too large, external aids (e.g., LORAN, OMEGA) may be used to reset or update the system. Optimal use of the data from external aids must account for the geometry of the update and also for the accuracy of the update relative to the accuracy of the inertial system. The Kalman filter, for example, is one of the computational procedures frequently applied for optimally combining data from different sources.

Errors can broadly be classified as:

1. *System heading error*: A misalignment angle in the heading of an object traveling with a velocity can cause serious errors. For example, a vehicle traveling with velocity of 500 km h^{-1} in the same direction with 0.1° initial heading error will be off the target by approximately 873 m at the end of 1 h travel.
2. *Scale error*: Error in scaling can accumulate. In order to minimize scale errors, a scale factor is used. The scale factor is the ratio between changes in the input and output signals. It simply translates the gyro output (counts per second in the case of RLG) into a corresponding angle rotation. Some instruments may have different scale factors for positive and negative inputs, known as *scale factor asymmetry*. (Scale factors are measured in $^\circ \text{ h}^{-1} \text{ mA}^{-1}$, $^\circ \text{ h}^{-1} \text{ Hz}^{-1}$, or $g \text{ Hz}^{-1}$.)
3. *Nonlinearity and composite errors*: In most cases, scale factors are not constant, but they can have second- or higher-order terms relating the output signals to the input. Statistical techniques can be employed to minimize these errors.
4. *Bias errors*: Zero offset or bias error is due to existence of some level of output signal for a zero input. Bias errors exist in accelerometers, gyros, tilt misalignments, etc.
5. *Random drift and random walk errors*: In some cases, the outputs of the devices can change due to disturbances inside the sensors, such as ball bearing noise in mechanical gyros. These disturbances may be related to temperature changes, aging, etc. White noise in optical gyros can cause a long-term accumulation in angle error known as the *random walk*.
6. *Dead band, threshold, resolution, and hysteresis errors*: These errors can be related to inherent operation of accelerometers and gyros. They can be due to stiction, minimum input required for an output, minimum measurable outputs, and nonrepeatability of variations in the output versus variations in the input.

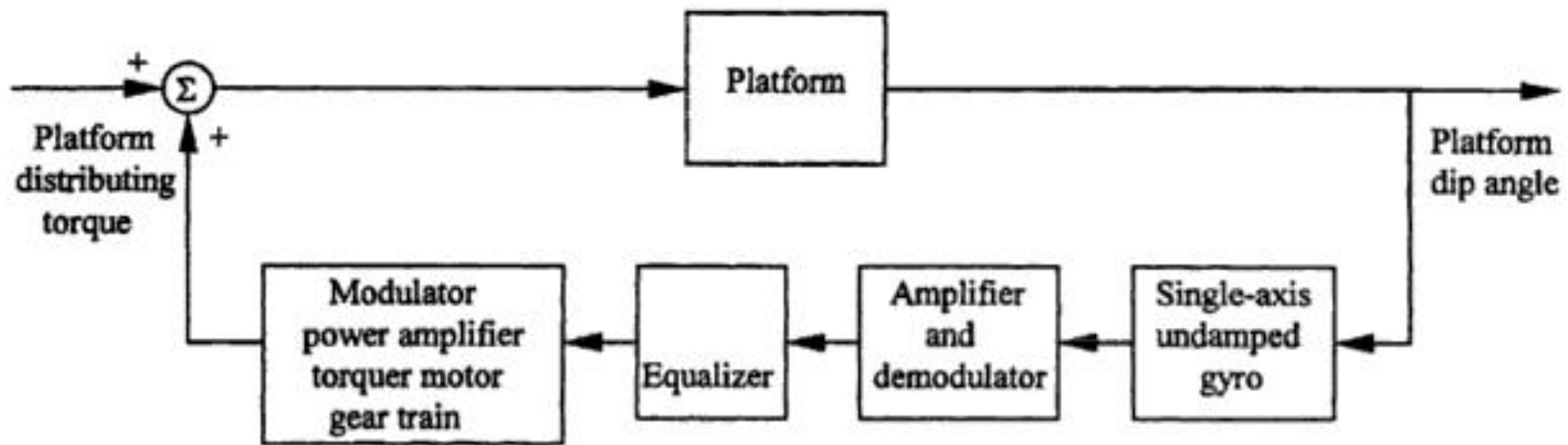


FIGURE 10.24 Stabilization is obtained using platforms designed to accurately maintain accelerometers and gyros leveled and oriented in the azimuth direction. In some cases, the platform is driven around its axis by servo amplifiers and electric motors. Sensitive pick-offs on the gyroscopes fed error signals are used to maintain a desired stability of the platform in the presence of disturbing torques.

It should be pointed out that this list is by no means exhaustive. Detailed error analysis can be found in the references cited.

Stabilization

The inertial navigation sensors must maintain angles within specified limits in spite of the disturbances imposed by the moving object. Accuracy requirements demand that the system must provide reliable and stable information in spite vibrations and other disturbing factors. One way of achieving stabilization is by using a stabilized platform. These platforms are designed to maintain accelerometers and gyros accurately leveled and oriented in the azimuth direction. In some cases, the platform is driven around its axis by servo amplifiers and electric motors. Usually, outputs of doubly integrating accelerometers are used directly to control the level-axis gyroscope precession rates. Sensitive pick-offs on the gyroscopes fed error signals are used to maintain a desired stable platform in the face of disturbing torques. The operation of a typical system, in block diagram form, is shown in Figure 10.24.

Unlike platform models, in a strapped-down system, gyroscopes and accelerometers are rigidly mounted to the vehicle structure so that they move with the vehicle. The accelerometers and gyroscopes are manufactured to measure accelerations and angles up to the maximum expected values. As the vehicle travels, the measured values are frequently transmitted to a computer. The computer uses these values to resolve the readings into the navigation axis sets and make deductions on the body axis sets.

Vehicular Inertial Navigation

In modern vehicular navigation, computerized maps and mobile communication equipment are integrated together with inertial and/or other electronic navigation systems. In recent years, in the wake of low-cost GPS systems, the vehicular navigation system has attracted much attention due to its large potential markets for consumer as well as business vehicles.

Automobile navigation systems are based on dead-reckoning, map matching, satellite positioning, and other navigational technologies. Map intelligent systems achieve high relative accuracy by matching dead-reckoned paths with road geometry encoded in a computerized map. This is also used to perform other functions such as vehicle routing and geocoding. Satellite-based navigation systems achieve high absolute accuracy with the support of dead-reckoning augmentation.

The capabilities and functions of automobile navigation systems depend on:

- Choosing the necessary technology
- Integrating the overall system
- Resolving driver interface
- Providing map data basis
- Coordinating mobile communications

Digital maps and mobile data communications combine together for full usefulness and effectiveness. The navigation systems are greatly enhanced in conjunction with stored digital maps combined with effective communications.

The usefulness of a navigation system is related to the accuracy in position determination. There are a number of methods available with varying accuracy; these include the following:

Dead-reckoning

Dead-reckoning is the process of determining vehicle location relative to an initial position by integrating measured increments and directions of travel. The devices include the odometer, the differential odometer, and a magnetic compass. Gyros and inertial systems prove to have limited applications in harsh automotive environments. Although, dead-reckoning systems suffer from error accumulation, they are widely used inertial navigation systems, particularly in robotics and vehicular applications. Even the most precise navigation system requires periodic reinitialization and continuous calibrations by computers.

Radiolocation

In radiolocation, the global positioning system (GPS) is used extensively. Nevertheless, LORAN is gaining popularity as means of tracking land vehicle location from a central location. But its modest accuracy limits its global application in automotive navigation.

Map Matching

Artificial intelligence concepts are applied to match dead-reckoned vehicle paths, which are stored in computers. In map matching, sensed mathematical features of the vehicle paths are continuously associated with those encoded in a map database. Thus, a vehicle's dead-reckoned location can be initialized automatically at every turn to prevent accumulation of dead-reckoning errors.

The first application of map matching technology was in the Automatic Route Control System (ARCS), which used a differential odometer for dead-reckoning. In another system, the Etak map matching system, a solid-state flux gate compass is used as well as a differential odometer to dead-reckon paths for matching with digitized maps and aerial photographs. Further details on these technologies can be found in the references given at the end of this chapter.

In a map matching system, as long as the streets and road connectivities are accurately defined, the process identifies position relative to the road network as visually perceived by the vehicle driver.

Most of the dead-reckoning equipment commercially available is sufficiently robust to support map matching when operating in a defined road network. However, a good dead-reckoning accuracy is required to achieve reinitialization through map matching upon returning to the road network after off-road operations.

Proximity Beacon

This approach uses strategically located short-range transmitters, and reception of their location coded signal infers the receiving vehicle's instantaneous location. There are several variations of the proximity approach; some versions involve two-way communications with the equipped vehicle. Typically, the driver enters the destination code on the vehicle panel, for automatic transmission to the roadside unit, as the vehicle approaches an instrumented intersection. The roadside unit, which can be networked with a traffic management system, analyzes the destination code and transmits route instructions to the display on the vehicle panel. Proximity beacon systems are being tested in Germany and Japan. One of the most popular system is the ALI-SCOUT (see references) proximity beacon system, which uses dead-reckoning and map matching techniques between beacons to download updated map and traffic data in Berlin.

The approach to the interface between an on-board navigation system and a vehicle operator must take into account ergonomics and safety considerations as well as functional requirements. As a result of intensive research, especially in the aerospace industry, display of information for the operator is a well-developed area. In a well-known European system, Philips' CARIN, a color CRT map display is used to show vehicle location relative to the surroundings. Many other systems use short visual messages, symbolic graphics, and voice.

the other four tones uniquely identify the station. Each station has a cesium clock that is calibrated within 1 μ s by exchanging portable atomic clocks with other Omega stations and with the U.S. Naval Observatory.

GPS systems give all vehicles on or near the Earth unprecedented navigation accuracy. A number of international airlines are equipped with confined GPS-Glonass receivers. Many experiments are now in progress to balance the cost versus accuracy of various combinations of inertial, Omega, Loran, GPS, and Transit equipment. Airborne receivers are designed that combine nav aids operating in a common radio band (e.g., GPS, DME, JTID, and IFF).

INMARSAT's "Marec" communication satellites serve ship traffic but are configured to serve air traffic by directly communicating with aircraft for approximately \$10 per call.

Underwater

The Ship's Inertial Navigational System (SINS) was originally developed for precision position-finding required by ballistic missile submarines in the late 1950s and early 1960s. The first deployment was on-board U.S. George Washington in 1960, and SINS are used today in submarines, aircraft carriers, and other surface warships. As the cost and size are continually decreasing, the system is also deployed in naval as well as merchant vessels. Another development of INS for underwater application is in the area of the autonomous underwater vehicle (AUV). In this section, a number of such products are described.

AUVs are used extensively for military and civilian purposes. Application examples are mapping, surveillance, ocean exploration, survey, and mining, all of which require precise position determination. The desired features of such systems are: low power, high accuracy, small volume, light weight, and low cost. Two typical examples are the LN family produced by Litton Guidance and Control Systems, and the system developed by the Harbor Branch Oceanographic Institution Inc. (HBOI). The specifications of some of these systems are briefly described below to give examples of underwater INS.

The Litton LN-100 System

Litton's LN-100 is an example of the strapdown INS. The LN-100 system consists of three Litton Zero-Lock Gyros (ZLG), a Litton A4 accelerometer triad, power supply, supporting electronics, and a JIWAG standard 80960 computer. The single-board computer performs all the control, navigation, and interface functions.

The HBOI System

The HBOI system was developed with Kearfott Guidance and Navigation (KGN) and utilizes a Monolithic Ring Laser Gyroscope (MRLG), motion sensors, GPS input, and control unit. The inertial measurement unit is based on the Kearfott's T16-B three-axis ring laser gyro and three accelerometers.

Robotics

Closely related to the autonomous underwater vehicles, autonomous mobile robots also use INS extensively as a self-contained, independent navigation system. Typical applications are mining, unknown terrain exploration, and off-line path planning. There are many commercially available inertial navigation systems suitable for cost-effective utilization in the navigation of robots. Some of these are: gyrocompasses, rate gyros, gyrochip, piezoelectric vibrating gyros, ring laser gyros, interferometric, and other types of fiber-optic gyros. Three popular systems will be explained here.

The Honeywell Modular Azimuth Position System (MAPS)

Honeywell's H-726 Modular Azimuth Position System (MAPS) is a typical example of an inertial navigation system for land-based vehicles. It consists of a Dynamic Reference Unit (DRU) that provides processed information from the inertial sensor assembly, a Control Display Unit (CDU) that is used for human-machine interface, and a Vehicle Motion Sensor (VMS) that monitors the vehicle's directional and distance information. The inertial sensor assembly comprises three Honeywell GG1342 ring-laser

Schwartz Electro-Optics, Inc.
3404 N. Orange Blossom Trail
Orlando, FL 32804
Tel: (407) 298-1802
Fax: (407) 297-1794

Siemens Co.
1301 Avenue of the Americas
New York, NY 10019

Southern Avionics Co.
5000-T Belmont
Beaumont, TX 77707
Tel: (800) 280-0322
Fax: (409) 842-2987

Sperry Marine Inc.
1070 T Seminole Tr.
Charlottesville, VA 22901
Tel: (804) 974-2000
Fax: (804) 974-2259

Systron Donner Inertial
Division
BEI Electronics
2700 Systron Drive
Concord, CA 94518-1399
Tel: (510) 682-6161
Fax: (510) 671-6590

Trackor Inc.
6500-T Trackor Lane
Austin, TX

Warren-Knight Inst. Co.
2045 Bennet Dr.
Philadelphia, PA 19116
Tel: (215) 484-9300

10.4 Satellite Navigation and Radiolocation

Halit Eren and C. C. Fung

Modern electronic navigation systems can be classified by range, scope, error, and cost. The range classifications are short, medium, and long ranges, within which exact limits are rather indefinite. The scope classifications can be either self-contained or externally supported, and active (transmitting) or passive (not transmitting) mode of operation.

Short-range systems include radiobeacons, radar, and Decca. Medium-range systems include Decca and certain types of extended-range radars. The long-range systems include Loran-C, Consol, and Omega. All these systems depend on active radio frequency (RF) transmissions, and all are externally supported with respect to the object being navigated, with the exception of the radar. In addition to these, there is another category of systems which are called *advanced navigation systems*; the transit satellite navigation systems, Glonass, and the Global Positioning Systems (GPS) are typical examples.

Utilization of electromagnetic radio waves is common to all navigation systems discussed here. Understanding of their behavior in the Earth's atmosphere is very important in the design, construction, and use of all kinds of navigation equipment — from satellites to simple hand-held receivers.

When an FM radio wave is generated within the Earth's atmosphere, the wave travels outward. The waves may be absorbed or reflected from surfaces of materials they encounter. The absorption and scattering of electromagnetic waves take place for many reasons, one of which is caused by excitation of electrons within the molecules in the propagation media. The behavior of an electromagnetic wave is dependent on its frequency and corresponding wavelength. Figure 10.26 shows the frequency spectrum of electromagnetic waves. They are classified as *audible waves* at the lower end of the spectrum, *radio waves* from 5 kHz to 300 GHz, and *visible light* and various other types of rays at the upper end of the spectrum.

For practical purposes, the radio wave spectrum is broken into eight bands of frequencies; these are: *very low frequency* (VFL) less than 30 kHz, *low frequency* (LF) 30 kHz to 300 kHz, *medium frequency* (MF) 300 kHz to 3 MHz, *high frequency* (HF) 3 MHz to 30 MHz, *very high frequency* (VHF) 30 MHz to 300 MHz, *ultra high frequency* (UHF) 300 MHz to 3 GHz, *super high frequency* (SHF) 3 GHz to 30 GHz, and *extremely high frequency* (EHF) 30 GHz to 300 GHz.

For easy identification, the frequencies above 1 GHz are further broken down by letter designators, as: L-band (1–2 GHz), S-band (2–4 GHz), C-band (4–8 GHz), X-band (8–12.5 GHz), and K-band (12.5–40 GHz). Since severe absorption of radar waves occurs near the resonant frequency of water vapor at 22.2 GHz, the K-band is subdivided into lower K-band (12.5–18 GHz) and upper K-band (26.5–40 GHz). Most navigation radars operate in the X- and S-bands, and many weapons fire control radars operate in the K-band range.

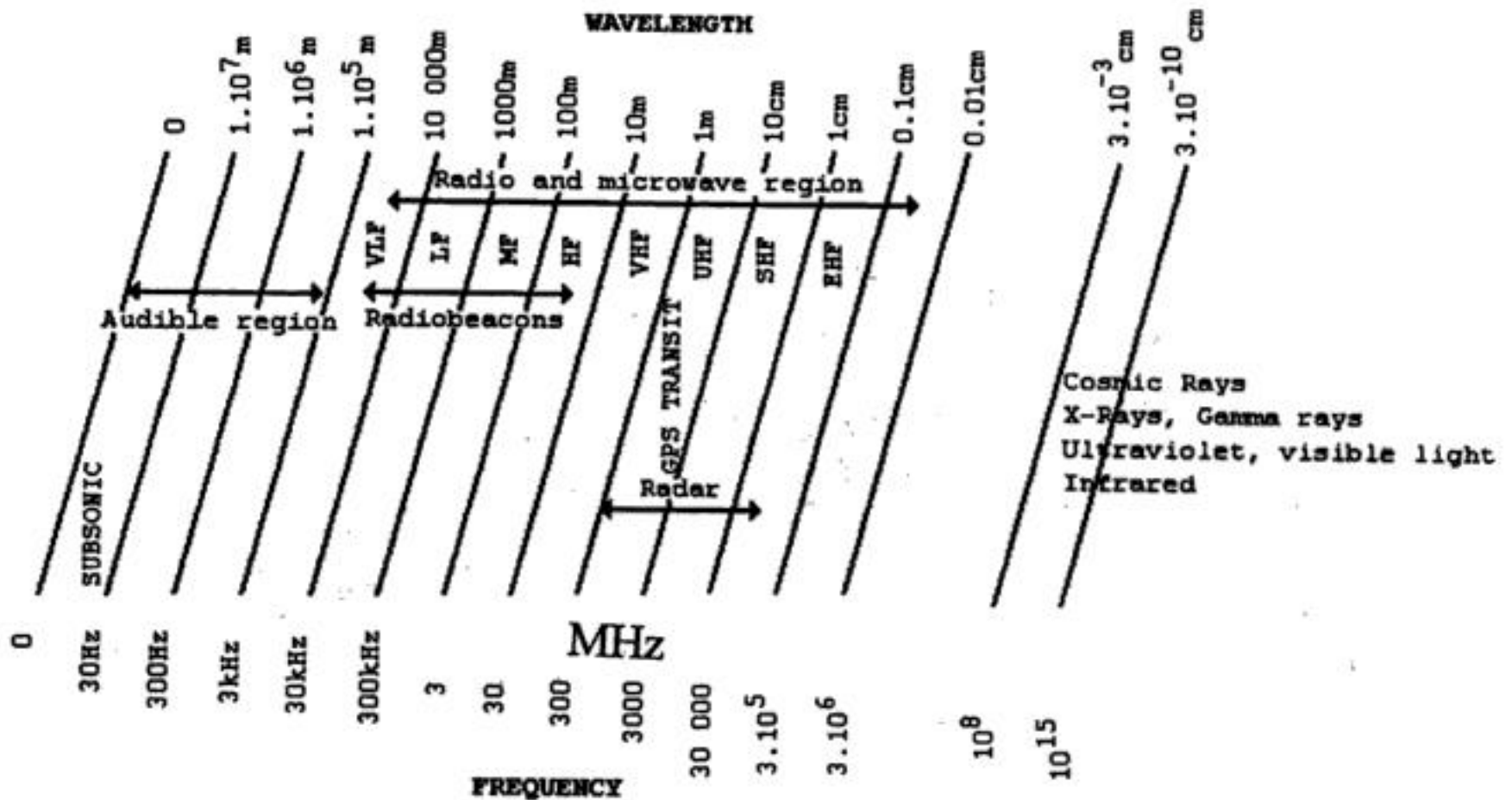


FIGURE 10.26 Electromagnetic wave frequency spectrum. Audible range can be heard if converted to sound waves. Radiobeacons operate in the VLF, LF, and MF ranges. Omega operating at VLF covers the entire world with only eight transmission stations. GPS, Transit, and Glonass use UHF frequencies. Wavelengths less than 10 cm are not suitable for satellite systems, but they are used in radars.

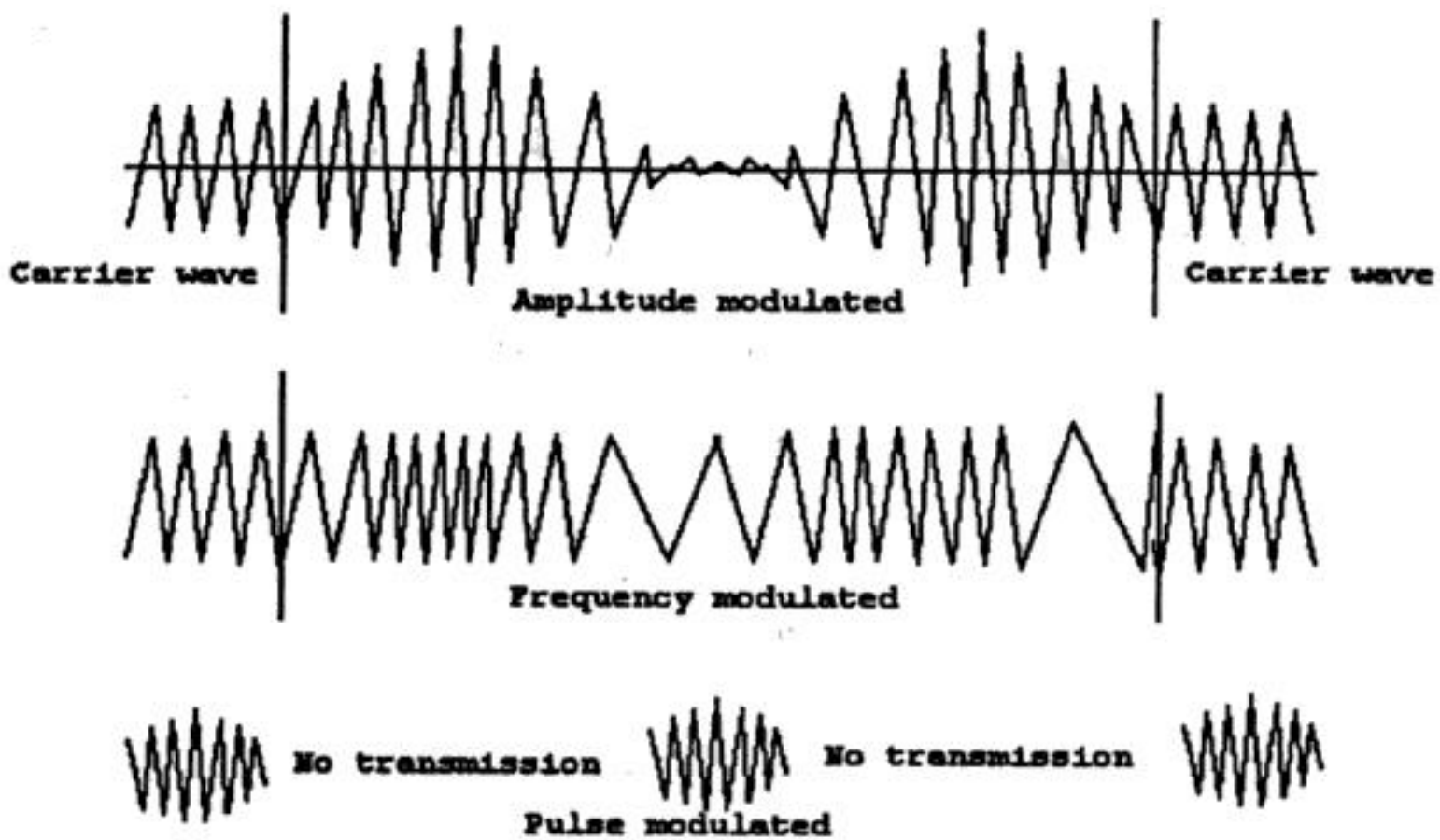


FIGURE 10.27 Amplitude, frequency, and pulse modulation of RF carrier waves. Amplitude modulation is suitable for broadcasting radio stations. Frequency modulation is used in commercial radio broadcasts. The pulse modulation is used in satellite systems, radars, and long-range navigation aids.

The radio waves are transmitted as continuous or modulated waves. A carrier wave (CW) is *modulated* to convey information in three basic forms: amplitude, frequency, and pulse modulation, as shown in Figure 10.27. The amplitude modulation (AM) modifies the amplitude of the carrier wave with a modulating signal. In frequency modulation (FM), the frequency of the carrier wave is altered in accordance with the frequency of the modulating wave. FM is used in commercial radio broadcasts and the sound

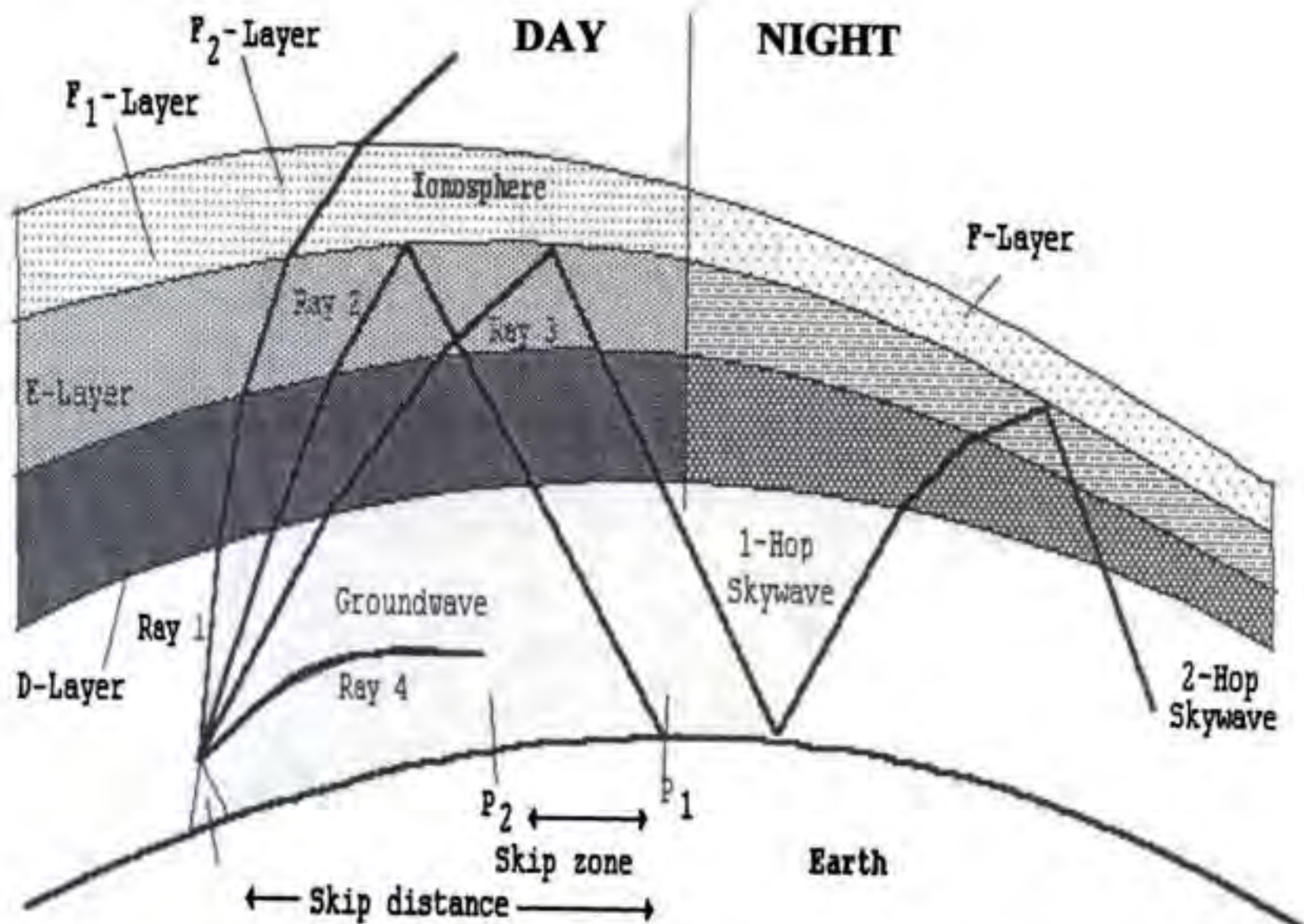


FIGURE 10.28 The four layers of the ionosphere and its effect on radio propagation. The four layers are produced by the ionization of molecules of particles in the atmosphere by ultraviolet rays of sun. The effect of ionosphere on the radio waves is shown by reflections, also termed hops. The frequency of the electromagnetic wave is important in its behavior through ionosphere.

portion of the television broadcast. Pulse modulation is different from AM and FM in that there is usually no impressed modulation wave employed. In this modulation, the continuous wave is broken up into very short bursts or "pulses," separated by periods of silence during which no wave is transmitted. This is used in satellite navigation systems, surface search radars, and long-range radio navigation aids such as Loran.

When an electromagnetic wave encounters an obstruction, *diffraction* takes place marked by a weak reception of the signal within the "shadow" zone. Two waves acting on the same point will also result in *interference*. The degree of interference depends on the phase and frequency relationship. For example, two waves of the same frequency with a 180° phase difference will result in a null at that point. Also, under certain conditions, a portion of the electromagnetic energy in radio waves may reflect back toward the Earth's surface to form the *ionosphere*. The ionosphere is a layer of charged particles located about 90 to 400 km high from Earth's surface; such reflected waves are called *sky waves*.

In the study of radio wave propagation, there are four *ionosphere layers* of importance, as shown in Figure 10.28. The D-layer is located about 60 km to 90 km and is formed during daylight. The E-layer is about 110 km. It persists through the night with decreased intensity. The F₁-layer is between 175 km and 200 km; it occurs only during daylight. The F₂-layer is between 250 km and 400 km; its strength is greatest in the day but it combines with the F₁-layer later to form a weak F-layer after dark. The layers in the ionosphere are variable, with the pattern seeming to have diurnal, seasonal, and sun spot periods. The layers may be highly conductive or may entirely hinder transmissions, depending on the frequency of the wave, its angle of incidence, height, and intensity on various layers at the time of transmission. In general, frequencies in the MF and HF bands are most suitable for ionosphere reflections during both day and night.

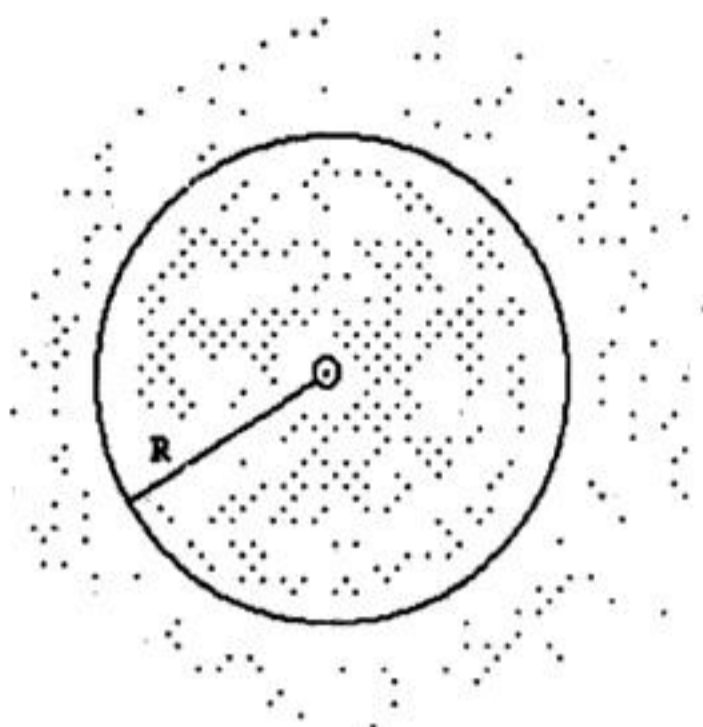


FIGURE 10.29 The rms radius circle that encompasses 68% of all measured positions. The variations in the measurements are due to a number of factors, including: ionosphere conditions, precise location of satellites, and inefficiencies in electronic circuits. (2 rms encompasses 95% of all indicated positions.)

Because of the higher resistance of the Earth's crust as compared to the atmosphere, the lower portions of radio waves parallel to the Earth's surface are slowed down, causing the waves to bend toward Earth. A wave of this type is termed a *ground wave*. The ground waves exist because they use the Earth's surface as a conductor. They occur at low frequency since LF causes more bending in conformity to Earth's shape. The ultimate range of such ground waves depends on the absorption effects. Sometimes, in the lower atmosphere, *surface ducting* occurs by multiple hopping, thus extending the range of a ground wave well beyond its normal limits. It is associated with higher radio and radar frequencies. This phenomenon is common in tropical latitudes. Behavior patterns of waves transmitted at various angles are illustrated in Figure 10.28.

Accuracy of Electronic Fix

There are a number of random effects that influence the accuracy of an electronic position determination; atmospheric disturbances along the transmission path, errors in transmitters and receivers, clocks, inaccuracy in electronic circuitry, gyro errors, etc. As a result, a series of positions determined at a given time and location usually results in a cluster of points near the true position. There are two measures commonly used to describe the accuracy: the first is the *circular error probable (CEP)* — a circle drawn on the true position whose circumference encompasses 50% of all indicated positions, and the second technique, more common, is the *root mean square (rms)*, where:

$$\text{rms} = \sqrt{\sum_{n=1}^N (E_n)^2 / N} \quad (10.105)$$

where E = the distance between actual and predicted positions
 N = the number of predicted positions

A circle, shown in Figure 10.29, with one rms value is expected to contain 68% of all the indicated positions. Another circle of radius equal to 2 rms should contain 95% of all the indicated positions, for isotropic scattering, or errors.

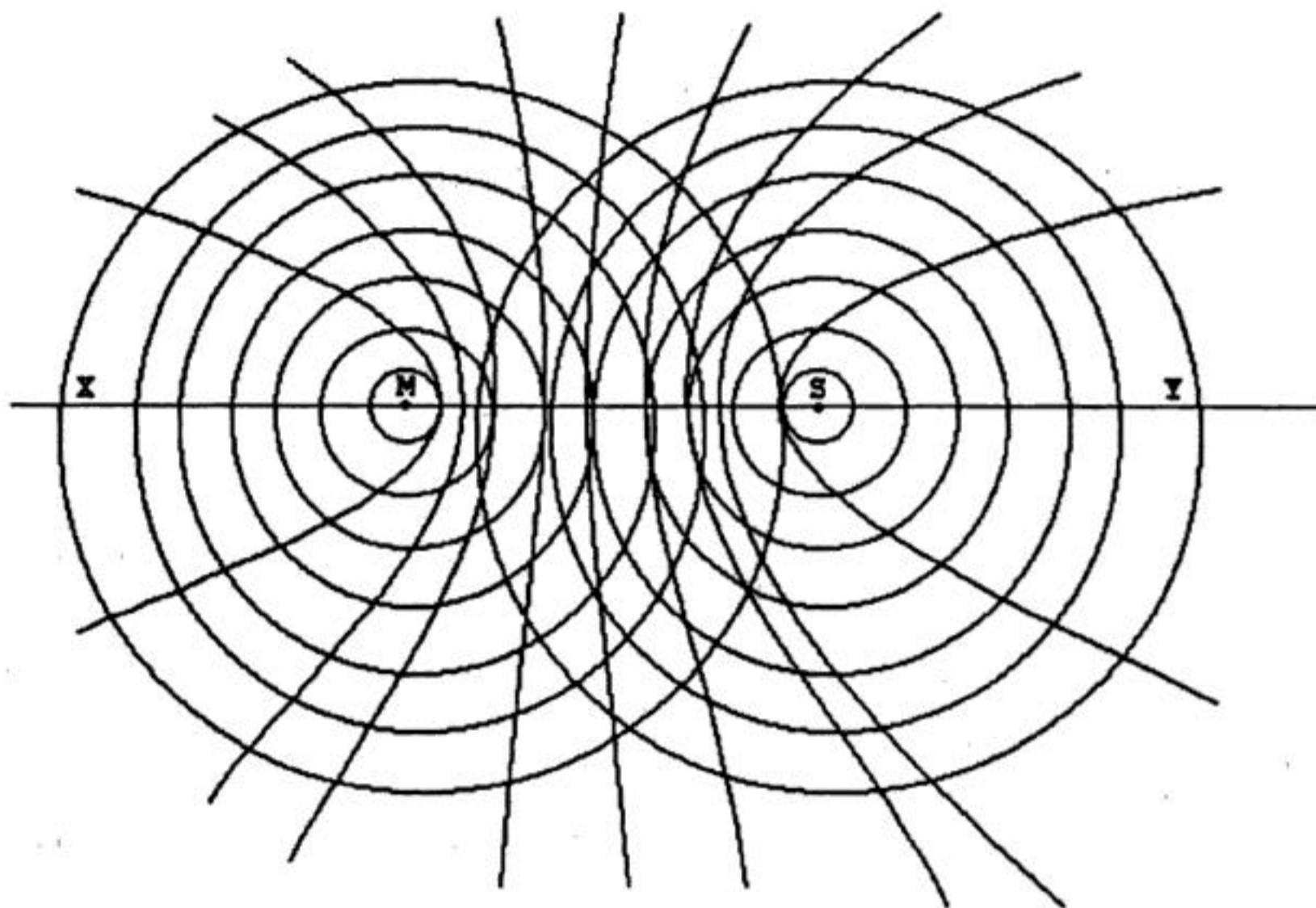


FIGURE 10.30 Hyperbolic patterns of two interacting radio waves propagated in opposite directions. These lines represent the locus of all points of a specified time difference between master and secondary pulses. The straight line MS is called the baseline. The maximum distance of the target object should not exceed 6 times the length of the baseline.

considered accurate only to within $\pm 2^\circ$; for example, under 200 km (120 miles) to the transmitter in favorable conditions and $\pm 5^\circ$ to 10° when conditions are unfavorable.

Information on the locations, ranges, and using procedures of radio beacons are given in a number of publications such as DMAHTC Publication No. 117 *The Radio Navigation Aids*. Correct radiobeacon bearings are plotted on Mercator charts for position estimation. Because of possible inaccuracies, radiobeacons are not used universally. Navigators such as small boats and merchant ships not equipped with other systems use radiobeacons.

Loran-C

Loran was developed in the 1940s to be one of the first systems implementing a long-range hyperbolic system for both ships and aircraft. The system used master and slave stations transmitting sequential radio waves in the upper MF band with frequencies between 1850 kHz and 1950 kHz. Loran-A featured ground wave coverage out to between 700 km and 1250 km from the baseline by day, and up to 2200 km by night. It was the most widely used electronic navigation system until 1970. Later, a system employing synchronized pulses for both time-difference and phase comparison measurements was developed, known as Loran-C. Loran-C was configured to operate in a *chain* form consisting of more than one slave station usually located in triangles.

All stations in the system transmit a signal on a common carrier frequency in mid-LF band of $100 \text{ kHz} \pm 10 \text{ kHz}$. The ground wave range is about 1900 km. One-hop sky waves have a range of about 3600 km, and two-hop signals were noted to have been received about 6250 km from the ground station. One-hop sky waves are produced both by day and by night, while two-hop sky waves are formed only at night. Present Loran-C chains have baseline distances between 1500 km and 2300 km. The accuracy of the system varies from about $\pm 200 \text{ m rms}$ near the baseline to $\pm 600 \text{ m rms}$ near the extreme ranges

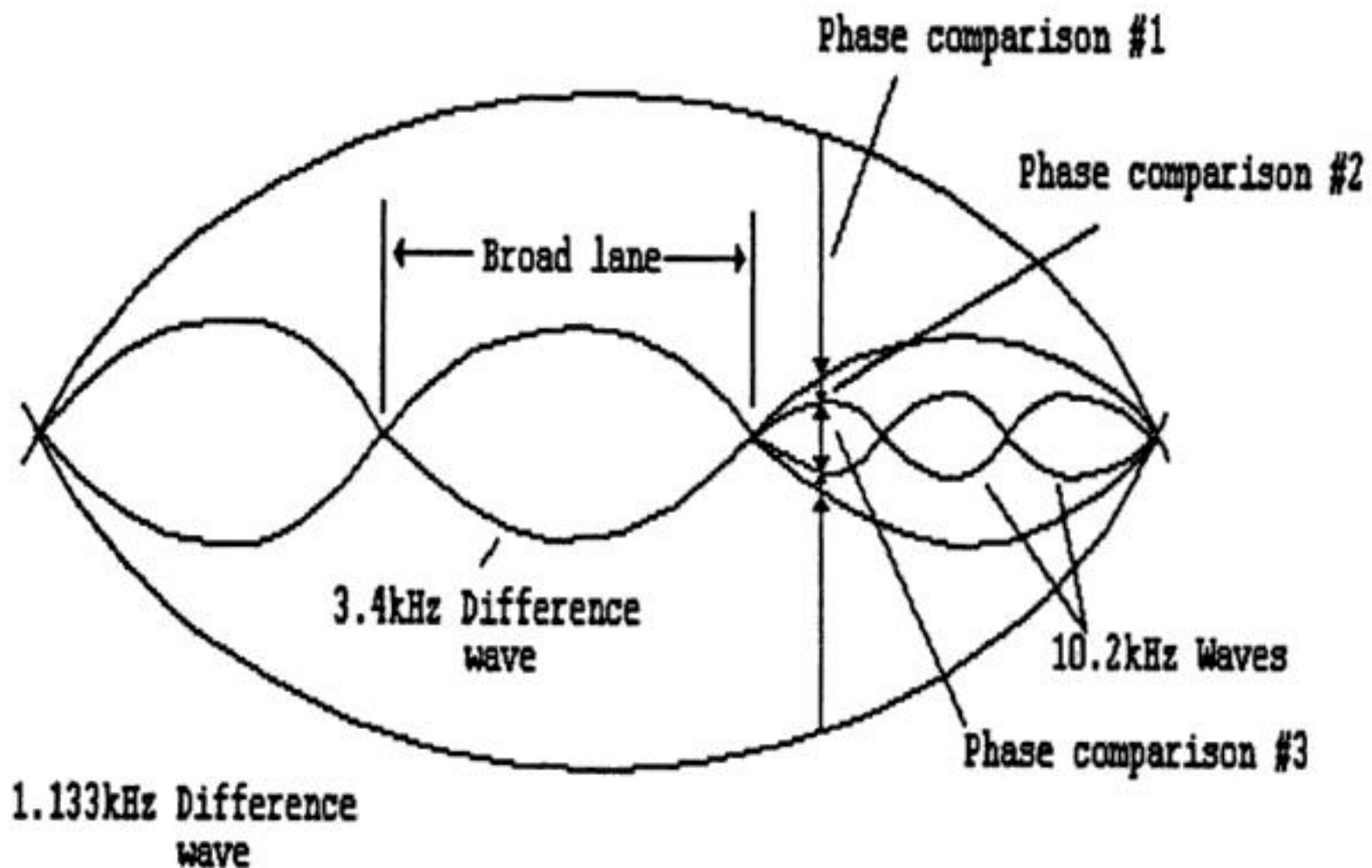


FIGURE 10.32 Three successive phase comparisons for lane resolution in Omega systems. Phase differences are compared in three stages with respect to three different signals transmitted by each station for accurate position finding. One wavelength is 25 km, representing two lanes. Accuracy of the Omega system is limited.

short signal interruptions and transmissions of harmonics simultaneously, zones and lanes are identified in a precise manner.

Consol

Consol is limited to the Eastern and Northern Atlantics. It is a hyperbolic system with extremely short baseline lengths, such that a collapsed hyperbolic pattern is formed. It employs three towers located three wavelengths apart. The operational frequencies are in the MF range between 250 kHz and 370 kHz. The system range is about 1900 km by day and 2400 km by night. The minimum range is about 40 km near the stations. One tower transmits a continuous wave, while other towers transmit waves with 180° phase shift by a *keying cycle*. The signals are modulated by dots and dashes such that receivers determine the position by counting them and printing on Consol grid patterns.

Omega

Omega is a hyperbolic navigation system that covers the entire world with only eight transmission stations located 7500 km to 9500 km apart. It transmits on frequencies in the VLF band from 10 kHz to 14 kHz at a power of 10 kW. The signals of at least three and usually four stations can be received at any position on Earth.

The 10 kHz to 14 kHz frequency band was chosen specifically to take advantage of several favorable propagation characteristics, such as: (1) to use the Earth's surface and ionosphere as a waveguide; (2) to enable submerged submarines to receive the signals; and (3) to form long baselines at 7500 km to 9500 km.

The basic frequency at which all eight stations transmit is 10.2 kHz. Each station transmits four navigation signals as well as a timing signal with atomic frequency standards ensuring that all stations are kept exactly in phase. Two continuous waves are in phase but traveling in opposite directions to produce a series of Omega lanes. Within each lane, a phase difference measurement would progress from 0° to 360° as the receiver moves across, as shown in Figure 10.32. Two Omega lanes complete one cycle, giving a wavelength of 25 km and lane of 12 km expanding as the distance from the baseline increases. Lanes are identified by three other signals transmitted by each station on a multiplexed basis. Omega

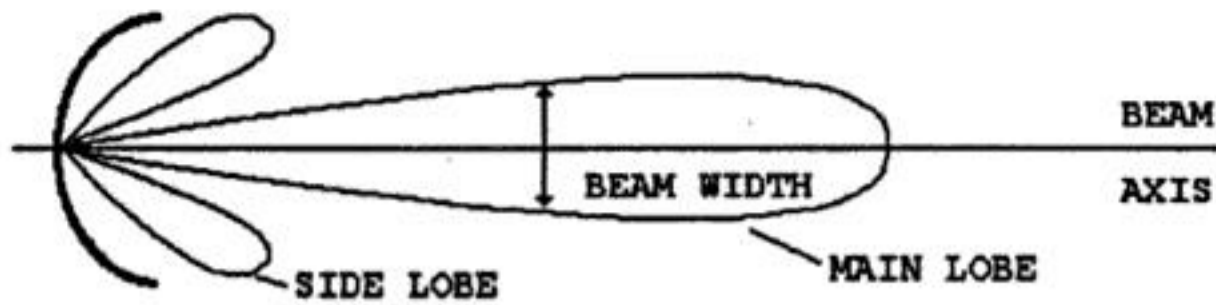


FIGURE 10.33 A surface search radar beam. High-frequency electromagnetic waves are formed by parabolic antenna. The receiving antenna is rotated 360° to scan the entire surrounding area. The location of the target is determined by the reflected back signals and the orientation of the antenna.

fixes have been accurate to within ± 1.5 km rms by day and ± 3 km rms by night. Differential techniques can greatly reduce this error.

Omega receivers are fully automated sets that provide direct lat/long readout for a cost of around U.S. \$1,000 for marine and aviation use. Nevertheless, since the precision of the system is not as good as others, they are mainly used as back-up systems.

Radar

The word is derived from *radio detection* and *ranging*. It works on the basic principle of reflection, which determines the time required for an RF pulse to travel from a reference source to a target and return as an echo. Most surface search and navigation radar high-frequency electromagnetic waves are formed by a parabolic antenna into a beam form, as shown in Figure 10.33. The receiving antenna is rotated to scan the entire surrounding area, and the bearings to the target are determined by the orientation of the antenna at the moment the echo returns. A standard radar is made up of five components: transmitter, modulator, antenna, receiver, and indicator. They operate on pulse modulation.

Radars are extremely important devices for air control applications. Nowadays, airborne beacon radar systems are well developed in traffic alert and collision avoidance systems (TCAS). In this system, each plane constantly emits an interrogation signal, which is received by all nearby aircraft that are equipped appropriately. The signal triggers a transponder in the target aircraft, which then transmits some information concerning 3-D location and identification.

Satellite Relay Systems

The use of satellites is a highly developed technology utilized extensively throughout the world. In the past 2 decades, it has progressed from quasi-experimental in nature to one with routine provisions of new services. They take advantage of the unique characteristics of *geostationary satellite orbits* (GSO). The design of satellite systems is well understood, but the technology is still dynamic. The satellites are useful for long-distance communication services, for services across oceans or difficult terrain, and point-to-multipoint services such as television distribution.

Frequency allocation for satellites is controlled by the International Telecommunication Union (ITU). In the U.S., the Federal Communications Commission (FCC) makes the frequency allocations and assignments for nongovernment satellite usage. The FCC imposes a number of conditions regarding construction and maintenance of in-orbit satellites.

There are many satellite systems operated by different organizations and different countries mainly developed for communications and data transmissions; these include: Iridium of Motorola, Globalstar of Loral Corporation, Intelsat, CS-series of Japan, Turksat of Turkey, Aussat of Australia, Galaxy and Satcom of the U.S., Anik of Canada, TDF of France, etc. Some of the communication satellite systems are suitable for navigation purposes. However, satellite systems specifically designed for navigation are limited in number. The most established and readily accessible by civilian and commercial users are the GPS system of the U.S. and the Glonass of Russia.

The first generation of the satellite system was the *Navy satellite system* (Navsat), which became operational in January 1964, following the successful launch of the first transit satellite into polar orbit.

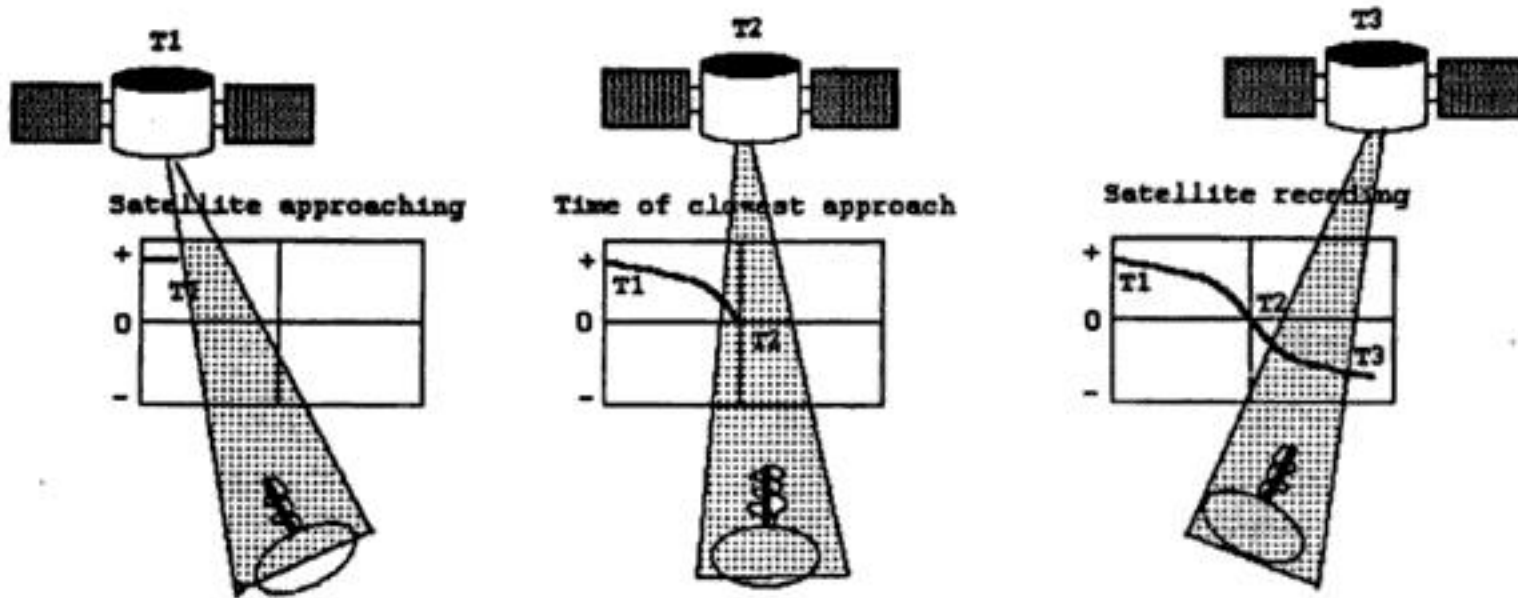


FIGURE 10.34 Transit satellite Doppler curve. As the satellite approaches the receiver, the frequency of the received signal increases due to Doppler shift. At the time of closest approach, the transmitted and received frequencies are the same. The frequencies received from a receding satellite result in lower values. This is also applicable in other position sensing satellites such as GPS, Glonass, Starfix, etc.

The system was declared open for private and commercial use in 1967. Civil designation of the name of the system is *Transit Navigation Satellite System*, or simply *Transit*. Later, this system evolved to become the modern Navsat GPS system, which will be discussed in detail. Most of the operational principles discussed here are inherited by the GPS system.

The Transit system consists of operational satellites, plus several orbiting spares, a network of ground tracking stations, a computing center, an injection station, naval observatory time signals, and receiver-computer combinations. The transit satellites are in circular polar orbits about 1075 km above ground with periods of revolution of about 107 min. Because of the rotation of the Earth beneath the satellites, every position on Earth comes within range of each satellite at least twice a day, at 12 h intervals. As originally intended, if at least five satellites are operational at any given time, the average time between fix opportunities would vary from about 95 min near the equator to about 35 min or less above 70° North and South.

The Transit system is based on the Doppler shift of two frequencies, 150 MHz and 400 MHz, transmitted simultaneously by each satellite moving its orbit at a tangential velocity of about 7.5 km s^{-1} . Two frequencies are used so that the effects of the ionosphere and atmospheric refraction on the incoming satellite transmission can be compensated for by the receivers. Each frequency is modulated by a repeating data signal lasting 2 min, conveying the current satellite time and its orbital parameters and other information. Within the receiver, a single Doppler signal is created by processing the two signals transmitted by the satellite. By plotting the frequency of this signal versus time, a characteristic curve of the type shown in Figure 10.34 is obtained. Since the frequency of the transmitted signal is compressed as the satellite approaches, according to what proportion of the velocity vector is seen by the user receiver, the curve begins at time T_1 at a frequency several cycles higher than the transmitted frequency.

Tracking stations record Doppler observations and memory readout received during each satellite pass to relay them to a computer center. Updated orbital position and time data communications are relayed to an "injection" station from the computer center for transmission to satellite in a burst once each 12 h. Enough data is supplied in this 15 s injection message to last for 16 h of consecutive 2 min broadcasts describing the current orbital positions of the satellite.

The system accuracy depends on the accuracy of the satellite orbit computation, the effect of ionosphere refraction, the precision of the receiver speed, and heading determination. Under optimal conditions, the system is capable of producing fixes with a maximum rms error of about 35 m for the stationary receivers anywhere on Earth. Nevertheless, if a site is occupied for several weeks, an accuracy better than 1 m can be achieved. The time signal transmitted as a "beep" at the end of each 2 min transmission cycle coincides with even minutes of Coordinated Universal Time, which can be used as a chronometer check.

There are other satellite systems, either already in existence or in the planning stages, suitable for navigation. Some of these are: Marec satellites operating at VHF and owned by the intergovernment consortium INMARSAT; privately owned Geostar provides services for oil industry; and many other systems offering transcontinental communication and navigation services as well as position sensing; examples include: SATCOM, ARINC's, Avsat, Starfix, etc.

Transponders

Transponders are transducers that respond to incoming signals by generating appropriate reply messages. Recent developments in technology have made the configuration of transponders possible using elaborate and powerful on-board signal processing. This enhanced the capacity by improving the link budgets, by adjusting the antenna patterns and by making the satellite resources available on a demand basis — called the “switch board in the sky concept.”

Increased interest in deep sea exploration has brought acoustic transponders to the forefront as an important navigation tool. They provide three-dimensional position information for subservience vehicles and devices.

Some transponders are based on radar signals that respond to radar illumination. Transponders are programmed to identify friend or foe or, in some cases, simply inform ground stations about the position of aircraft.

Transponders are used for emergency warning. The U.S. and Canadian satellites carry Sarsat transponders, and Russian satellites carry Cospas transponders. They are also used as warning devices in collision avoidance systems in aircraft and land vehicles.

Global Satellite Navigation Systems

The GPS System

The Global Satellite Navigation Systems are second-generation satellites evolved primarily from the Naval Global Positioning System. They provide a continuous three-dimensional position-finding capability (i.e., latitude, longitude, and altitude), in contrast to the periodic two-dimensional information of the Transit system. Twenty-four operational satellites, as shown in Figure 10.35, constitute the system. Each satellite orbit is circular, about 2200 km high, and inclined at angles of 55° with respect to Earth's axis.

The position determination using the GPS system is based on the ability of the receivers to accurately determine the distance to the GPS satellites above the user's horizon at the time of fix. If accurate distances of two such satellites and the heights are known, then the position can be determined. In order to do this, the receiver would need to know the exact time at which the signal was broadcast and the exact time that it was received. If the propagation speeds through the atmosphere are known, the resulting range can be calculated. The measured ranges are called *pseudoranges*. Nowadays, normally, information is received from at least four satellites, leading to accurate calculations of the fix. The time errors plus propagation speed errors result in range errors, common to all GPS receivers. Time is the fourth parameter evaluated by the receiver if at least four satellites can be received at a given time. If a fifth satellite is received, an error matrix can be evaluated additionally.

Each GPS satellite broadcasts simultaneously on two frequencies for the determination and elimination of ionosphere and other atmospheric effects. The Navstar frequencies are at 1575.42 MHz and 1227.6 MHz, designated as L1 and L2 in the L-band of the UHF range. Both signals are modulated by 30 s navigation messages transmitted at 50 bits s^{-1} . The first 18 s of each 30 s frame contain *ephemeris* data for that particular satellite, which defines the position of the satellite as a function of time. The remaining 12 s is the *almanac* data, which define orbits and operational status of all satellites in the system. The GPS receivers store and use the ephemeris data to determine the pseudorange, and the almanac data to help determine the four best satellites to use for positional data at any given time. However, the “best four” philosophy has been overtaken slowly by an all-in-view philosophy.

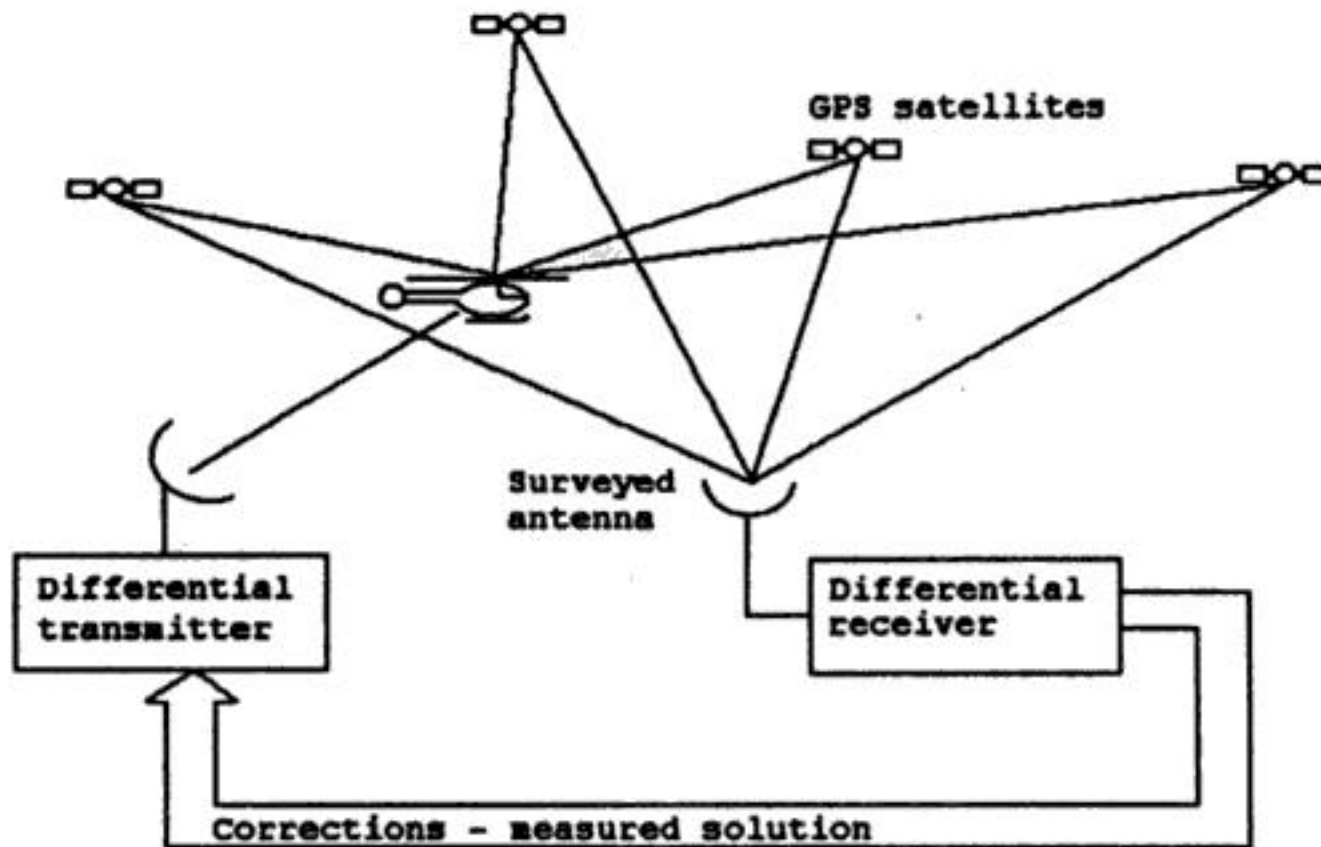


FIGURE 10.36 Differential GPS operation. Satellite and target positions are sensed by ground-fixed receivers or mobile receivers with exact known positions. Errors of the target position due to signals received from the satellites are corrected using the information from the fixed receivers. Using this method, accuracy in submeter range can be obtained.

with occasional interrupts to collect ephemeris and almanac data. Once the data is received, computation is carried out within 5 s, making this system suitable for stationary or near-stationary fixes.

Fast sequencing receivers have two channels: one for making continuous pseudorange measurements, and the other collection of the ephemeris and almanac data. This type is used in medium dynamic applications such as ground vehicles.

Continuous tracking receivers employ multiple channels (at least five) to track, compute, and process the pseudoranges to the various satellites being utilized simultaneously, thus obtaining the highest possible degree of accuracy and making it suitable for high dynamic applications such as aircraft and missiles. The parallel channel receivers are so cost effective nowadays that other types are likely to disappear.

A number of companies produce highly sophisticated GPS receivers. EURONAV GPS receivers operate on two and five channels for military applications. They provide features such as precise time, interfacing with digital flight instruments, RS-422 interface, altimeter input, self initialization, etc.

Software implementation satellite management functions, having different features, are offered by many manufacturers. In the case of DoD NAVSTAR GPS receivers, for example, three functional requirements are implemented: (1) database management of satellite almanac, ephemeris, and deterministic correction data; (2) computation of precise satellite position and velocity for use by navigation software; and (3) using satellite and receiver position data to periodically calculate the constellation of four satellites with optimum geometry for navigation. The DoD receivers are divided for three functions as Satellite Manager (SM), Satellite-Data-Base-Manager (SDBM) SV-Position Velocity Acceleration (SVPVA), and Select-Satellites (SS).

Differential navigation is also applied where one user set is navigating relative to another user set via a data link. In some cases, one user has been at a destination at some prior time and is navigating relative to coordinates measured at that point. The true values of this receiver's navigation fix are compared against the measured values, and the differences become the differential corrections. These corrections are transmitted to area user sets in real-time, or they may be recorded for post-mission use so that position fixes are free of GPS-related biases.

Differential navigation and GPS systems find applications in enroute navigations for commercial and civil aviation, military application, navigation of ships especially in shallow waters, in station keeping of aircraft, seismic geophysical explorations, land surveying, transport vehicles and traffic controls, etc.

10.5 Occupancy Detection

Jacob Fraden

Occupancy sensors detect the presence of people in a monitored area. Motion detectors respond only to moving objects. A distinction between the two is that the occupancy sensors produce signals whenever an object is stationary or not, while the motion detectors are selectively sensitive to moving objects. The applications of these sensors include security, surveillance, energy management (electric lights control), personal safety, friendly home appliances, interactive toys, novelty products, etc. Depending on the applications, the presence of humans may be detected through any means that is associated with some kind of a human body's property or actions [1]. For example, a detector may be sensitive to body weight, heat, sounds, dielectric constant, etc. The following types of detectors are presently used for the occupancy and motion sensing of people:

1. *Air pressure sensors*: detect changes in air pressure resulting from opening doors and windows
2. *Capacitive*: detectors of human body capacitance
3. *Acoustic*: detectors of sound produced by people
4. *Photoelectric*: interruption of light beams by moving objects
5. *Optoelectric*: detection of variations in illumination or optical contrast in the protected area
6. *Pressure mat switches*: pressure-sensitive long strips used on floors beneath the carpets to detect the weight of an intruder
7. *Stress detectors*: strain gages embedded into floor beams, staircases, and other structural components
8. *Switch sensors*: electrical contacts connected to doors and windows
9. *Magnetic switches*: a noncontact version of switch sensors
10. *Vibration detectors*: react to the vibration of walls or other building structures; may also be attached to doors or windows to detect movements
11. *Glass breakage detectors*: sensors reacting to specific vibrations produced by shattered glass
12. *Infrared motion detectors*: devices sensitive to heat waves emanating from warm or cold moving objects
13. *Microwave detectors*: active sensors responsive to microwave electromagnetic signals reflected from objects
14. *Ultrasonic detectors*: similar to microwaves, except that instead of electromagnetic radiation, ultrasonic waves are used
15. *Video motion detectors*: video equipment that compares a stationary image stored in memory with the current image from the protected area
16. *Laser system detectors*: similar to photoelectric detectors, except that they use narrow light beams and combinations of reflectors
17. *Triboelectric detectors*: sensors capable of detecting static electric charges carried by moving objects

One of the major aggravations in detecting occupancy or intrusion is a false positive detection. The term "false positive" means that the system indicates an intrusion when there is none. In some noncritical applications where false positive detections occur once in a while, for example, in a toy or a motion switch controlling electric lights in a room, this may be not a serious problem: the lights will be erroneously turned on for a short time, which will unlikely do any harm. In other systems, especially those used for security purposes, the false positive detections, while generally not as dangerous as false negative ones (missing an intrusion), may become a serious problem. While selecting a sensor for critical applications, consideration should be given to its reliability, selectivity, and noise immunity. It is often good practice to form a multiple sensor arrangement with symmetrical interface circuits; this can dramatically improve the reliability of a system, especially in the presence of external transmitted noise. Another efficient way to reduce erroneous detections is to use sensors operating on different physical principles [2]; for example, combining capacitive and infrared detectors is an efficient combination as they are receptive to different kinds of transmitted noise.

Ultrasonic Sensors

Ultrasonic detectors are based on transmission to the object and receiving reflected acoustic waves. Ultrasonic waves are mechanical — they cover frequency range well beyond the capabilities of human ears, i.e., over 20 kHz. However, these frequencies may be quite perceptible by smaller animals, like dogs, cats, rodents, and insects. Indeed, the ultrasonic detectors are the biological ranging devices for bats and dolphins.

When the waves are incident on an object, part of their energy is reflected. In many practical cases, the ultrasonic energy is reflected in a diffuse manner. That is, regardless of the direction where the energy comes from, it is reflected almost uniformly within a wide solid angle, which may approach 180°. If an object moves, the frequency of the reflected waves will differ from the transmitted waves. This is called the Doppler effect (see below). To generate any mechanical waves, including ultrasonic, the movement of a surface is required. This movement creates compression and expansion of the medium, which can be a gas (air), a liquid, or a solid. The most common type of the excitation device that can generate surface movement in the ultrasonic range is a piezoelectric transducer operating in the so-called *motor* mode [3]. The name implies that the piezoelectric device directly converts electrical energy into mechanical energy.

Microwave Motion Detectors

Microwave detectors offer an attractive alternative to other detectors, when it is required to cover large areas and to operate over an extended temperature range under the influence of strong interferences (e.g., wind, acoustic noise, fog, dust, moisture, etc.). The operating principle of the microwave detector is based on radiation of electromagnetic radio frequency (RF) waves toward a protected area. The most common frequencies are 10.525 GHz (X-band) and 24.125 GHz (K-band). These wavelengths are long enough ($\lambda = 3$ cm at X-band) to pass freely through most contaminants, such as airborne dust, and short enough to be reflected by larger objects.

The microwave part of the detector consists of a Gunn oscillator, an antenna, and a mixer diode. The Gunn oscillator is a diode mounted in a small precision cavity that, on application of power, oscillates at microwave frequencies. The oscillator produces electromagnetic waves, part of which is directed through an iris into a waveguide and focusing antenna that directs the radiation toward the object. Focusing characteristics of the antenna are determined by the application. As a general rule, the narrower the directional diagram of the antenna, the more sensitive it is (the antenna has a higher gain). Another general rule is that a narrow beam antenna is much larger, while a wide-angle antenna can be quite small. A typical radiated power of the transmitter is 10 mW to 20 mW.

An antenna transmits the frequency f_0 , which is defined by the wavelength λ_0 as:

$$f_0 = \frac{c_0}{\lambda_0} \quad (10.106)$$

where c_0 is the speed of light. When the target moves toward or away from the transmitting antenna, the frequency of the reflected radiation will change. Thus, if the target is moving away with velocity v , the reflected frequency will decrease, and it will increase for the approaching targets. This is called the *Doppler effect*, after the Austrian scientist Christian Johann Doppler (1803–1853). While the effect was first discovered for sound, it is applicable to electromagnetic radiation as well. However, in contrast to sound waves that may propagate with velocities dependent on movement of the source of the sound, electromagnetic waves propagate with speed of light, which is an absolute constant. The frequency of reflected electromagnetic waves can be predicted by the theory of relativity as:

$$f_r = f_0 \frac{\sqrt{1 - (v/c_0)^2}}{1 + v/c_0} \quad (10.107)$$

For practical purposes, however, the quantity $(v/c_0)^2$ is very small compared with unity; hence, it can be ignored. Then, the equation for the frequency of the reflected waves becomes identical to that for the acoustic waves:

$$f_r = f_0 \frac{1}{1+v/c_0} \quad (10.108)$$

Due to a Doppler effect, the reflected waves have a different frequency f_r . A mixing diode combines the radiated (reference) and reflected frequencies and, being a nonlinear device, produces a signal that contains multiple harmonics of both frequencies.

The Doppler frequency in the mixer can be found from:

$$\Delta f = f_0 - f_r = f_0 \frac{1}{c_0/v+1} \quad (10.109)$$

and since $c_0/v \gg 1$, the following holds after substituting Equation 10.106:

$$\Delta f \approx \frac{v}{\lambda_0} \quad (10.110)$$

Therefore, the signal frequency at the output of the mixer is linearly proportional to the velocity of a moving target. For example, if a person walks toward the detectors with a velocity of 0.6 m s^{-1} , a Doppler frequency for the X-band detector is $\Delta f = 0.6/0.03 = 20 \text{ Hz}$.

Equation 10.110 holds true only for movements in the normal direction. When the target moves at angles Θ with respect to the detector, the Doppler frequency is:

$$\Delta f \approx \frac{v}{\lambda_0} \cos \Theta \quad (10.111)$$

Micropower Impulse Radar

In 1993, Lawrence Livermore National Laboratory developed a *micropower impulse radar* (MIR), which is a low-cost, noncontact ranging sensor [3]. The operating principle of the MIR is fundamentally the same as a conventional pulse radar system, but with several significant differences. The MIR consists of a noise generator whose output signal triggers a pulse generator. Each pulse has a fixed short duration, while the repetition of these pulses is random, according to triggering by the noise generator. The pulses are spaced randomly with respect to one another in a Gaussian noise-like pattern. It can be said that the pulses have the pulse frequency modulation (PFM) by white noise with maximum index of 20%. In turn, the square-wave pulses cause amplitude modulation (AM) of a radio transmitter. The radio transmitter produces short bursts of high-frequency radio signal that propagate from the transmitting antenna to the surrounding space. The electromagnetic waves reflect from the objects and propagate back to the radar. The same pulse generator that modulates the transmitter, gates (with a predetermined delay) the radio receiver to enable the output of the MIR only during a specific time window. Another reason for gating the receiver is to reduce its power consumption. The reflected pulses are received, demodulated (the square-wave shape is restored from the radio signal), and the time delay with respect to the transmitted pulses is measured. Since the pulses are spaced randomly, practically any number of identical MIR systems can operate in the same space without a frequency division (i.e., they work at the same carrier frequency within the same bandwidth). There is little chance that bursts from the interfering transmitters overlap and, if they do, the interference level is significantly reduced by the averaging circuit.

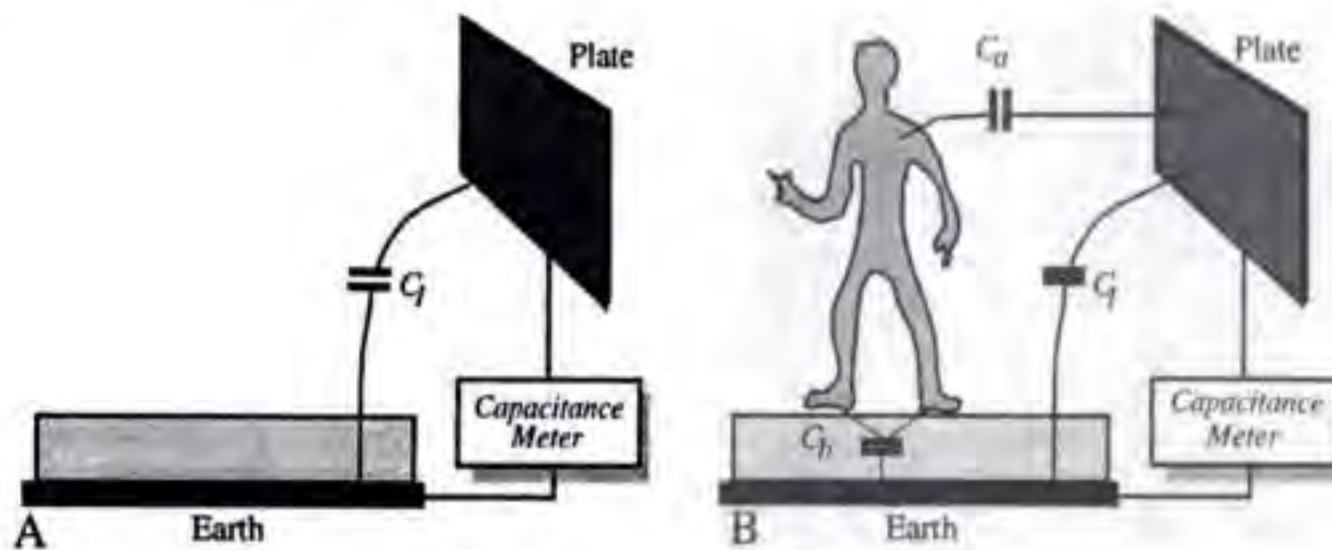


FIGURE 10.37 An intruder brings in additional capacitance to a detection circuit.

Capacitive Occupancy Detectors

Being a conductive medium with a high dielectric constant, a human body develops a coupling capacitance to its surroundings. This capacitance greatly depends on such factors as body size, clothing, materials, type of surrounding objects, weather, etc. However wide the coupling range is, the capacitance can vary from a few picofarads to several nanofarads. When a person moves, the coupling capacitance changes, thus making it possible to discriminate static objects from the moving ones. In effect, all objects form some degree of capacitive coupling with respect to one another. If a human (or for that purpose, anything) moves into the vicinity of the objects whose coupling capacitance with each other has been previously established, a new capacitive value arises between the objects as a result of the presence of an intruding body [3]. Figure 10.37 shows that the capacitance between a test plate and Earth is equal to C_1 . When a person moves into the vicinity of the plate, it forms two additional capacitors: one between the plate and its own body C_a , and the other between the body and the Earth, C_b . Then, the resulting capacitance C between the plate and the Earth becomes larger by ΔC .

$$C = C_1 + \Delta C = C_1 + \frac{C_a C_b}{C_a + C_b} \quad (10.112)$$

With the appropriate apparatus, this phenomenon can be used for occupancy detection [3]. What is required is to measure a capacitance between a test plate (the probe) and a reference plate (the Earth).

Figure 10.38 illustrates a circuit diagram for detecting variations in the probe capacitance C_p [4]. That capacitance is charged from a reference voltage source V_{ref} through a gate formed by transistor Q_1 when the output voltage of a control oscillator goes low. When it goes high, transistor Q_1 closes while Q_2 opens. The probe capacitance C_p discharges through a constant-current sink constructed with a transistor Q_3 . A capacitor C_1 filters the voltage spikes across the transistor. The average voltage, e_p , represents a value of the capacitor C_p . When an intruder approaches the probe, the latter's capacitance increases, which results in a voltage rise across C_1 . The voltage change passes through the capacitor C_2 to the input of a comparator with a fixed threshold V_T . The comparator produces the output signal V_{out} when the input voltage exceeds the threshold value.

When a capacitive occupancy (proximity) sensor is used near or on a metal device, its sensitivity can be severely reduced due to capacitive coupling between the electrode and the device's metallic parts. An effective way to reduce that stray capacitance is to use driven shields [3].

Triboelectric Detectors

Any object can accumulate, on its surface, static electricity. These naturally occurring charges arise from the triboelectric effect; that is, a process of charge separation due to object movements, friction of clothing

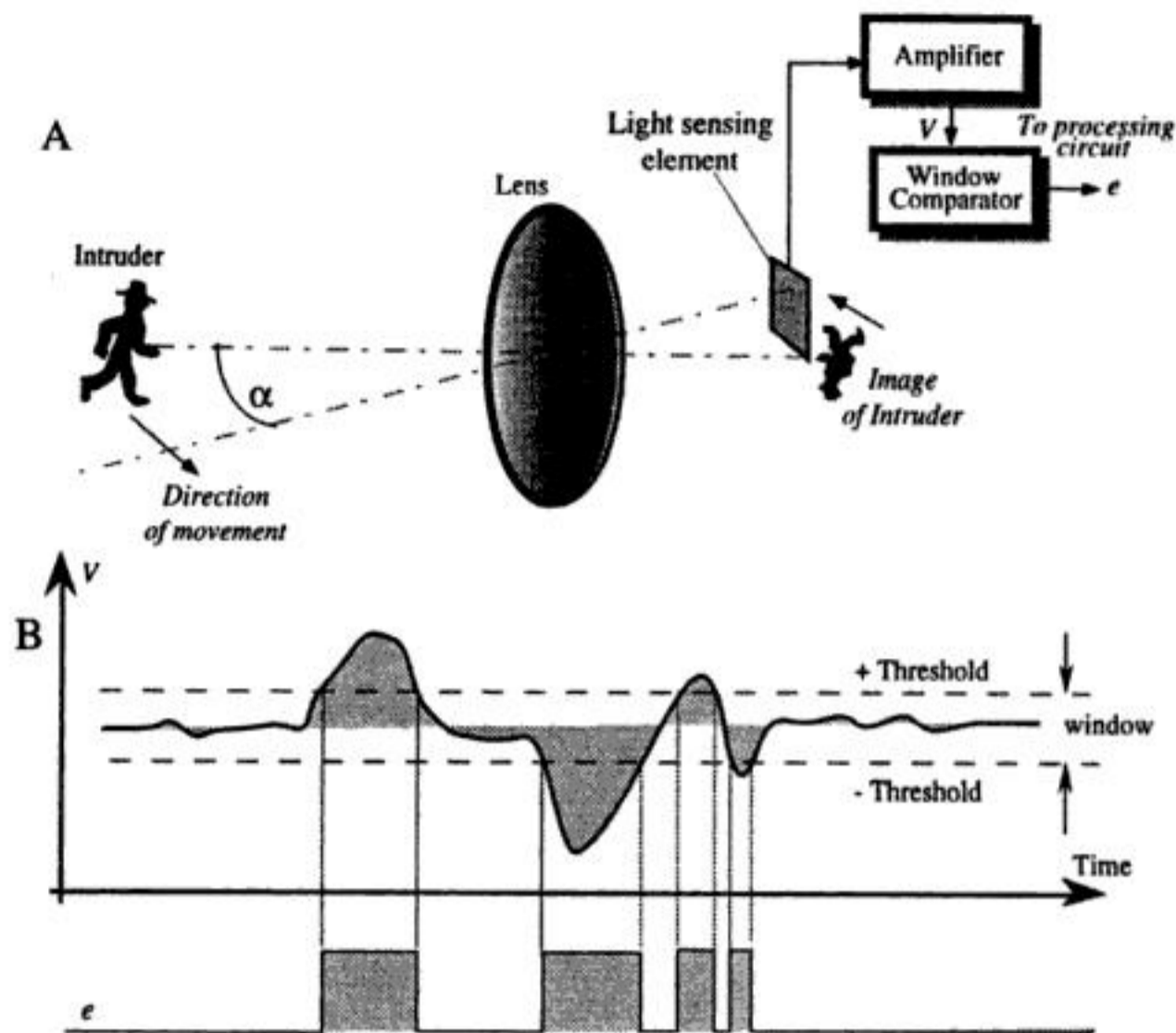


FIGURE 10.39 General arrangement of an optoelectronic motion detector. A lens forms an image of a moving object (intruder). When the image crosses the optical axis of the sensor, it superimposes with the sensitive element (A). The element responds with the signal that is amplified and compared with two thresholds in the window comparator (B). (From J. Fraden, *Handbook of Modern Sensors*, 2nd ed., Woodburg, NY: AIP Press, 1997. With permission.)

Most of the objects (apart from very hot) radiate electromagnetic waves only in the mid- and far-infrared spectral ranges. Hence, visible and near-infrared light motion detectors must rely on an additional source of light to illuminate the object. The light is reflected by the object's body toward the focusing device for subsequent detection. Such illumination can be sunlight or the invisible infrared light from an additional near-infrared light source (a projector).

The major application areas for the optoelectronic motion detectors are in security systems (to detect intruders), in energy management (to turn lights on and off), and in the so-called "smart" homes where they can control various appliances such as air conditioners, cooling fans, stereo players, etc. They can also be used in robots, toys, and novelty products. The most important advantage of an optoelectronic motion detector is simplicity and low cost.

Sensor Structures

A general structure of an optoelectronic motion detector is shown in Figure 10.39(A). Regardless what kind of sensing element is employed, the following components are essential: a focusing device (a lens or a focusing mirror), a light detecting element, and a threshold comparator. An optoelectronic motion detector resembles a photographic camera. Its focusing components create an image of its field of view on a focal plane. While there is no mechanical shutter like in a camera, in place of the film, a light sensitive element is used. The element converts the focused light into an electric signal. A focusing lens creates an image of the surroundings on a focal plane where the light sensitive element is positioned. If the area is unoccupied, the image is static and the output signal from the element is steady stable. When an "intruder" enters the room and keeps moving, his/her image on the focal plane also moves. In a certain moment, the intruder's body is displaced for an angle α and the image overlaps with the element. This

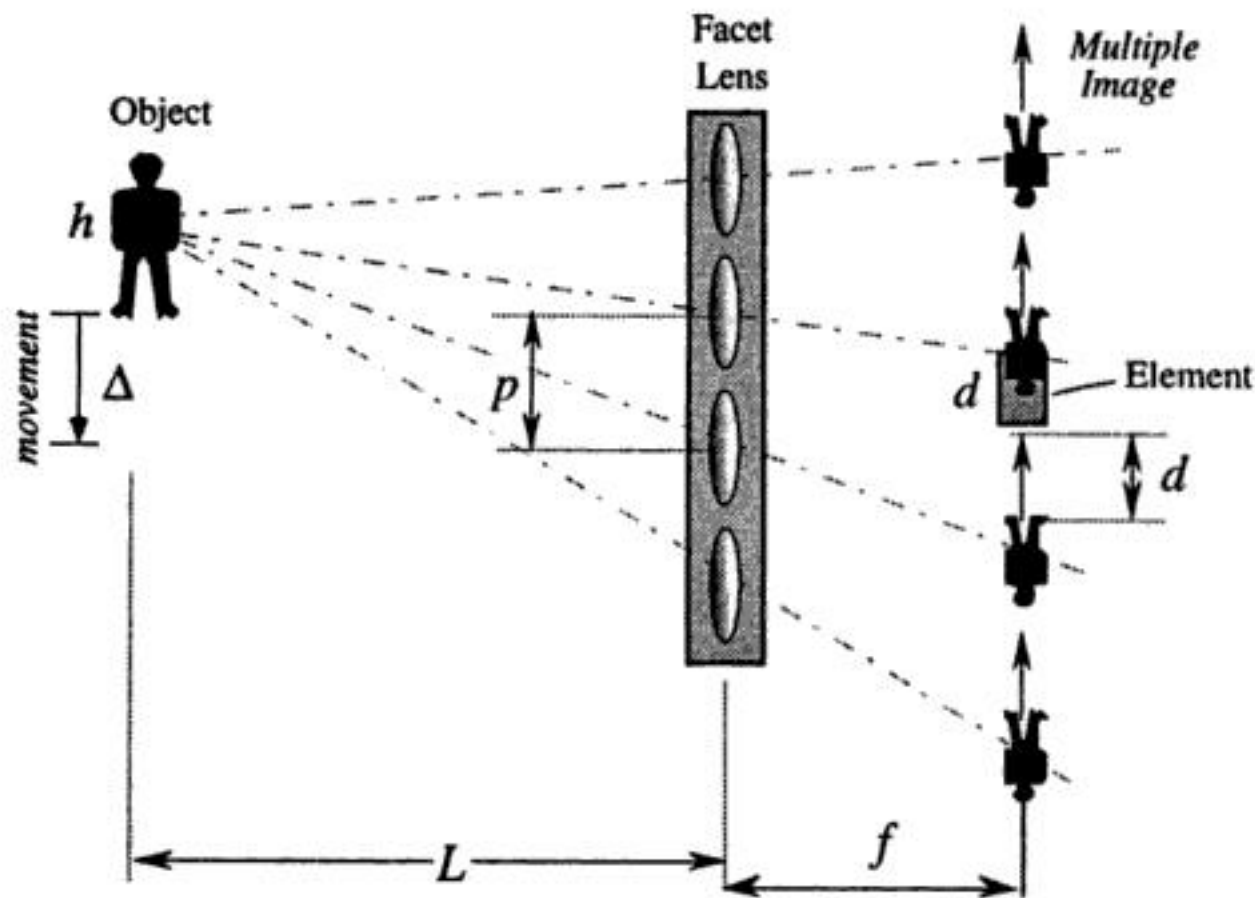


FIGURE 10.41 Facet lens creates multiple images near the sensing element.

resulting in an alternate signal. By combining multiple facets, it is possible to create any desirable detecting pattern in the field of view, in both horizontal and vertical planes. Positioning of the facet lens, focal distances, number, and a pitch of the facets (a distance between the optical axes of two adjacent facets) can be calculated in every case by applying rules of geometrical optics [3]. In the far-infrared spectral range (thermal radiation sensors), the polyethylene facet Fresnel lenses are used almost exclusively, thanks to their low cost and relatively high efficiency.

For the visible portion of the optical spectrum, a simple, very inexpensive, yet efficient motion detector can be developed for nondemanding applications, like light control or interactive toys, using simple photoresistors and pinhole lenses [3, 6, 7].

Far-Infrared Motion Detectors

A motion detector that perceives electromagnetic radiation that is naturally emitted by any object operates in the optical range of thermal radiation, also called far-infrared (FIR). Such detectors are responsive to radiative heat exchange between the sensing element and the moving object. The principle of thermal motion detection is based on the physical theory of emission of electromagnetic radiation from any object whose temperature is above absolute zero (see Chapter 32, Section 6, on *Infrared Thermometers*).

For IR motion detection, it is essential that a surface temperature of an object be different from that of the surrounding objects, so a thermal contrast would exist. All objects emanate thermal radiation from their surfaces and the intensity of that radiation is governed by the Stefan–Boltzmann law. If the object is warmer than the surroundings, its thermal radiation is shifted toward shorter wavelengths and its intensity becomes stronger. Many objects whose movement is to be detected are nonmetals, hence they radiate thermal energy quite uniformly within a hemisphere. Moreover, the dielectric objects generally have a high emissivity. Human skin is one of the best emitters, with emissivity over 90%, while most fabrics also have high emissivities, between 0.74 and 0.95 [3]. Below, two types of far-infrared motion detectors are described. The first utilizes a passive infrared (PIR) sensor, while the second has active far-infrared (AFIR) elements.

PIR Motion Detectors

These detectors became very popular for security and energy management systems. The PIR sensing element must be responsive to far-infrared radiation within a spectral range from $4\ \mu\text{m}$ to $20\ \mu\text{m}$ where most of the thermal power emanated by humans is concentrated. There are three types of sensing elements

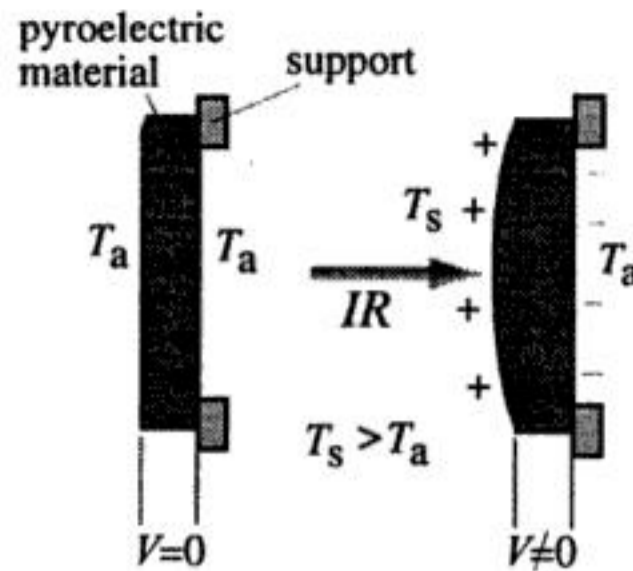


FIGURE 10.42 A simplified model of a pyroelectric effect as a secondary effect of piezoelectricity. Initially, the element has a uniform temperature (A); upon exposure to thermal radiation, its front side expands, causing a stress-induced charge (B).

that are potentially useful for that detector: thermistors, thermopiles, and pyroelectrics; however, pyroelectric elements are used almost exclusively for the motion detection thanks to their simplicity, low cost, high responsivity, and broad dynamic range. A pyroelectric effect is described in Chapter 32, Section 7 on *Pyroelectric Thermometers*. How this effect can be used in practical sensor design is discussed here.

A pyroelectric material generates an electric charge in response to thermal energy flow through its body. In a simplified way it may be described as a secondary effect of thermal expansion (Figure 10.42). Since all pyroelectrics are also piezoelectrics, the absorbed heat causes the front side of the sensing element to expand. The resulting thermally induced stress leads to the development of a piezoelectric charge on the element electrodes. This charge is manifested as voltage across the electrodes deposited on the opposite sides of the material. Unfortunately, the piezoelectric properties of the element also have a negative effect. If the sensor is subjected to a minute mechanical stress due to any external force, it also generates a charge that in most cases is indistinguishable from that caused by the infrared heat waves. Sources of such mechanical noise are wind, building vibrations, loud sound, etc.

To separate thermally induced charges from the piezoelectrically induced charges, a pyroelectric sensor is usually fabricated in symmetrical form (Figure 10.43(A)). Two identical elements are positioned inside the sensor's housing. The elements are connected to the electronic circuit in such a manner as to produce the out-of-phase signals when subjected to the same in-phase inputs. The idea is that interferences produced by, for example, the piezoelectric effect or spurious heat signals are applied to both electrodes simultaneously (in phase) and thus will be canceled at the input of the circuit, while the variable thermal radiation to be detected will be absorbed by only one element at a time, thus avoiding a cancellation.

One way to fabricate a differential sensor is to deposit two pairs of electrodes on both sides of a pyroelectric element. Each pair forms a capacitor that can be charged either by heat or by mechanical stress. The electrodes on the upper side of the sensor are connected together forming one continuous electrode, while the two bottom electrodes are separated, thus creating the opposite-serially connected capacitors. Depending on the side where the electrodes are positioned, the output signal will have either a positive or negative polarity for the thermal influx. In some applications, a more complex pattern of the sensing electrodes is required (for example, to form predetermined detection zones), so that more than one pair of electrodes is needed. In such a case, for better rejection of the in-phase signals (common mode rejection), the sensor should still have an even number of pairs where positions of the pairs alternate for better geometrical symmetry. Sometimes, such an alternating connection is called an interdigitized electrode.

A differential sensing element should be mounted in such a way as to ensure that both parts of the sensor generate the same signal if subjected to the same external factors. At any moment, the optical component must focus a thermal image of an object on the surface of one part of the sensor only, which is occupied by a single pair of electrodes. The element generates a charge only across the electrode pair that is subjected to a heat flux. When the thermal image moves from one electrode to another, the current

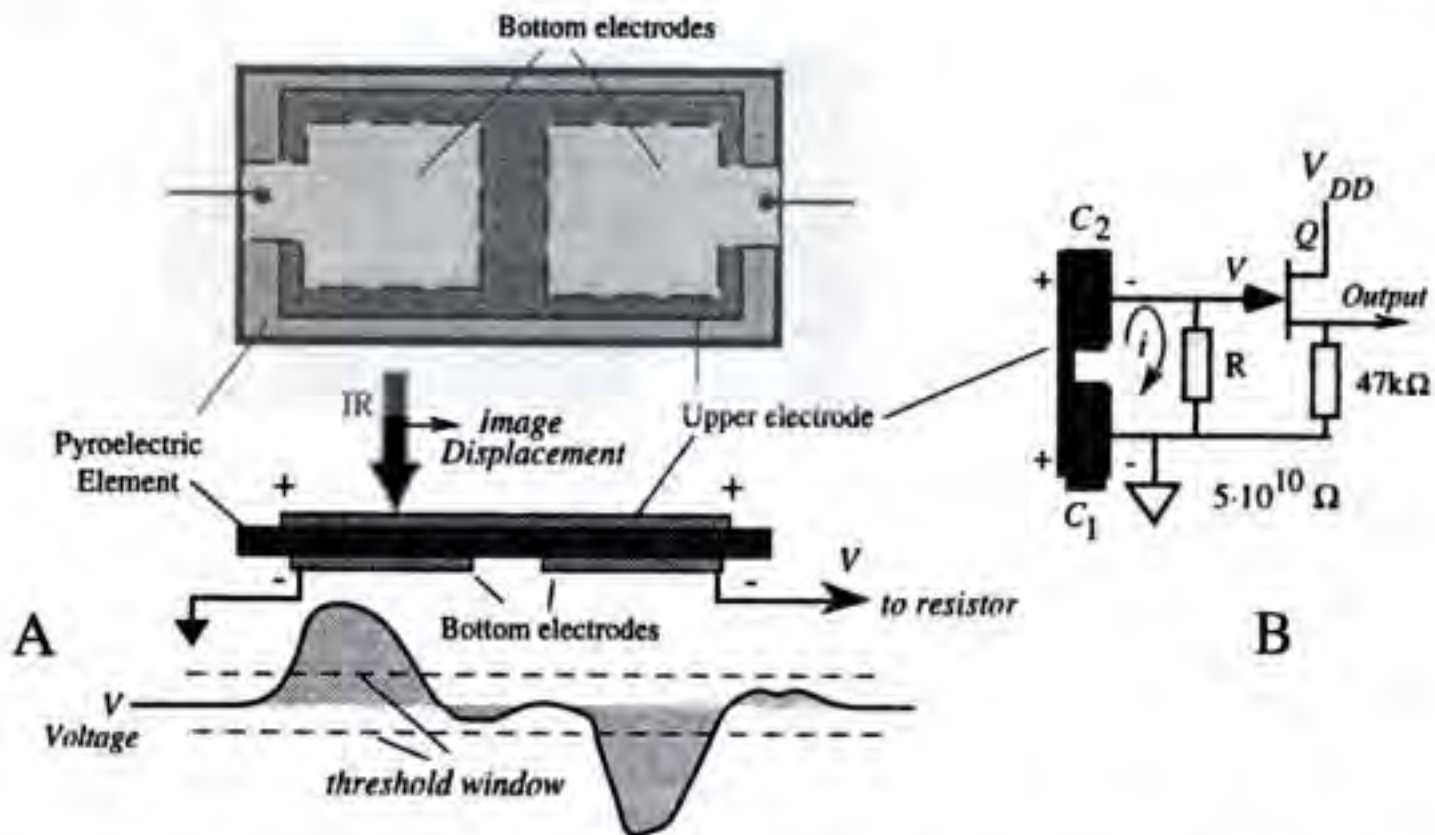


FIGURE 10.43 Dual pyroelectric sensor. (A) A sensing element with a front (upper) electrode and two bottom electrodes deposited on a common crystalline substrate. (B) A moving thermal image travels from the left part of the sensor to the right, generating an alternating voltage across bias resistor, R .

i flowing from the sensing element to the bias resistor R (Figure 10.43(B)) changes from zero, to positive, then to zero, to negative, and again to zero (Figure 10.43(A) lower portion). A JFET transistor Q is used as an impedance converter. The resistor R value must be very high. For example, a typical alternate current generated by the element in response to a moving person is on the order of 1 pA (10^{-12} A). If a desirable output voltage for a specific distance is $v = 50 \text{ mV}$, then according to Ohm's law, the resistor value is $R = v/i = 50 \text{ G}\Omega$ ($5 \times 10^{10} \Omega$). Such a resistor cannot be directly connected to a regular electronic circuit; hence, transistor Q serves as a voltage follower (the gain is close to unity). Its typical output impedance is on the order of several kilohms.

The output current i from the PIR sensor can be calculated on the basis of the Stefan-Boltzmann law as [3]:

$$i = \frac{2P\sigma a\gamma}{\pi hc} bT_s^3 \frac{\Delta T}{L^2} \quad (10.113)$$

where $\Delta T = (T_b - T_a)$ is the temperature gradient between the object and its surroundings, P is the pyroelectric coefficient, σ is the Stefan-Boltzmann constant, a is the lens area, γ is the lens transmission coefficient, h is the thickness, and c is the specific heat of the pyroelectric element, respectively, and L is the distance to the object.

There are several conclusions that can be drawn from Equation 10.113. The first part of the equation (the first ratio) characterizes a detector, while the rest relates to an object. The pyroelectric current i is directly proportional to the temperature difference (thermal contrast) between the object and its surroundings. It is also proportional to the surface area of the object that faces the detector. A contribution of the ambient temperature T_a is not as strong as it might appear from its third power. The ambient temperature must be entered in kelvin, hence its variations become relatively small with respect to the scale. The thinner the sensing element, the more sensitive the detector. The lens area also directly affects signal magnitude. On the other hand, pyroelectric current does not depend on the sensor's area as long as the lens focuses an entire image on a sensing element.

AFIR Motion Detectors

The AFIR motion detector is a new class of thermal sensors whose operating principle is based on balancing thermal power supplied to the sensing element [8, 9]. Contrary to a passive motion detector

that absorbs thermal radiation from a warmer object, an AFIR motion detector is active; that is, it radiates heat waves *toward* the surroundings. The sensor's surface temperature (T_s) is maintained somewhat above ambient. The element is combined with a focusing system, very much like the PIR detector; however, the function of that system is inverse to that of the passive detectors. A focusing part in the AFIR detector projects a thermal image of the warm sensing element into its surroundings. The AFIR sensors have a significant advantage over the PIR: immunity against many interferences (such as RFI and microphonics).

The output voltage from the AFIR motion detector can be described by the following equation [3]:

$$\Delta V = -\frac{R}{V_0} \frac{\sigma \alpha \gamma}{\pi} b T_s^3 \frac{\Delta T}{L^2} \quad (10.114)$$

where R is the resistance of the sensor's heater and V_0 is the heating voltage. The minus sign indicates that for warmer moving objects, the output voltage decreases. There is an obvious similarity between Equations 10.113 and 10.114; however, sensitivity (detection range) of the AFIR sensor can be easily controlled by varying R or V_0 . For better sensitivity, the temperature increment above ambient can be maintained on a fairly low level. Practically, the element is heated above ambient by only about 0.2°C .

References

1. S. Blumenkrantz, *Personal and Organizational Security Handbook*, Government Data Publications, Washington, D.C.: 1989.
2. P. Ryser and G. Pfister, Optical fire and security technology: sensor principles and detection intelligence, *Transducers'91. Int. Conf. Solid-State Sensors Actuators*, 1991, 579-583.
3. J. Fraden, *Handbook of Modern Sensors*, 2nd ed., Woodburg, NY: AIP Press, 1997.
4. N. M. Calvin, *Capacitance proximity sensor*. U.S. Patent No. 4,345,167, 1982.
5. J. Fraden, *Apparatus and method for detecting movement of an object*, U.S. Patent No. 5,019,804, 1991.
6. J. Fraden, *Motion discontinuance detection system and method*. U.S. Patent No. 4,450,351, 1984.
7. J. Fraden, *Toy including motion-detecting means for activating same*. U.S. Patent No. 4,479,329, 1984.
8. J. Fraden, *Active infrared motion detector and method for detecting movement*. U.S. Patent No. 4,896,039, 1990.
9. J. Fraden, Active far infrared detectors, in *Temperature. Its Measurement and Control in Science and Industry*, Vol. 6, Woodburg, NY: American Institute of Physics, 1992, Part 2, 831-836.

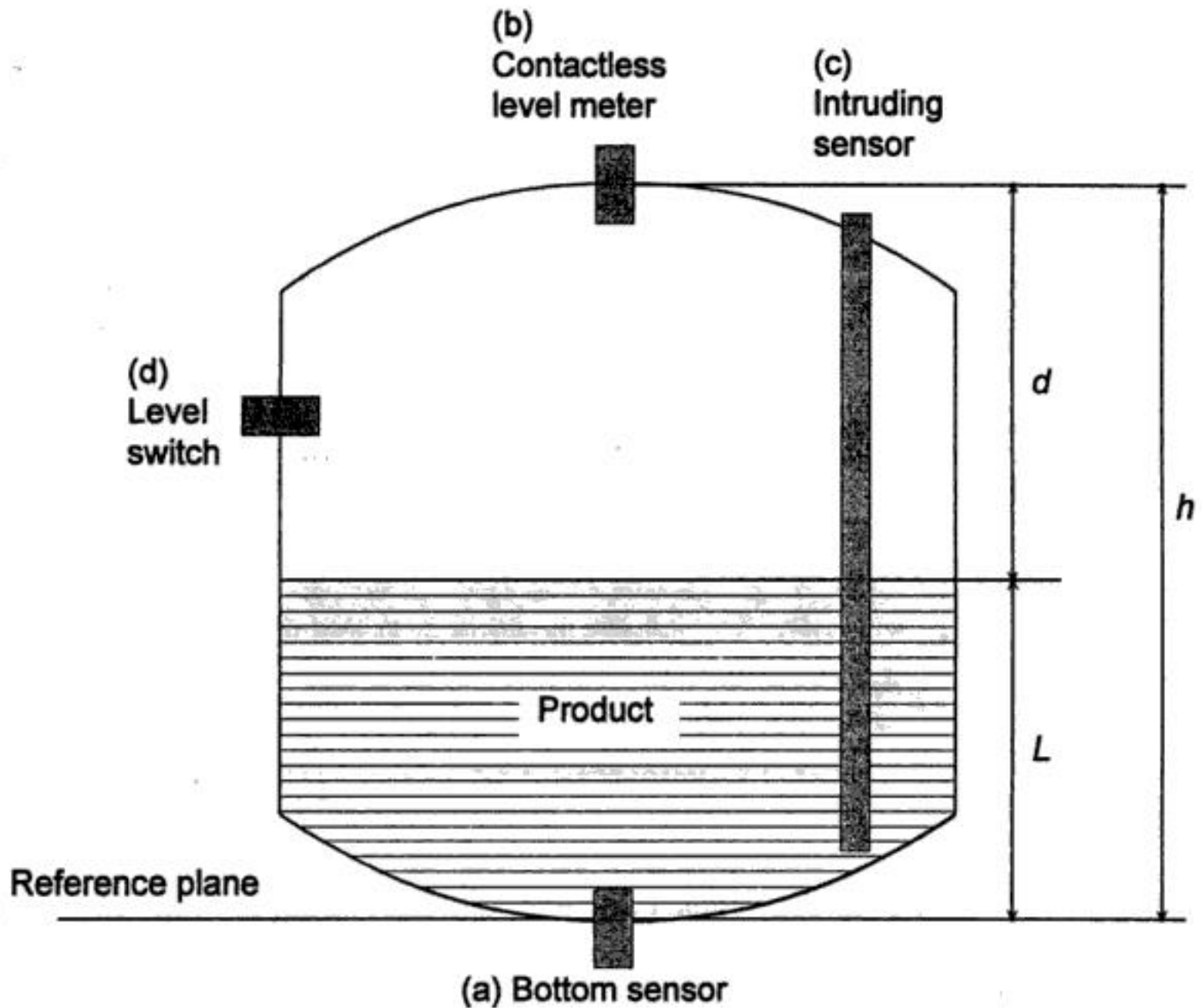


FIGURE 11.1 Representation of a tank with a liquid or solid material (hatched area), the product to be measured. The level sensor can be mounted (a) contacting product at the bottom, (b) as a contactless instrument on top, (c) as an intrusive sensor, or (d) at the sides as a level switch.

11.1 Measurements Using the Effects of Density

All methods described in this chapter have in common that the product in the tank has an effect due to its density ρ , (1) producing buoyancy to a solid submerged into the liquid, or (2) executing a force due to its weight.

Displacer

Displacers measure the buoyancy of a solid body that is partially submerged in the liquid. The change in weight is measured. Figure 11.2 illustrates the parameters used for these calculations. The cross section A of the body is assumed to be constant over its length b . The weight of force F_G due to gravity g and mass m is:

$$F_G = g m = g A b \rho_D \quad (11.2)$$

The buoyant force F_B accounts for the partial length L_d that is submerged with the remainder of the body in the atmosphere:

$$F_B = g A L_d \rho_L + g A (b - L_d) \rho_A \quad (11.3)$$

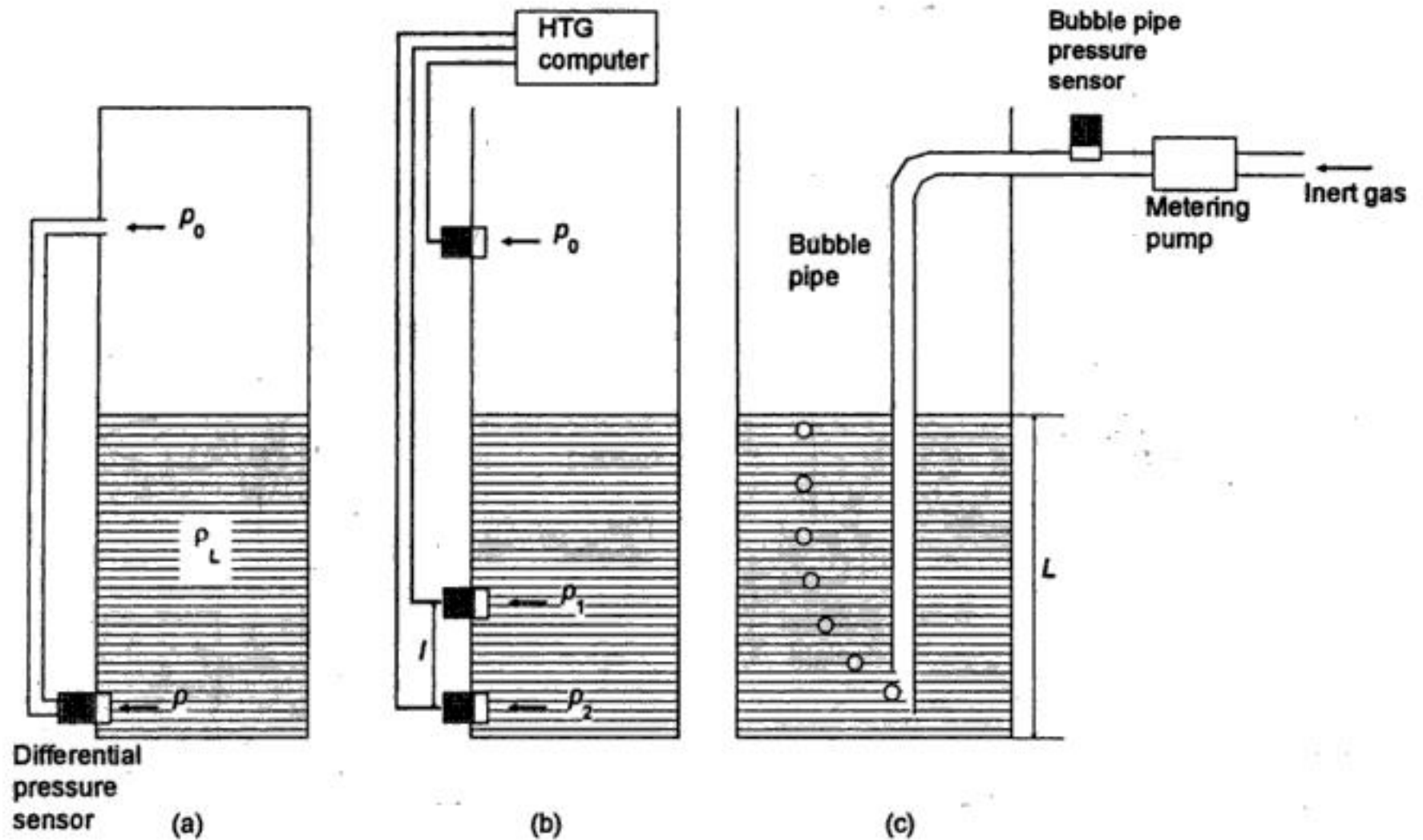


FIGURE 11.5 Level gaging by hydrostatic pressure measurement. The bottom pressure p is proportional to level. (a) The atmospheric pressure p_0 can be taken into consideration by a differential measurement. The low side of the differential pressure sensor is connected via a thin pipe to the top of the tank. (b) A differential measurement within the liquid is called "hydrostatic tank gaging, HTG" and can be used to compensate errors due to density variations of the liquid. The signals from all three sensors are evaluated by a computer. (c) With a so-called "bubble tube," the sensor can be mounted on the top of the tank: an inert gas is injected into the tube such that bubbles of gas escape from the end of the tube. The flow rate of the gas is constant so the head pressure in the system can be measured at the inlet of the pipe.

$$\rho_L = \frac{p_2 - p_1}{gl} \Rightarrow L = \frac{p_2 - p_0}{p_2 - p_1} l \quad (11.9)$$

A system of this configuration is often called "hydrostatic tank gaging" (HTG). Figure 11.5(c) shows a further arrangement, called "bubble tube," in which the bottom pressure is transmitted to the top of the tank. This is often used for level gaging if the sensor cannot be mounted at the bottom of the tank. It requires a tank with pressure equalization due to the steady insertion of inert gas.

Balance Method

Here simply the weight F of the complete tank is measured, dependent on the level L :

$$F = F_0 + gAL\rho_L \quad (11.10)$$

where F_0 is the weight of the empty tank and A the cross-sectional area, which is assumed to be constant throughout the tank height. In order to measure the weight force correctly, it is necessary to isolate the complete tank mechanically. For precise measurements, the buoyancy in air must be taken into consideration:

$$F = F_0 + gAL(\rho_L - \rho_A) \Leftrightarrow L = \frac{F - F_0}{gA(\rho_L - \rho_A)} \quad (11.11)$$

For techniques of weight measurement, refer to Chapter 20 of this handbook.

This method has severe disadvantages when the tank is not inside a building. Outside, wind forces and the weight of snow and rain can cause errors.

11.2 Time-of-Flight Measurements

An indirect measurement of level is evaluating the time-of-flight of a wave propagating through the atmosphere above the liquid or solid. This is primarily a distance measurement; the level can then be calculated accordingly. The increasing demand of industry for nonintrusive continuous level gaging systems has been instrumental in accelerating the development of technologies using time-of-flight measurements [7].

Basic Principle

Although different types of physical waves (acoustic or electromagnetic) are applied, the principle of all these methods is the same: a modulated signal is emitted as a wave toward the product, reflected at its surface and received by a sensor, which in many cases is the same, (e.g., the ultrasonic piezoelectric transducer or the radar antenna). Figure 11.6 demonstrates the principle of operation. The measuring system evaluates the time-of-flight t of the signal:

$$t = \frac{2d}{v} \quad (11.12)$$

where v is the propagation velocity of the waves.

One can generate an unmodulated pulse, a modulated burst as in Figure 11.6(b), or special forms. Table 11.1 lists the main properties of the three preferred types of waves, used for time-of-flight level gaging.

The very short time spans of only a few nanoseconds for radar and laser measurement techniques require the use of time expansion by sampling (see Chapter 85 of this handbook) or special evaluation methods (see below).

Ultrasonic

Ultrasonic waves are longitudinal acoustic waves with frequencies above 20 kHz. Ultrasonic waves need a propagation medium, which for level measurements is the atmosphere above the product being measured. Sound propagates with a velocity of about 340 m s^{-1} in air; but this value is highly dependent on temperature and composition of the gas, and also on its pressure (see Chapter 6, Section 7 of this handbook). In vacuum, ultrasonic waves cannot propagate. In practice, the reflection ratio is nearly 100% at the product's surface (e.g., at transitions gas/liquid or gas/solid). Piezoelectric transducers (see Chapter 26, Section 3 of this handbook) are utilized as emitter and detector for ultrasonic waves, a membrane coupling it to the atmosphere. The sensor is installed as in Figure 11.1(b), the signal form is as in Figure 11.6(b). Level gaging is, in principle, also possible with audible sound 16 Hz to 20 kHz or infrasonic waves less than 16 Hz.

Another procedure is to propagate the waves within the liquid by a sensor mounted at the bottom of the tank. The velocity of sound in the liquid must be known, considering the dependence on temperature and type of liquid. This method is similar to an echo sounder on ships for measuring the water depth. For more information about time-of-flight ultrasound evaluation techniques, refer to Chapter 6, Section 7 of this handbook.

Microwaves

Microwaves are generally understood to be electromagnetic waves with frequencies above 2 GHz and wavelengths of less than 15 cm. For technical purposes, microwave frequencies are used up to max. 120 GHz; in practice, the range around 10 GHz (X-band) is preferred.

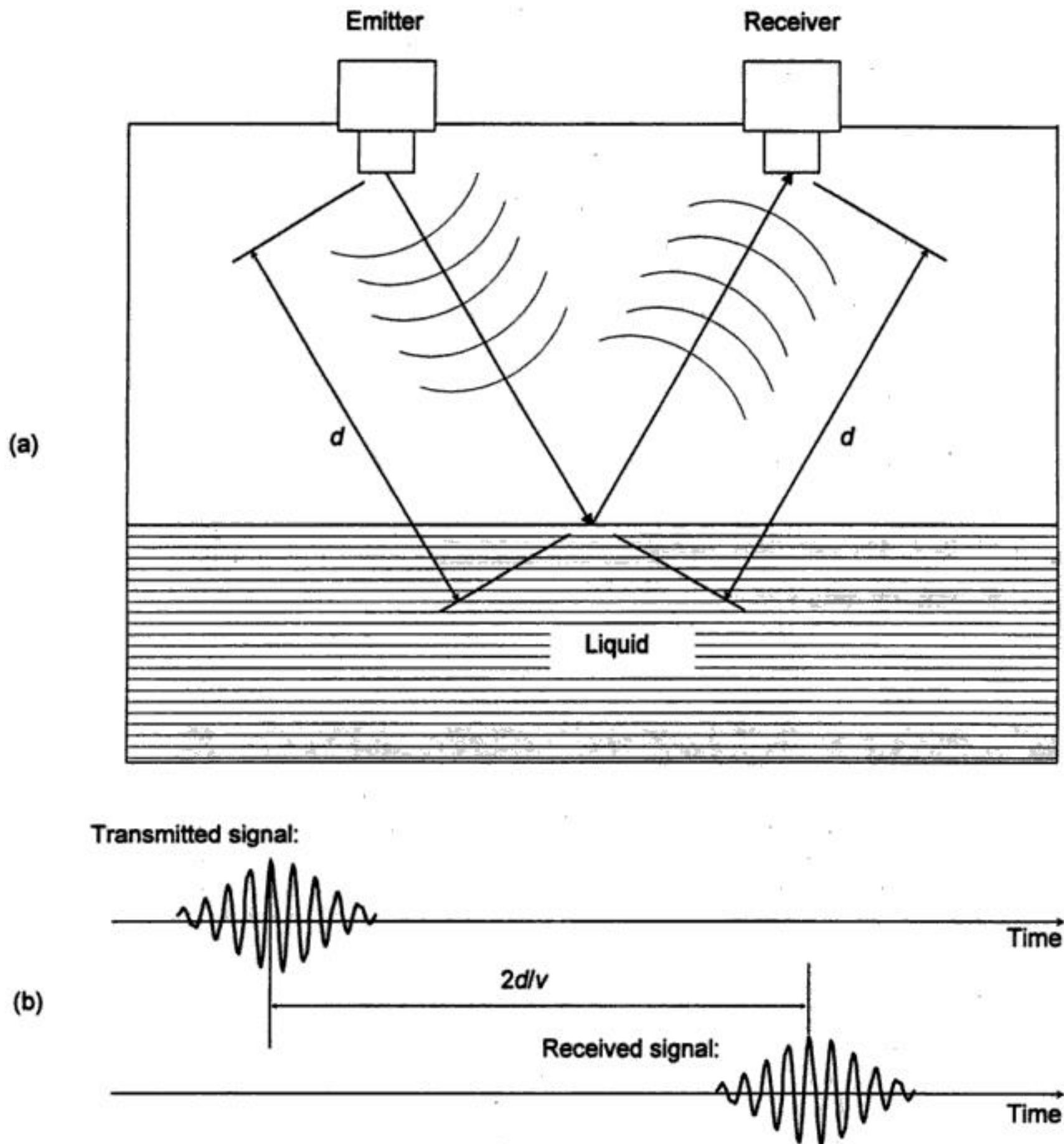


FIGURE 11.6 (a) Representation of time-of-flight measurements. The emitter couples a wave (ultrasonic or electromagnetic) into the atmosphere that propagates the wave toward the liquid. Its surface reflects the wave and a sensor receives it. (b) Due to the propagation velocity v , a time delay is measured between emission and receipt of the signal. This example is characterized by a modulated burst. The time scale is arbitrary.

TABLE 11.1 Properties of the wave types for time-of-flight measuring.

Principle	Wave Velocity	Avg. Carrier Frequency	Wavelength	Avg. Burst Time
Ultrasonic	340 m s^{-1}	50 kHz	7 mm	1 ms
Radar	$300,000 \text{ km s}^{-1}$	10 GHz	3 cm	1 ns
Laser	$300,000 \text{ km s}^{-1}$	300 THz	1 μm	1 ns

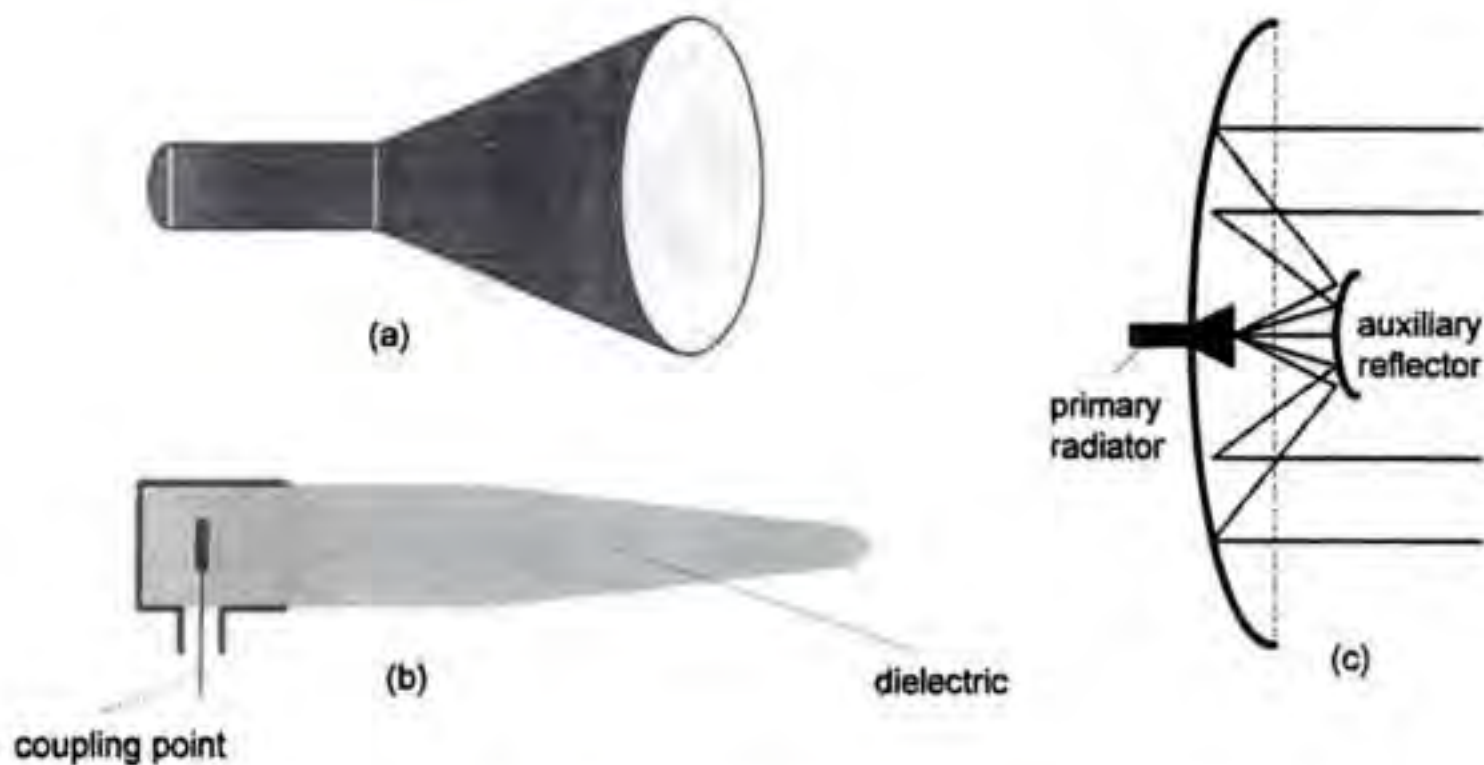


FIGURE 11.7 Practical antenna forms used for radar level instruments: (a) conical horn antenna, (b) dielectric rod antenna, and (c) parabolic mirror with a small antenna as primary radiator and an auxiliary reflector giving a very small beam angle (so-called Cassegrain model).

The usually applied time-of-flight measurements with microwaves are RADAR-based [8, 9]. The term “RADAR” is generally understood to mean a method by means of which short electromagnetic waves are used to detect distant objects and determine their location and movement. It is an acronym from RADio Detection And Ranging. Figure 11.7 shows preferred antenna forms. They are usually combined with a compact sensor, as in Figure 11.8. For level measuring systems, a small radiation angle is desirable in order to avoid interfering reflections from the tank wall or tank internals as much as possible. The larger the aperture area, the smaller the radiation angle and the higher the antenna gain. The power balance is given by the general radar equation:

$$P_R = \frac{P_T G_T R G_R}{D^2} \tag{11.13}$$

- where
- P_R = received power
 - P_T = transmitted power
 - G_T = transmitting antenna gain
 - R = reflection factor of target
 - G_R = receiving antenna gain
 - D^2 = propagation loss to and from the surface, due to power density decrease and atmospheric influences

The reflection factor R of the product’s surface is dependent on the dielectric permittivity ϵ_r of the liquid or bulk material:

$$R = \frac{(\sqrt{\epsilon_r} - 1)^2}{(\sqrt{\epsilon_r} + 1)^2} \tag{11.14}$$



FIGURE 11.8 Design of a compact industrial level radar system. The converter above the flange includes the complete microwave circuitry, signal processing stages, microprocessor control, display, power supply, and output signal [6].

In level measurement situations, the reflecting area is so large that it intersects the beam cross section completely; therefore, D^2 is approximately proportional with distance d^2 . Thus also, the received power decreases proportionately with d^2 , as derived in [8]:

$$P_R \propto \frac{1}{d^2} \quad (11.15)$$

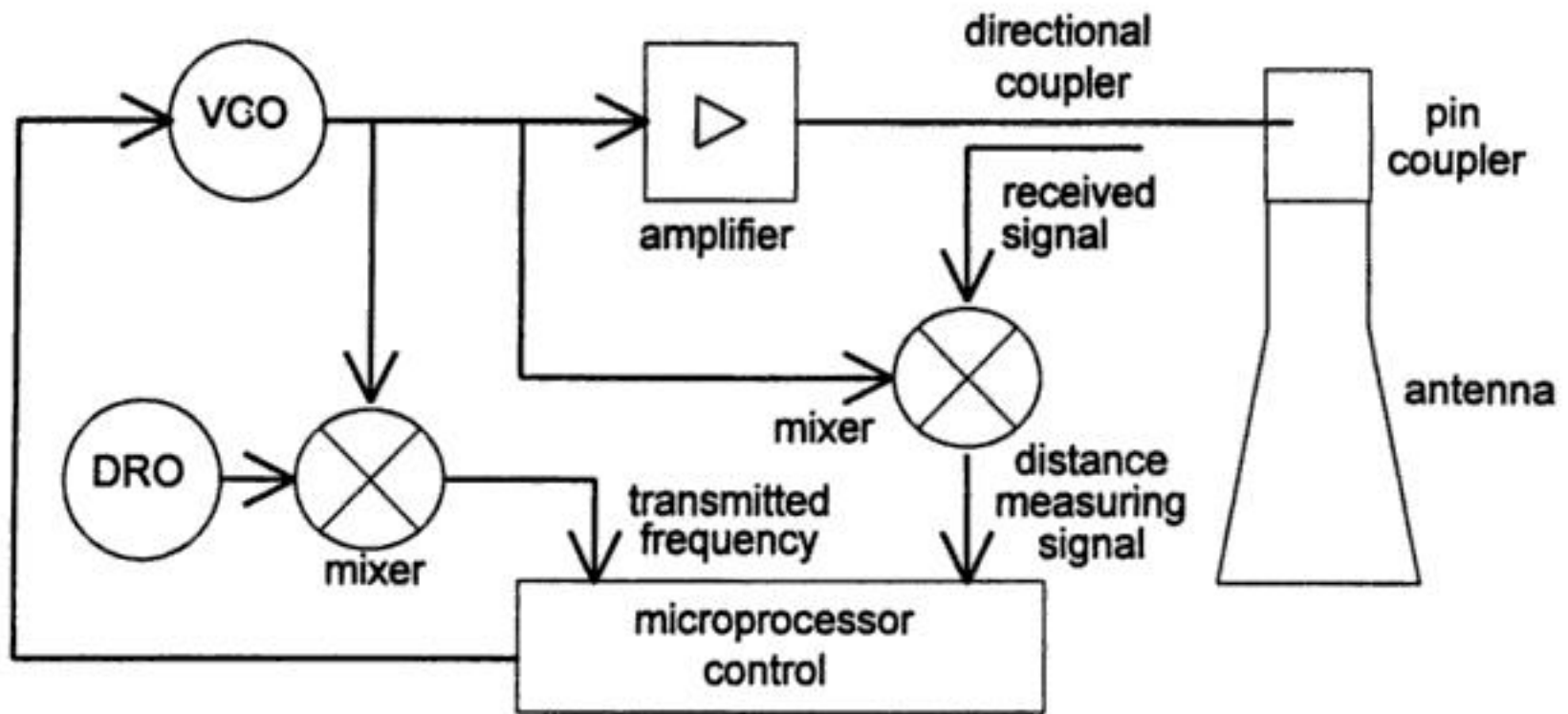


FIGURE 11.11 Basic circuit block diagram of an FMCW radar system: a microprocessor controls a voltage-controlled oscillator (VCO), such that the desired frequency sweep is obtained. This signal is amplified and fed into the transmitting antenna. The instantaneous frequency must be measured in order to ensure good sweep linearity. This is accomplished by counting the frequency after it has been mixed with the known frequency of a dielectric resonance oscillator (DRO). A directional coupler decouples the received signal, which is mixed with the transmission signal and processed by the microprocessor.

evaluation is by means of digital signal processing. For more information about signal processing techniques using spectrum analysis, refer to Chapter 83 of this handbook.

Time-of-Flight Through Product

Alternatively, the propagation time of the waves through a weakly absorbing liquid or bulk material of low permittivity ϵ_r can be measured, as well as the propagation through the air. In cases where the reflection from the interface between air and the upper surface of the product is poor, part of the signal travels through the liquid and is reflected a second time at the tank bottom or at an interface between two liquids (e.g., oil on water).

Figure 11.12 demonstrates this technique. The evaluation is done in the following four steps:

1. Where microwaves in the tank atmosphere of height d are propagated at the speed of light c , microwaves in the medium (relative permittivity = ϵ_r , height L) are propagated at a slower velocity v .
2. Hence, the reflection r_2 from the tank bottom appears to be shifted downward, and the apparent tank height h_v is greater than the true height h .
3. The transit time in the medium is $t_1 = L/v$, where for the same distance in an empty tank would be $t_0 = L/c$. The ratio of apparent "thickness layer" ($h_v - d$) to true filling height ($h - d$) therefore corresponds to the ratio of the wave propagation rates:

$$\frac{h_v - d}{h - d} = \frac{c}{v} = \sqrt{\epsilon_r} \quad (11.17)$$

4. When ϵ_r , h , and h_v are known, distance d and, from that, filling height L can be calculated exactly:

$$L = h - d = \frac{h_v - h}{\sqrt{\epsilon_r} - 1} \quad (11.18)$$

By this method, a direct level measurement — not a distance measurement — is attained. It can even be applied when signal r_1 from the surface of the medium is no longer measurable. The evaluation of

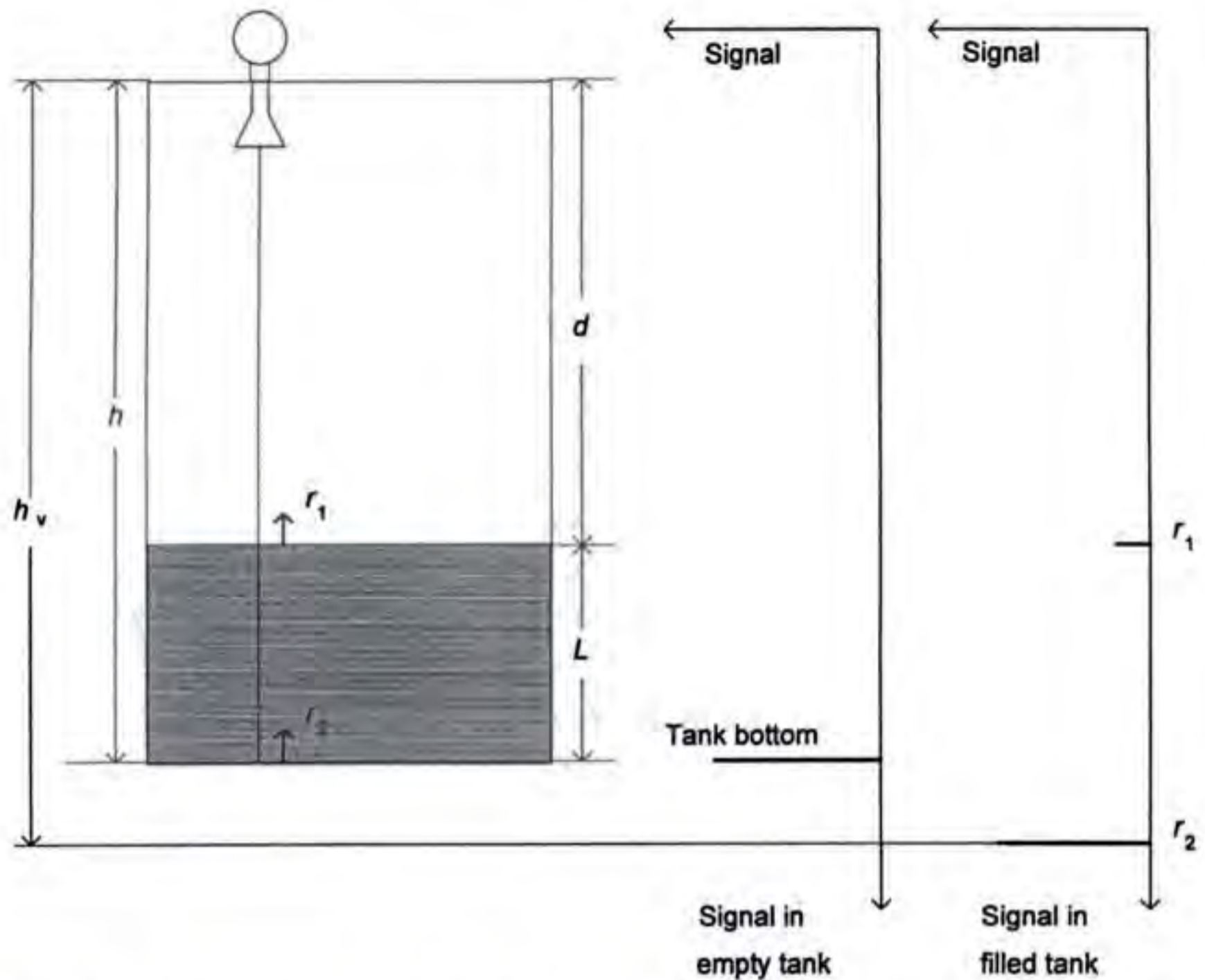


FIGURE 11.12 Representation of time-of-flight measurement through liquid: the wave is reflected once (r_1) at the product's surface and a second time (r_2) at the tank bottom. Due to the reduced wave velocity within the liquid, the reflection r_2 appears below the geometric position of the bottom. From that shift, the level can be calculated; see Equations 11.17 and 11.18.

the tank bottom reflection signal is known as "tank bottom tracing," and is used in the radar level system offered by Krohne in Figure 11.8.

11.3 Level Measurements by Detecting Physical Properties

To measure level, one can detect physical parameters that are significantly different between the atmosphere and the product; for example, conductivity, viscosity, or attenuation of any type of radiation. Two types are possible: (1) continuous measurement with an integral sensor, or (2) switching by a point measurement when the sensor comes in contact with the product.

Electrical Properties

The sensor must be in direct or indirect contact with the product to detect its electrical properties. For continuous measurement, only part of the intrusive sensor must be in contact with the product to detect the difference in dielectric permittivity or conductivity.

Capacitive

In most applications, a rod electrode is arranged vertically in the tank. The electrode can be (1) noninsulated if the liquid is nonconductive, or (2) insulated. The metallic vessel acts as a reference electrode.

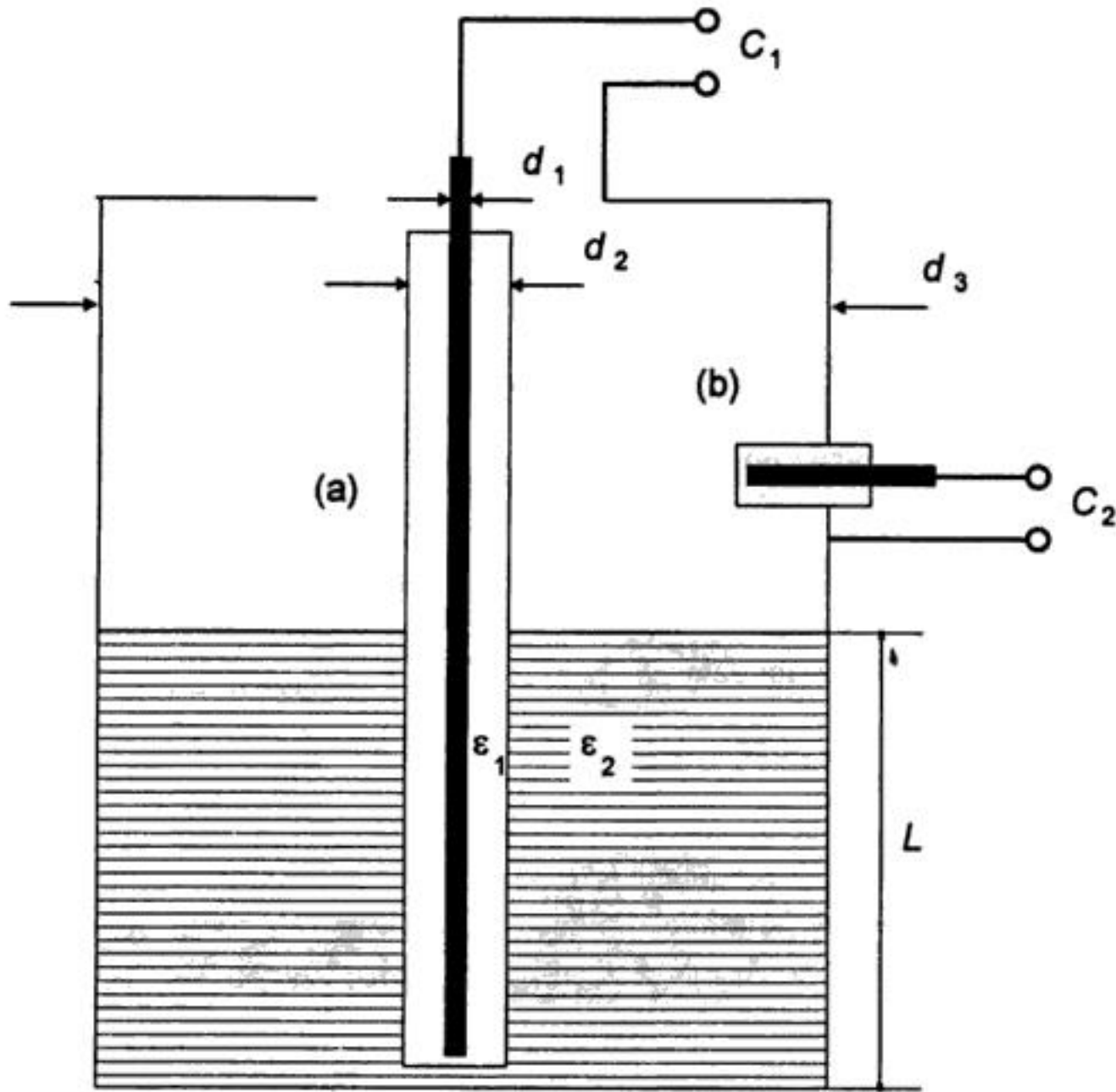


FIGURE 11.13 Principle of operation for a capacitance-type level device. (a) An insulated electrode protrudes into the liquid. The capacitance between the inner conductor and the tank walls is measured. (b) As a capacitance level switch, the electrode can be mounted at the appropriate position directly into the tank wall.

The result depends on the permittivity ϵ_2 of the product. Figure 11.13(a) shows an electrode concentrically mounted on a cylindrical tank. For such a rotationally symmetrical configuration, the capacitance C of an insulated electrode changes with level L according to:

$$C = \frac{2\pi\epsilon_0 L}{\frac{1}{\epsilon_1} \ln \frac{d_2}{d_1} + \frac{1}{\epsilon_2} \ln \frac{d_3}{d_2}} \Leftrightarrow L = \frac{C \left(\frac{1}{\epsilon_1} \ln \frac{d_2}{d_1} + \frac{1}{\epsilon_2} \ln \frac{d_3}{d_2} \right)}{2\pi\epsilon_0} \quad (11.19)$$

ϵ_0 is the dielectric constant of vacuum (8.85×10^{-12} As $V^{-1}m^{-1}$); ϵ_1 and ϵ_2 are the relative permittivities of the insulation material and the liquid, respectively.

If the liquid itself is highly conductive, Equation 11.19 simplifies to:

$$C = \frac{2\pi\epsilon_0 \epsilon_1 L}{\ln \frac{d_2}{d_1}} \Leftrightarrow L = \frac{\ln \frac{d_2}{d_1}}{2\pi\epsilon_0 \epsilon_1} \quad (11.20)$$

If the electrode is not insulated, the following equation is valid:

$$C = \frac{2\pi\epsilon_0\epsilon_2 L}{\ln \frac{d_3}{d_1}} \Leftrightarrow L = \frac{\ln \frac{d_3}{d_1}}{2\pi\epsilon_0\epsilon_2} \quad (11.21)$$

When arranged horizontally, as in Figure 11.13(b), a capacitive sensor can act as a level switch.

For the electrical measurement of capacitance, refer to Chapter 6.3 of this handbook. For a more precise measurement of conductive liquids, a method measuring the complex impedance is helpful.

Conductive

The resistance of the liquid between two electrodes is measured with (1) a strip line with two parallel electrodes similar to Figure 11.9(a), or (2) a rod electrode with the metal vessel as the reference electrode, similar to Figure 11.13(a) without insulator.

Radiation Attenuation

All radiation (e.g., gamma rays, ultrasonic waves, electromagnetic waves) is attenuated to some degree in any medium. In general, attenuation in liquids or bulk materials is higher than in gases. This effect is used to measure level or limits, without direct contact of the sensor.

Radiometric

The intensity I of gamma rays is attenuated by the liquid according to its damping factor α :

$$I = I_0 \exp(-\alpha d) \quad (11.22)$$

The source can be a radioactive material Co-60 or Cs-137, having half-lives of 5.23 years and 29.9 years, respectively. Emitter and sensor may take the form of (1) a point emitting the rays radially in all directions, (2) a rod emitting radially from a line, or (3) an array consisting of several point emitters in a row. Any combination of point/rod/array emitter with point/rod/array detector is possible. Figure 11.14 shows two different configurations. Radiation protection regulations must be considered. A real-time clock in the system must compensate for the decrease of intensity (dose rate) I by time t according to the half-life T_H of the applied material:

$$I = I_0 2^{-t/T_H} \quad (11.23)$$

For more information about radiation detection methods, refer to Chapter 66 of this handbook. Plastic scintillators and counting tubes are preferred for radiometric level gaging. The level-intensity characteristic is nonlinear, so the equipment should be calibrated in place. Mengelkamp [10] describes the radiometric techniques in more detail.

Ultrasonic Switch

A short ultrasonic transmission path can be used to detect products that dampen sonar waves. For instance, this method is applicable for the detection of slurries or to determine the interface between two different liquids. When combined with a servo system, the vertical profile of ultrasonic attenuation can be measured. Another application uses a noncontact sensor mounted on the outside of the vessel. It measures the acoustic impedance through the vessel wall that changes if liquid or gas is present behind the wall.

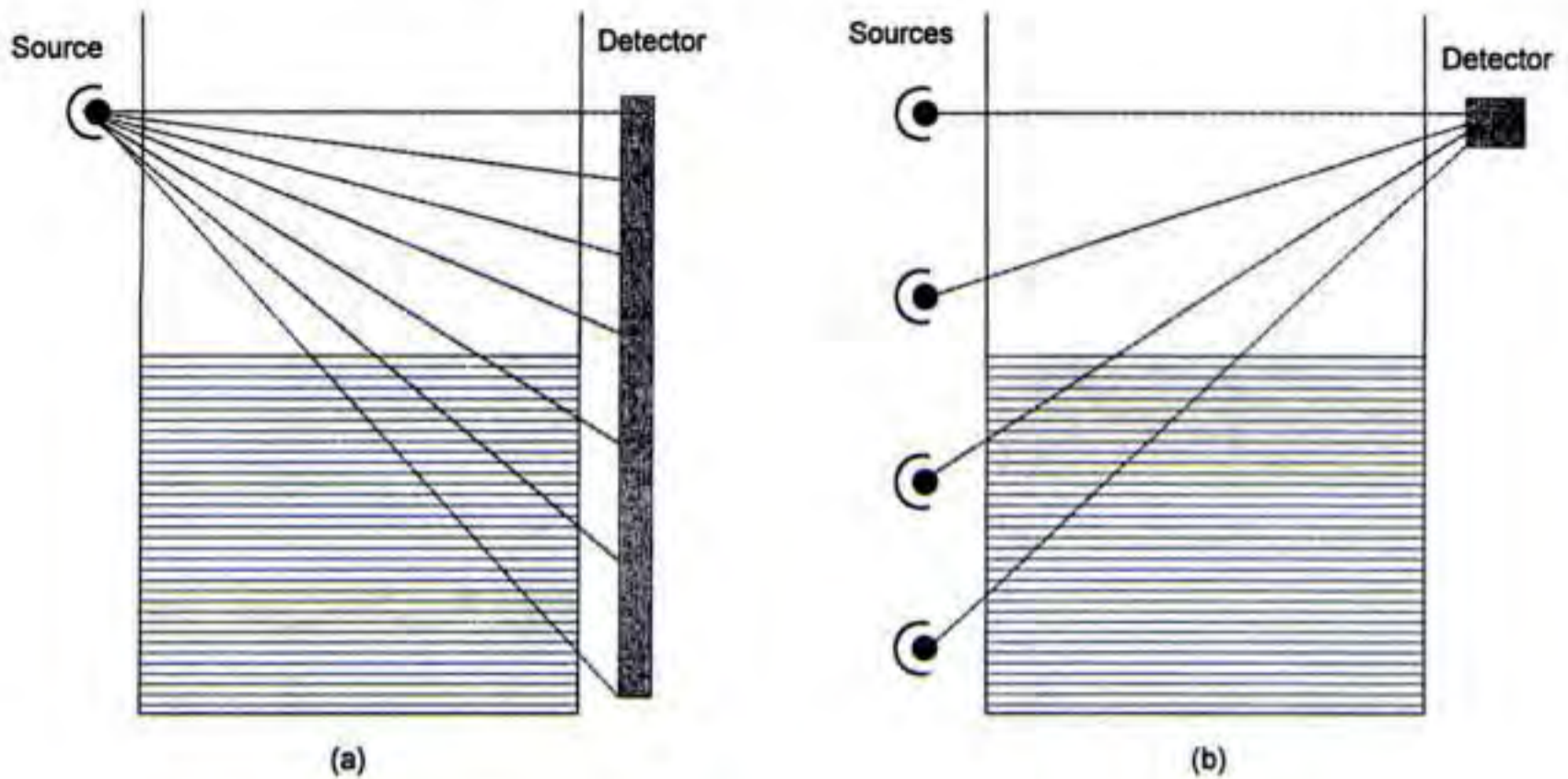


FIGURE 11.14 Representation of a radiometric continuous level gage. The rays are emitted by a radioactive source, propagate through the tank walls, and are damped by the liquid. In (a), a point source is combined with a rod detector (e.g., scintillator rod); (b), a source array is combined with a point detector.

Microwave Switch

Liquids and solids dampen microwaves in many cases, sometimes absorbing them completely. A simple unmodulated microwave source and an accompanying receiver are sufficient for level switching.

Photoelectric Barrier

A photoelectric barrier can act as a level switch for liquids that are not transparent, as well as most solids. But in closed nontransparent tanks, the coupling of the photoelectric components to the tank will not be possible in most cases.

Thermal and Mechanical

For some special applications, level sensors utilize the different heat dissipation properties and viscosities of the media.

Thermal

A self-heated resistor with a high temperature coefficient is immersed into the liquid. Heat dissipation causes the temperature to drop somewhat in the region where the liquid covers the sensor. Therefore, the resistance change is nearly linear with the level. This method is often used in automotive applications. In some applications with heated liquids (e.g., chemical reaction vessels), a simple temperature sensor can be used as a level switch by emitting a signal when the liquid contacts the sensor and heats it.

Viscosity

The effect of viscosity, which is significantly higher for liquids than for gases, dampens the movement of a body. These level sensors measure the degree of damping of a vibrating fork when dipped in a liquid. Normally, it is only used as a point level switch. Figure 11.15 shows such a “tuning fork,” named according to the typical form with two or three vibrating paddles. The integrated electronics evaluate the power loss or the frequency shift of the mechanical resonance system. For solids, a sensor with a rotating paddle that stops when contacting the product is useful.

3. Verein Deutscher Ingenieure, Verband Deutscher Elektrotechniker (VDI/VDE), *Füllstandmessung von Flüssigkeiten und Feststoffen (Level Measurement of Liquids and Solids)*, VDI/VDE 3519, Part 1, Berlin: Beuth, 1984.
4. Verein Deutscher Ingenieure, Verband Deutscher Elektrotechniker (VDI/VDE), *Füllstandmessung von Flüssigkeiten und Feststoffen (Level Measurement of Liquids and Solids)*, VDI/VDE 3519, Part 2, Berlin: Beuth, 1984.
5. K.W. Bonfig (ed.), *Technische Füllstandsmessung und Grenzstandskontrolle*, Ehningen: Expert, 1990.
6. Krohne Messtechnik, *Technical Data Sheets of Level Measurement Products*, Duisburg: Krohne, 1996.
7. K. Blundy, Radar systems — setting a practical approach, *Control & Instrum.*, July 1996.
8. D. Brumbi, *Fundamentals of Radar Techniques for Level Gauging*, Duisburg: Krohne, 1995.
9. D. Brumbi, Measuring process and storage tank level with radar technology, *Int. Radar Conf. IEEE*, 256-260, 1995.
10. B. Mengelkamp, *Radiometrie, Füllstand- und Dichtemessung*, Berlin: Elitera, 1972.

12

Area Measurement

Charles B. Coulbourn

*Los Angeles Scientific
Instrumentation Co.*

Wolfgang P. Buerner

*Los Angeles Scientific
Instrumentation Co.*

12.1	Theory	12-1
	Planimeter • Digitizer • Grid Overlay	
12.2	Equipment and Experiment	12-6
12.3	Evaluation.....	12-11

One must often measure the area of enclosed regions on plan-size drawings. These areas might be either regular or irregular in shape and describe one of the following:

- Areas enclosed by map contours
- Cross section of the diastolic and systolic volumes of heart cavities
- Farm or forest land shown in aerial photographs
- Cross sections of proposed and existing roads
- Quantities of materials used in clothing manufacture
- Scientific measurements
- Swimming pools
- Quantities of ground cover

Tools for this type of measurement include planimeters, digitizer-computer setups, digitizers with built-in area measuring capability, and grid overlay transparencies.

12.1 Theory

Planimeter

A planimeter is a mechanical integrator that consists of a bar (tracer arm), a measuring wheel with its axis parallel to the bar, and a mechanism that constrains the movement of one end of the bar to a fixed track, Figure 12.1. The opposite end of the bar is equipped with a pointer for tracing the outline of an area. The measuring wheel, Figure 12.2, is calibrated with 1000 or more equal divisions per revolution. Each division equals one count. It accumulates counts, P , according to:

$$P = \frac{K}{\pi D} \int \sin \phi ds \quad (12.1)$$

where K = number of counts per revolution of the measuring wheel

D = diameter of the measuring wheel

ϕ = angle between the measuring wheel axis and the direction of travel

s = traced path

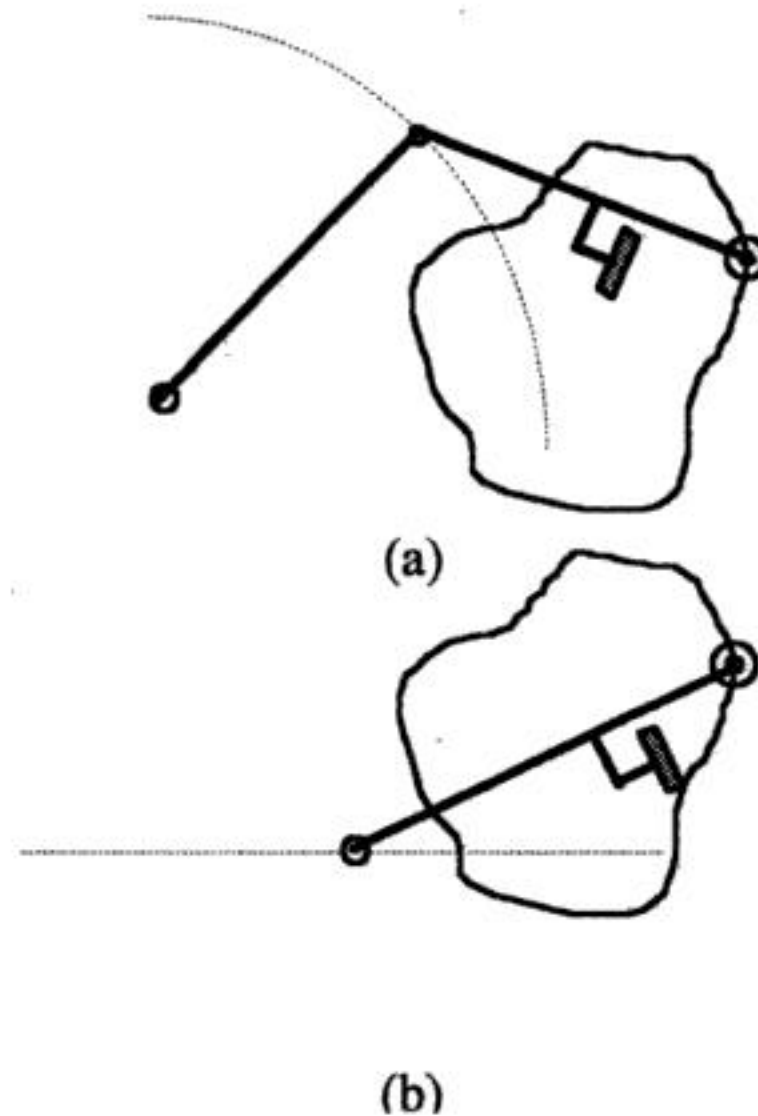


FIGURE 12.1 The constrained end of a polar planimeter (a) follows a circular path; the constrained end of a linear planimeter (b) follows a straight line path.

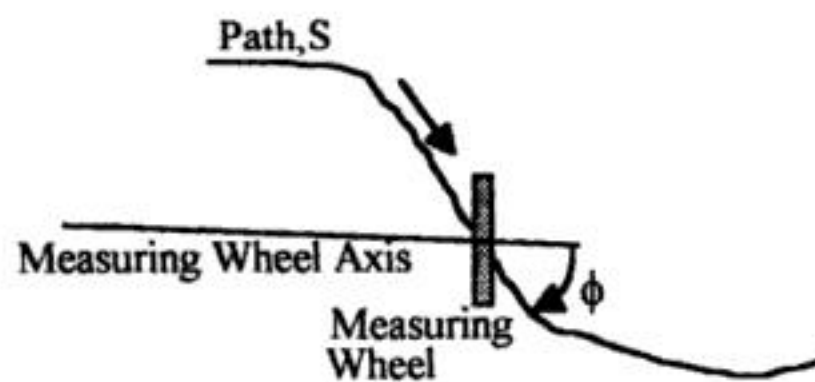


FIGURE 12.2 The rotation of a measuring wheel is proportional to the product of distance moved and the sine of the angle between the wheel axis and direction of travel.

The size of an area, A , traced is:

$$A = \frac{P}{K} \times \pi D \times L \quad (12.2)$$

where L = length of bar

P = accumulated counts (Equation 12.1)

Figure 12.3 shows how a basic wheel and bar mechanism determines the area of a parallelogram. The traced path is along the sloped line; however, the wheel registers an amount that is a function of the product of the distance traveled and the sine of the angle between the direction of travel and the axis of the measuring wheel (Equation 12.1). This is the altitude of the parallelogram. The product of the altitude (wheel reading converted to distance) and base (bar length) is the area.

Figure 12.4(a) illustrates the operation of a planimeter when the area of a four-sided figure is measured. Figures 12.4(b), (c), (d), and (e) show the initial and final positions of the bar as each side of the figure

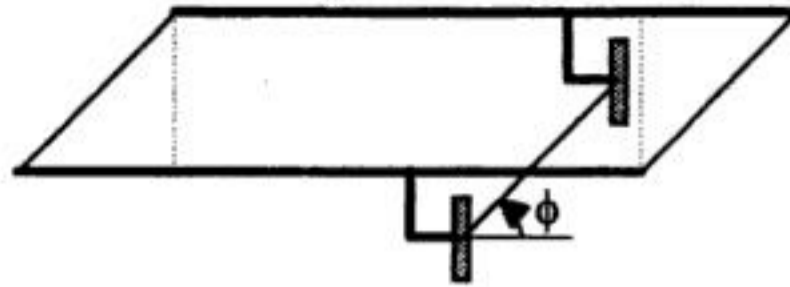
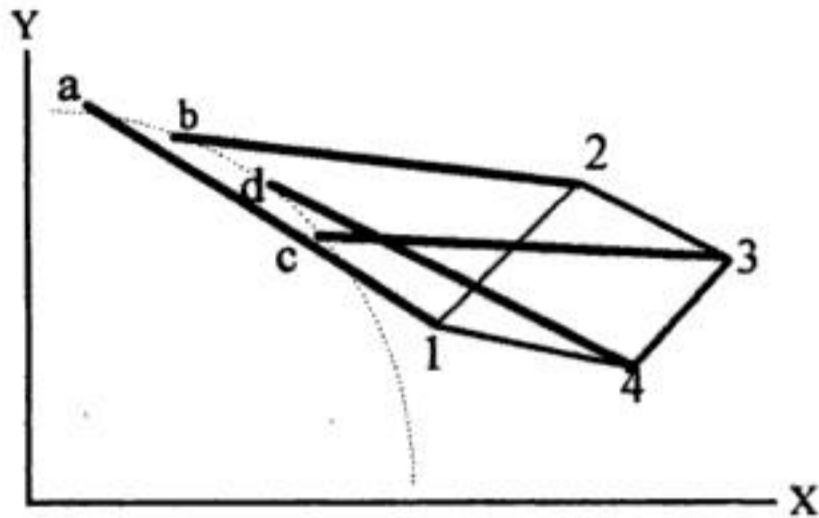
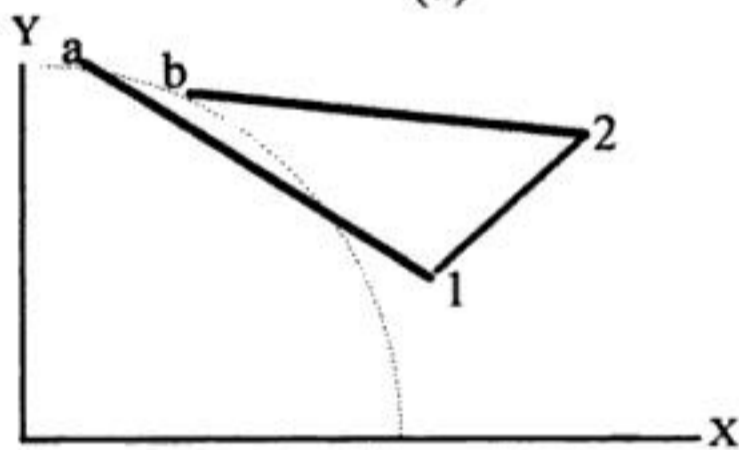


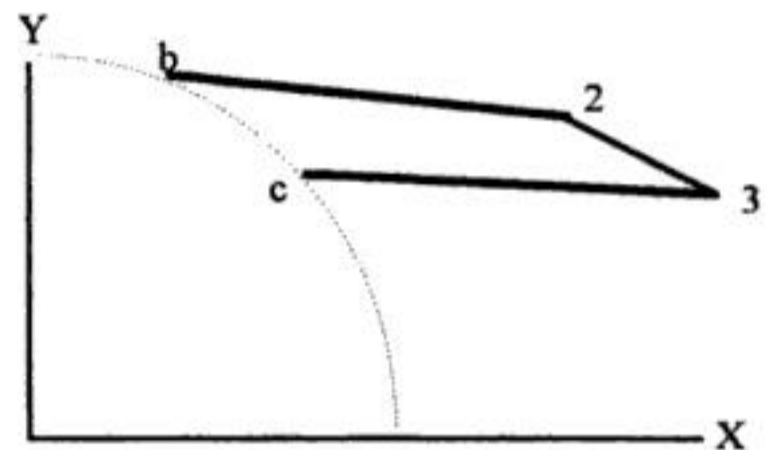
FIGURE 12.3 The area of this parallelogram is proportional to the product of tracer arm length and measuring wheel revolutions.



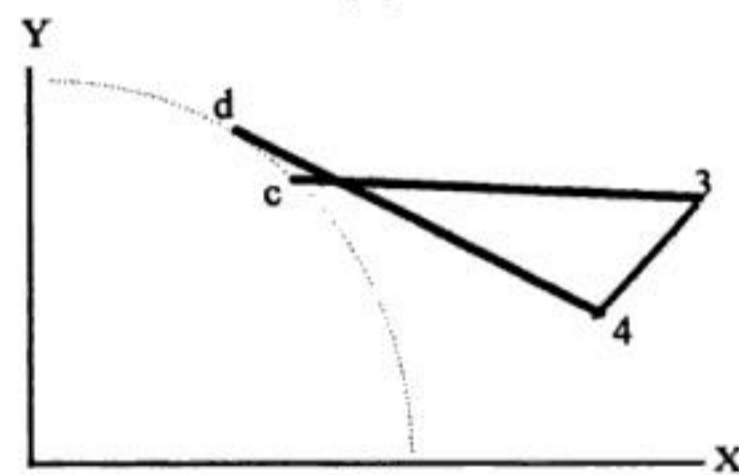
(a)



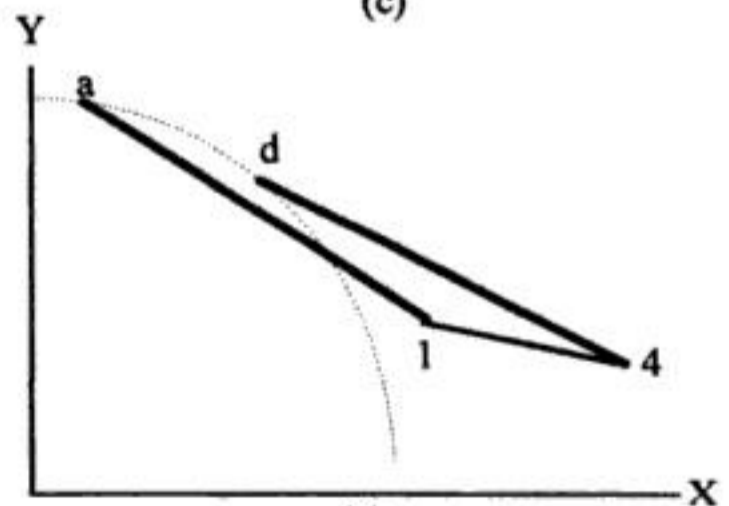
(b)



(c)



(d)



(e)

FIGURE 12.4 This schematic shows a planimeter pointing to each junction of a four-sided figure being traced in (a) and at the ends of each of the segments in (b) through (e). The constrained end of the tracer arm follows a circle.

is traced. Applying the general expression for the area under a curve, $A = \int f(x)dx$ for each of these partial areas gives:

$$A_s = \left(\int_a^1 + \int_1^2 + \int_2^b + \int_b^a \right) f(x)dx \tag{12.3}$$

$$A_b = \left(\int_b^2 + \int_2^3 + \int_3^c + \int_c^b \right) f(x) dx \quad (12.4)$$

$$A_c = \left(\int_c^3 + \int_3^4 + \int_4^d + \int_d^c \right) f(x) dx \quad (12.5)$$

$$A_d = \left(\int_d^4 + \int_4^1 + \int_1^a + \int_a^d \right) f(x) dx \quad (12.6)$$

where A_a , A_b , A_c , and A_d are the four partial areas.

The total area of the figure is the sum of the four partial areas. Combining the terms of these partial areas and rearranging them so that those defining the area traced by the left end of the bar are in one group, those defining the area traced by the other end of the bar are in a second group, and those remaining are in a third group results in the following:

$$A = \left\{ \left[\int_b^a + \int_c^b + \int_d^c + \int_a^d \right] f(x) dx + \left[\int_a^1 + \int_2^b + \int_3^2 + \int_3^c + \int_4^3 + \int_4^d + \int_d^4 + \int_1^a \right] f(x) dx \right. \\ \left. + \left[\int_1^2 + \int_2^3 + \int_3^4 + \int_4^1 \right] f(x) dx \right\} \quad (12.7)$$

The first four integrals describe the area traced by the left end of the bar. Since this end runs along an arc, it necessarily encloses an area equal to zero. The final four integrals describe the area traced by the right end of the bar. This is the four-sided figure. The remaining eight integrals cancel out since $\int_a^1 + \int_1^a = 0$, etc. Thus, the total area equals the area traced by the right end of the bar. Note that the same reasoning applies to figures of any number of sides and of any shape.

Digitizer

A digitizer converts a physical location on a map to digital code representing the (x, y) coordinates of the location. The digital code is normally converted to a standard ASCII or binary format and transmitted to a computer where computations are made to determine such things as area or length. Certain digitizers can also compute areas and lengths without the use of a computer.

Area, A , can be computed using the coordinate pairs that define the area boundary.

$$A = \frac{1}{2}(y_1 + y_2)(x_2 - x_1) + \frac{1}{2}(y_2 + y_3)(x_3 - x_2) + \dots + \frac{1}{2}(y_{n-1} + y_n)(x_n - x_{n-1}) \\ + \frac{1}{2}(y_n + y_1)(x_1 - x_n) \quad (12.8)$$

where x_1, x_2, x_3 , etc. = sequentially measured x coordinates along the boundary.

y_1, y_2, y_3 , etc. = are corresponding y coordinates



FIGURE 12.6 Mechanical planimeters are normally preferred when occasional use is required.

One of the simplest and least costly area measuring devices is the manual polar planimeter (Figure 12.6). It consists of a weight, two arms, a measuring wheel, and a pointer. The *weight* secures one end of a *pole arm*, allowing the other end to rotate along a fixed arc. The rotating end of the pole arm attaches to one end of a *tracer arm* and constrains its movement to the arc. At the other end of the tracer arm is a *pointer* used for tracing the periphery of an unknown area. The *measuring wheel* is located in the box at one end of the tracer arm. The location of the measuring wheel is not critical; however, its axle must be parallel to the tracer arm.

The length of both arms of the planimeter shown in Figure 12.6 can be adjusted. The length of the pole arm has absolutely no effect on measurement accuracy and is adjustable only for convenience. The effective length of the tracer arm directly affects the reading; a shorter arm results in a larger reading. By adjusting the tracer arm length, one can achieve a very limited range of scaling; however, this is usually not done. Rather, the arm length is adjusted according to the general size of areas to be measured: a shorter arm for smaller areas and a longer arm for larger areas. A shorter arm results in a greater number of counts per unit area, which is needed for smaller areas. Scaling is usually done by multiplying the result by an appropriate value.

One zeros the measuring wheel of the planimeter shown in Figure 12.6 by turning a small knurled wheel attached to the measuring wheel axle. On other models, one pushes a plunger to zero the wheel. The latter method is easier but is sensitive to misalignment.

A planimeter pointer is usually a lens with a small circle engraved in the center, although some planimeter models use a needle as the pointer. The best type of pointer is a matter of personal preference although lens pointers are much more popular.

Figure 12.7 shows two electronic planimeters with digital readouts. With these planimeters, one can measure length as well as area. To measure length, snap out the measuring wheel housing, attach an auxiliary handle, and roll the wheel along the line to be measured. Extremely high accuracy can be achieved.

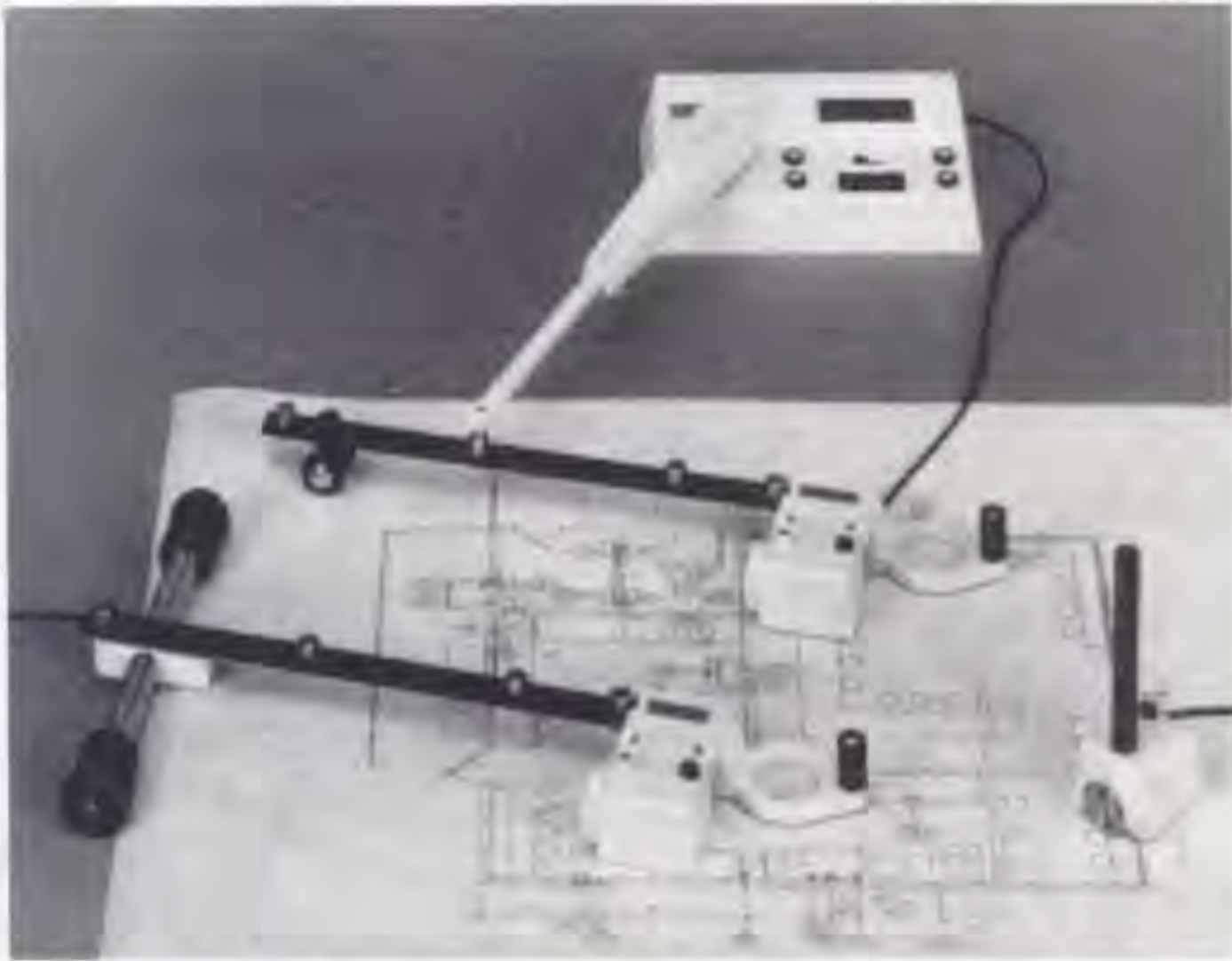


FIGURE 12.7 Electronic planimeters are easier to use and read, and are preferred especially when frequently used.

The upper planimeter in Figure 12.7 is a polar type. It consists of the same parts as the model in Figure 12.6, except that the measuring wheel is attached to a small optical encoder and the weight is packaged in the digital readout. The encoder provides two square-wave outputs that are in phase quadrature, that is, one is 90° out of phase with the other. Both outputs are fed to a processor that counts the pulses and uses the phase difference to determine count direction.

The processor has an electronic scale feature that translates the planimeter reading to real measurement units such as square feet, square meters, acres, or hectares. One can transmit processor data to a computer using an auxiliary interface unit (not shown).

The planimeter at the bottom of Figure 12.7 is a linear model since the path traveled by the constrained end of the tracer arm is a straight line. The straight line path is maintained by a rigid connection between the two carriage wheels and their axle. Other linear planimeters use an actual rail to guide the constrained end of the tracer arm along a straight line path.

The polar planimeter shown at the top of Figure 12.8 and the linear planimeter shown at the bottom of Figure 12.8 are both compact battery-operated models. The measuring wheel is built into the processor, which is attached to the pole arm of the planimeter. The effective length of the tracer arm for both instruments extends from the axis of the constrained end to the pointer, which for these instruments is a small circle engraved in the center of a lens. These planimeters provide electronic scaling and averaging of multiple readings. They cannot be used to measure length.

Figure 12.9 shows an arm digitizer that can be used either as a stand-alone area and length measuring device or to digitize a map or drawing. When operating as a digitizer, the arm digitizer displays the (x, y) coordinates as well as transmits them to a computer. The digitizer has a built-in interface and can transmit using any of over 24 different ASCII and binary codes, each with a choice of parameters. It can be set to measure x and y coordinates in either English or metric units.

Three other modes, in addition to the digitizer mode, are available for computing and displaying either area and length, area and item count, or item count and length. Measurements can also be made in either English or metric units. Any displayed item can be transmitted to a computer through the built-in interface.

ultimate area measuring device will consist of a detached cursor for pointing and a small calculator-like device for operating on the results, displaying them, storing them, and sending them to a computer.

Scanners and associated software will also impact the field of area measurement, particularly as their coverage increases and their price decreases.

Defining Terms

Planimeter: A mechanical integrator used for measuring the area of enclosed regions on maps, plans, etc.

Pole arm: One of the two bars comprising a polar planimeter. One end of the pole arm is fixed and the other end is free to rotate. The length of the pole arm has no effect on the planimeter reading.

Tracer arm: The bar of a planimeter to which is attached the measuring wheel. One end of the tracer arm is constrained to a fixed path, while the other end traces the perimeter of an enclosed region whose area is being measured. The length of the pole arm is indirectly proportional to the planimeter reading.

Polar planimeter: A planimeter with a tracer arm whose constrained end follows a circle.

Linear planimeter: A planimeter with a tracer arm whose constrained end follows a straight line.

Measuring wheel: The planimeter wheel whose degree of rotation is directly proportional to area.

Digitizer: A device to convert data or an image to digital form. The digitizers discussed here convert images to digital form and are categorized as graphic digitizers.

Pointer: The part of a planimeter or digitizer that is used to follow the line being traced.

Resolution element: The smallest elemental area that can be discerned. When referred to in connection with area measurement, it is an area with a value of 1.

References

1. P. A. Santi, J. Fryhofer, and G. Hansen, Electronic planimetry, *Byte*, March 1980, 113–122.
2. K. Mandelberg, Anonymous, *Internet*, 4–96.
3. E. Jones, '89 Planning guide-digitizers, *Architectural & Engineering*, July 1980, 37–40.
4. P. E. Maybaum, Digitizing and computer graphics, *Keyboard*, Sept. 1978, 1–3.

Further Information

F. A. Willers, *Mathematische Instrumente*, Munchen und Berlin: Verlag von R. Oldenbourg, 1943.

13

Volume Measurement

René G. Aarnink
University Hospital Nijmegen
Hessel Wijkstra

13.1	Plethysmography Theory.....	13-2
	Air/Water Plethysmography or Chamber Plethysmography • Electrical Plethysmography	
13.2	Numerical Integration with Imaging.....	13-6
13.3	Indicator Dilution Methods	13-11
	Thermodilution • Radionuclide Techniques • Gas Dilution	
13.4	Water Displacement Volumetry	13-14
13.5	Equipment and Experiments	13-16
13.6	Evaluation.....	13-16

For simple geometric shapes, volume measurements can be performed analytically by measuring the dimensions of the object in question and using the appropriate formula for that shape. Volume formulae are found in geometry and calculus textbooks as well as in reference books such as CRC Press's *Handbook of Mathematical Science* or other references.

Volume can also be measured by fluid displacement. The object whose volume is to be measured is placed in a container filled with fluid and the initial and final volumes are measured. The object's volume is equal to the final volume minus the initial volume. This technique is especially useful for irregularly shaped objects.

Fluids can also be used to measure volume of cavities within objects by filling the cavity entirely with a fluid and then measuring the volume of the fluid after it is removed from the cavity.

The remainder of this chapter is dedicated to the more specific problem of volume measurements in medical applications.

Quantitative volume information can be of importance to the clinician in the diagnosis of a variety of diseases or abnormalities. It can also improve the understanding of the physiology of the patient. Information on volume may be used in various applications such as cardiac monitoring, diagnosis of prostate diseases, follow-up during organ transplantation, surgery for tumor resection, blood flow measurements, plastic surgery, follow-up of preterm infants, sports performance analysis, etc. Because of this wide spectrum of applications, various techniques for volume measurements have been developed, some of which are useful in determining the amount of blood flow to the organ (dynamic), while others are used to obtain the size of the object (static). The techniques can be invasive or noninvasive, and are based on either direct or indirect measurements. Each technique has its own advantages and disadvantages, and the application determines the selection of volume measurement method.

One of the earliest techniques to measure (changes of) body volume was *plethysmography*, originally developed by Glisson (1622) and Swammerdam (1737) to demonstrate isovolumetric contraction of isolated muscle. The measuring technique consists of surrounding the organ or tissue with a rigid box filled with water or air. The displacement of the fluid or the change in air pressure indicates the volume

changes of the tissue due to arterial in-flow. Two major types of plethysmography exist, and these can be distinguished by the technique used to measure the volume change. These are *volume plethysmography* (direct-measurement displacement plethysmography including water and air types), and *electrical plethysmography* (strain-gages, inductive and impedance plethysmographs). The physical condition in which the measurements should be performed determines the plethysmographic method chosen.

Advances in medical imaging have provided new possibilities for noninvasively extracting quantitatively useful diagnostic information. Images are constructed on a grid of small picture elements (pixels) that reflect the intensity of the image in the array occupied by the pixel. Most medical images represent a two-dimensional projection of a three-dimensional object.

Currently, the most commonly used medical imaging modalities are ultrasound, nuclear magnetic resonance imaging, X-rays and X-ray computer tomography. *Ultrasound imaging* is based on the transmission and reflection of high-frequency acoustic waves. Waves whose frequencies lie well above the hearing range can be transmitted through biological tissue and will be reflected if they cross a boundary between media of different acoustic properties. These reflected signals can be reconverted to electrical signals and displayed to obtain a two-dimensional section. *Magnetic resonance imaging* is based on the involvement of the interaction between magnetic moment (or spin) of nuclei and a magnetic field. The proton spins are excited by an external radio frequency signal and the return to an equilibrium distribution is used to construct cross-sectional images. *Computer tomography* or *CT-scanning* uses a rotating source of X-rays. The X-ray source and detector are scanned across a section of the object of interest, and the transmission measurements are used to reconstruct an image of the object.

These imaging techniques display cross-sections of views that can be used to estimate the size of specific components. One way to estimate the volume of internal objects is make certain assumptions concerning the shape and to apply formulae to estimate the volume with dimensions of the object such as length, height, and width. A more accurate technique is *step-section planimetry*, a clinical application of numerical integration. During this procedure, cross-sections of the object are recorded with a certain (fixed) interval. The area of the object is determined in every section, and the total volume is calculated by multiplying the contribution of each section with the interval and summarizing all contributions.

The volume of a fluid-filled region can be calculated if a known quantity of indicator is added to the fluid and the concentration measured after it has been dispersed uniformly throughout the fluid. The selection of the indicator used to measure the volume depends on the application and can be based on temperature, color, or radioactivity.

Finally, volume (changes) can also be performed directly using *water-displacement volumetry*, a sensitive but time-consuming method to measure the volume of an extremity. This method is not suitable for patients in the immediate postoperative period.

13.1 Plethysmography Theory

Fluid in-flow can be measured by blocking the out-flow and then measuring the change in volume due to the in-flow over a defined time interval. Plethysmography enables the volume change to be determined in a reasonably accurate manner. A direct-volume displacement plethysmograph uses a rigid chamber usually filled with water, into which the limb or limb segment is placed. This type of plethysmograph, sometimes called *chamber plethysmography*, can be used in two ways. It can be used to measure the sequence of pulsations proportional to the individual volume changes with each heart beat (arterial plethysmography or rheography). Also, the total amount of blood flowing into the limb or digit can be measured by venous occlusion: by inflating the occluding cuff placed upstream of the limb or digit just above the venous pressure 5 kPa to 8 kPa (40 mm Hg to 60 mm Hg), arterial blood can enter the region but venous blood is unable to leave. The result is that the limb or digit increases its volume with each heart beat by the volume entering during that beat.

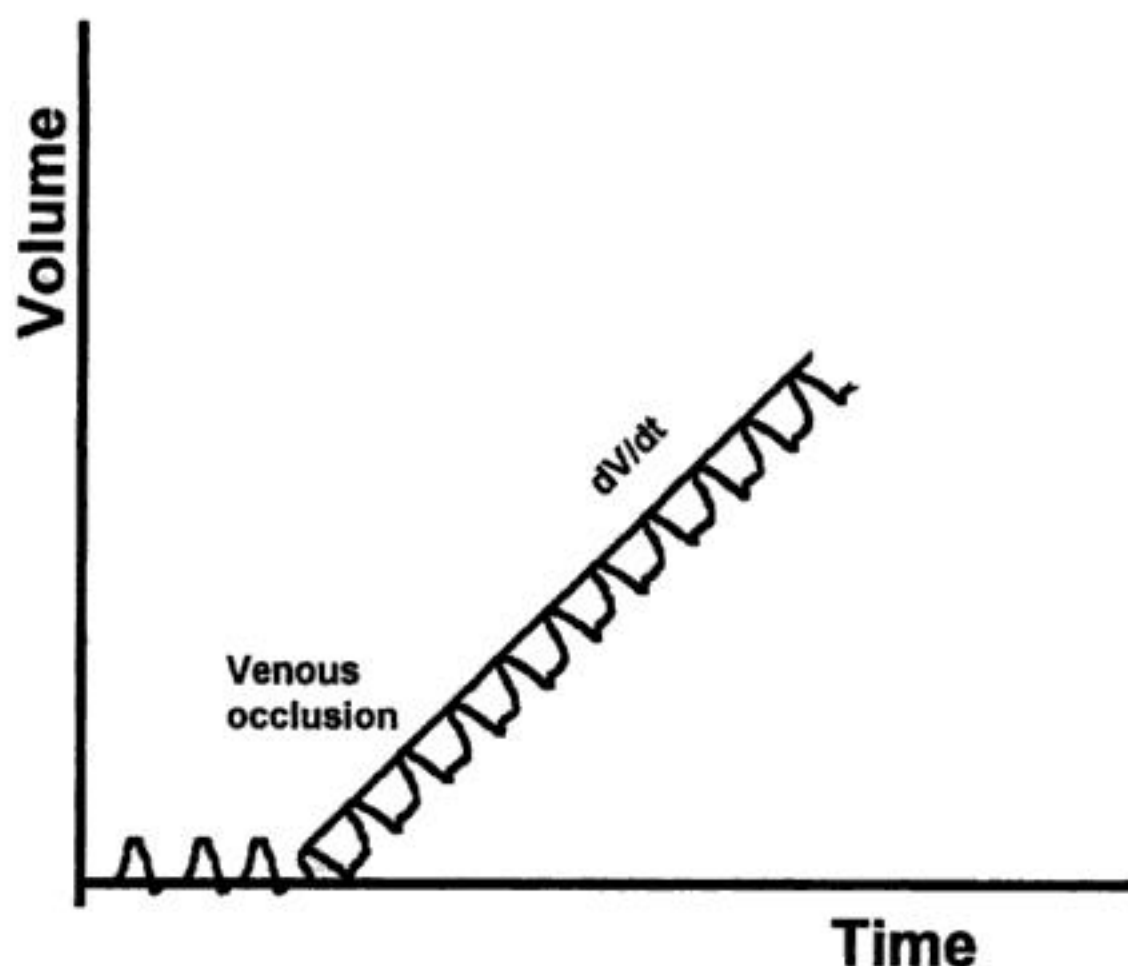


FIGURE 13.1 Typical recording from air plethysmograph during venous occlusion, revealing the changes in volume over time caused by arterial in-flow.

Air/Water Plethysmography or Chamber Plethysmography

Air plethysmography uses the relation between the volume change of a gas in a closed system and the corresponding pressure and temperature [1]. The relation between pressure and volume is described by Boyle's law, which can be written as:

$$P_i V_i = P_f V_f \quad (13.1)$$

with P_i and V_i the initial pressure and volume, respectively, and P_f and V_f the final pressure and volume, measured at constant temperature: the displacement of the fluid or the compressing of the air is a direct measure of the blood flow or original volume. The air plethysmograph uses the change in pressure that occurs in a cuff wrapped around the segment of interest due to the change in volume. By inflating the cuff to about 60 mm Hg, the arterial in-flow causes small increases in pressure. These small changes in pressure over the cardiac cycle can be monitored (see Figure 13.1). The measurement of blood flow is achieved by comparing the pressure changes to changes caused by removing known amounts of air from the system. The second measurement uses volume changes at various pressures between systolic and diastolic pressures, and the peak deflection is compared with the deflection caused by removal of known amounts of air from the system. In segmental plethysmography, two cuffs are used to measure the volume changes in a segment of a limb. Venous occlusion is established by the first cuff, while the second is inflated to a pressure that will exclude blood flow from those parts that should not be included in the measurement.

The technique can also be used for whole-body plethysmography [2], a common technique used to measure residual volume, functional residual capacity (FRC), and total lung volume, the parameters usually determined during pulmonary function testing [3]. In body plethysmography, the patient sits inside an airtight box, and inhales or exhales to a particular volume (usually the functional residual capacity, the volume that remains after a normal exhalation). Then, a shutter closes the breathing tube,

while the subject tries to breathe through the closed tube. This causes the chest volume to expand and decompress the air in the lungs. The increase in chest volume slightly reduces the air volume of the box, thereby increasing the pressure. First, the change in volume of the chest is quantified using Boyle's law with the initial pressure and volume of the box, and the pressure in the box after expansion. The change in volume of the box is equal to the change in volume of the chest. A second measurement using the initial volume of the chest (unknown) and the initial pressure at the mouth, and the inspiratory volume (the unknown chest volume and the change in volume obtained in the first measurement) together with the pressure at the mouth during the inspiratory effort. By solving Boyle's law for the unknown volume, the original volume of gas present in the lungs when the shutter was closed is obtained, which is normally the volume present at the end of a normal exhalation or FRC.

Alternative methods for measurement of volume changes with plethysmography have been introduced to overcome the disadvantages of chamber plethysmography such as cost, complexity, and awkwardness of use. Important alternatives are elastic-resistance strain-gage plethysmography, impedance plethysmography, inductive plethysmography, and photoelectric plethysmography.

Electrical Plethysmography

Strain-Gage Plethysmography

In 1953, Whitney described the elastic-resistance strain-gage plethysmograph [4]. The strain-gage instrument quantifies the change in resistance of the gage as it is stretched due to the enlargement of the object (e.g., limb segment). The gage is typically a small elastic tube filled with mercury or an electrolyte or conductive paste. The ends of the tube are sealed with copper plugs or electrodes that make contact with the conductive column. These plugs are connected to a Wheatstone bridge circuit.

To illustrate the principle of the strain-gage method, a strain-gage of length l_0 is placed around the limb segment, a circular cross section with radius r_0 . The length of the strain-gage can thus also be expressed as $l_0 = 2\pi r_0$. Expansion of the limb gives a new length $l_1 = 2\pi r_1$, with r_1 the new radius of the limb. The increase in length of the strain-gage is thus $\delta l = 2\pi(r_1 - r_0)$, while the change in cross-sectional area of the limb δA is $\pi(r_1^2 - r_0^2)$. This can also be expressed as:

$$\delta A = \pi \left[2r_0(r_1 - r_0) + (r_1 - r_0)^2 \right] \quad (13.2)$$

Since $r_1 - r_0$ is usually small, $(r_1 - r_0)^2$ can be neglected. Consequently, δA can be written as $2\pi r_0(r_1 - r_0)$ which, on dividing by $A = \pi(r_0)^2$, gives:

$$\frac{\delta A}{A} = 2 \frac{r_1 - r_0}{r_0} \quad (13.3)$$

But since $\delta V/V = \delta A/A$ and $\delta l/l = \delta r/r$:

$$\frac{\delta V}{V} = 2 \frac{l_1 - l_0}{l_0} \quad (13.4)$$

Thus, the percentage increase in volume can be obtained by measuring the initial gage length and the change in length. In practice, this change in length of the strain-gage is recorded over a relatively short period to overcome the problem of back pressure that builds up in the limb or digit because of the venous occlusion, and the measurement is usually expressed as milliliters per minute per 100 g of tissue, or: volume flow rate = $2(\delta l/l_0) \times (100/t)$, where δl is the increase in gage length during time t . To measure $\delta l/t$, the gage must be calibrated by removing it from the limb and stretching it on a measuring jig until

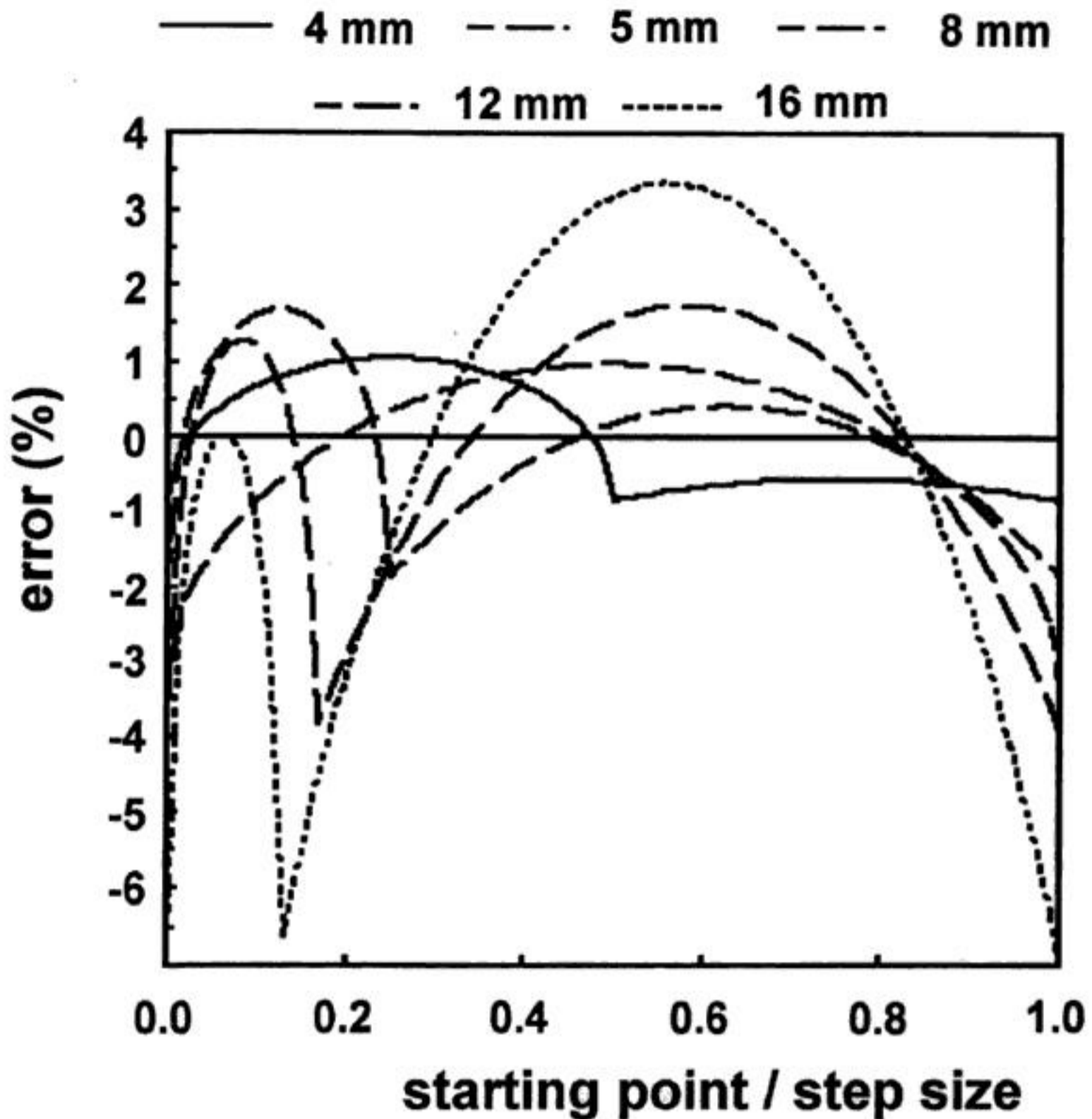


FIGURE 13.7 The errors in surface estimates by numerical integration for different step sizes h presented as percentage of the exact solution, as a function of the position of the first section, located between 0 and h . (From R. G. Aarnink et al., *Physiol. Meas.*, 16, 141-150, 1995. With permission.)

error in volume estimation is less than 5%. Although clinical application might introduce additional errors such as movement artifacts, this rule provides a good indication of the intersection distance that should be used for good approximations with numerical integration.

13.3 Indicator Dilution Methods

The principle of the indicator dilution theory to determine gas or fluid volume was originally developed by Stewart and Hamilton. It is based on the following concept: if the concentration of an indicator that is uniformly dispersed in an unknown volume is determined, and the volume of the indicator is known, the unknown volume can be determined. Assuming a single in-flow and single out-flow model, all input will eventually emerge through the output channel, and the volumetric flow rate can be used to identify the volume flow. It can be described with two equations, of which the second is used to determine the volume using the result of the first equation [16]:

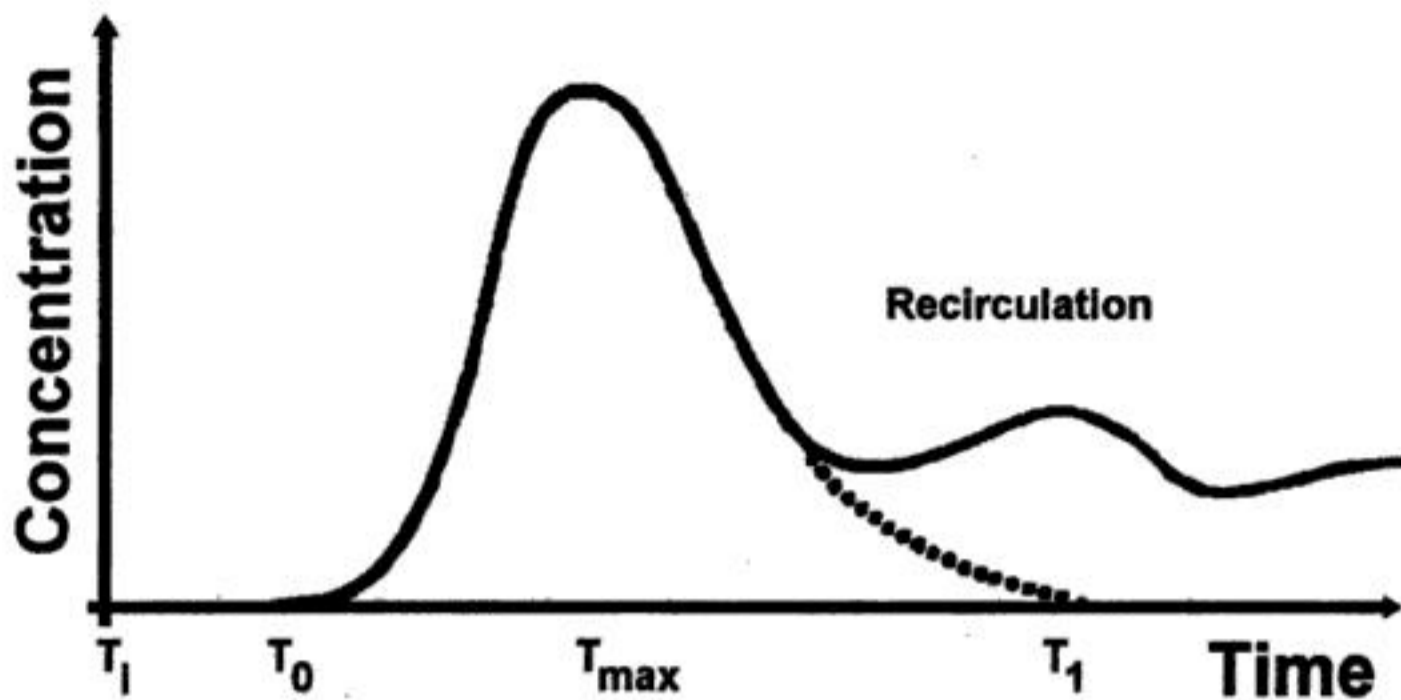


FIGURE 13.8 The time-concentration curve of the indicator is used to obtain the unknown volume of the fluid. The area under the time-concentration curve is needed to estimate the volumetric flow rate, and this flow rate can be used to estimate the volume flow. For accurate determination of the flow rate, it is necessary to extrapolate the first pass response before the area under the curve is estimated.

$$F = \frac{I_i}{\int_0^{\infty} C(t) dt} \quad (13.18)$$

$$\text{Volume} = F t_m \quad (13.19)$$

where F = flow rate
 I_i = total quantity of injected indicator
 $C(t)$ = concentration of the contrast as function of time
 Vol = volume flow
 t_m = mean transit time

The integrated value of the concentration over time can be obtained from the concentration-time curve (Figure 13.8) using planimetry and an approximation described by the Riemann sum expression (Equation 13.13). Care must be taken to remove the influence of the effect of recirculation of the indicator. An indicator is suitable for use in the dilution method if it can be detected and measured, and remains within the circulation during at least its first circuit, without doing harm to the object. Subsequent removal from the circulation is an advantage. Dyes such as Coomassie blue and indocyanine green were initially used, using the peak spectral absorption at certain wavelengths to detect the indicator. Also, radioisotopes such as I-labeled human serum albumin have been used. Currently, the use of cooled saline has become routine in measurement of cardiac output in the intensive care unit.

Thermodilution

To perform cardiac output measurements, a cold bolus (with a temperature gradient of at least 10E) is injected into the right atrium and the resulting change in temperature of the blood in the pulmonary artery is detected using a rapidly responding thermistor. After recording the temperature-time curve, the blood flow is calculated using a modified Stewart-Hamilton equation. This equation can be described as follows [17]:

$$V_i SW_i SH_i (T_b - T_i) = F \int T_b (dt) SW_b SH_b \quad (13.20)$$

The terms on the left represent the cooling of the blood (with temperature T_b), as caused by the injection of a cold bolus with volume V_i , with temperature T_i , specific weight SW_i and specific heat SH_i . The same amount of "indicator" must appear in the blood (with specific weight and heat SW_b and SH_b , respectively) downstream to the point of injection, where it is detected in terms of the time-course of the temperature T_b in the flow F by means of a thermistor. The flow can be described as:

$$F = \frac{T_b - T_i}{\int T_b (dt)} K \quad (13.21)$$

where K (the so-called calculation constant) is described as:

$$K = V_i \frac{SW_i - SH_i}{SW_b - SH_b} C \quad (13.22)$$

This constant is introduced to the computer by the user. T_i is usually measured at the proximal end of the injection line; therefore, a "correction factor" C must be introduced to correct for the estimated losses of cold saline in the catheter, due to the dead space volume. Furthermore, the warming effect of the injected fluid during passing through the blood must be corrected; this warming effect depends on the speed of injection, the length of immersion of the catheter, and the temperature gradient. This warming effect reduces the actual amount of injected fluid, so the correction factor C is less than 1 [17].

Radionuclide Techniques

Radionuclide imaging is a common noninvasive technique used in the evaluation of cardiac function and disease. It can be compared to X-ray radiology, but now the radiation emanating from inside the human body is used to construct the image. Radionuclide imaging has proven useful to obtain ejection fraction and left ventricular volume. As in thermodilution, a small volume of indicator is injected peripherally; in this case, radioisotopes (normally indicated as radiopharmaceuticals) are used. The radioactive decay of radioisotopes leads to the emission of alpha, beta, gamma, and X radiation, depending on the radionuclide used. For *in vivo* imaging, only gamma- or X-radiation can be detected with detectors external to the body, while the minimal amount of energy of the emitted photons should be greater than 50 keV. Using a photon counter such as a Gamma camera, equipped with a collimator, the radioactivity can be recorded. Two techniques are commonly used: the first pass method and the dynamic recording method [11].

First Pass Method

To inject a radionuclide bolus into the blood system, a catheter is inserted in a vein, for example, an external jugular or an antecubital vein. The camera used to detect the radioactive bolus is usually positioned in the left anterior oblique position to obtain optimal right and left ventricular separation, and is tilted slightly in an attempt to separate left atrial activity from that of the left ventricle. First, the background radiation is obtained by counting the background emissions. Then, a region of interest is determined over the left ventricle and a time-activity curve is generated by counting the photon emissions over time. The radioactivity count over time usually reveals peaks at different moments; the first peak occurs when the radioactivity in the right ventricle is counted, the second peak is attributed to left ventricular activity. More peaks can occur during recirculation of the radioactivity. After correction for

background emissions, the time-intensity curve is evaluated to determine the ejection fraction (EF) of the heart, which is given by:

$$EF = \frac{c_d - c_s}{c_d - c_b} \quad (13.23)$$

where c_d is the end diastolic count, c_s is the end systolic count, and c_b is the background count [18]. Using a dynamic cardiac phantom, the accuracy of the EF measurement and LV volume estimation by radio-nuclide imaging has been determined [18]. The count-based method was found to give accurate results for evaluating cardiac EF and volume measurements.

Dynamic Recording Method

During dynamic recording, the data acquisition is triggered by another signal, usually the ECG, to obtain the gated activity during ventricular systole and diastole. The contraction of the left ventricle of the heart is used to align the acquisition data from different cardiac cycles. During dynamic recording, information about the cardiac function is averaged over many heart beats during a period when the radiopharmaceuticals are uniformly distributed in the blood pool. Estimations of end systolic and end diastolic volumes are made and the ejection fraction is calculated (see above). Also, other parameters can be obtained; for example, the amplitude and phase of the ejection from the left ventricle. These parameters may show malfunction in the cardiac cycle. These dynamic recordings may also be applied to the brain, lungs, liver, kidney, and vascular systems and they can be important to judge the function of these organs (e.g., after kidney transplantation).

Gas Dilution

Gas dilution is a method to determine lung volumes because standard techniques for lung measurements measure only the inhaled or exhaled air as a function of time (spirometry) and cannot be used to assess the absolute lung volume. The subject is connected to a spirometer that contains a known concentration of gas (e.g., helium). The subject is then asked to breathe for several minutes, to equalize the concentration of helium in the lung and the spirometer. Using the law of conservation of matter, the volume of the lung can be calculated. Since the total amount of helium is the same before and after measurement, the fractional concentration times the volume before equals the fractional concentration times the volume after: $C_1 \times V_1 = C_2 \times (V_1 + V_2)$. The volume after the measurement can be extracted from this equation and, by subtracting the volume of the spirometer, the lung volume is calculated.

13.4 Water Displacement Volumetry

The measurement of the volume of an extremity such as arm or leg can be measured using a water tank, which is illustrated schematically in Figure 13.9. The advantage of water displacement volumetry is the possibility for direct measurement of objects with an irregular form. In the clinical situation, volume determination of the leg can be valuable for monitoring the severity of edema or hematoma after surgery or trauma. A setup for water displacement can be developed by the user, and an example of such a system is described by Kaulesar Sukul et al. [19]. A special tank for water displacement volumetry of the leg was constructed consisting of two overflow tubes. The tank was filled to the lower overflow tube, and the overflow tube was subsequently closed. The patient was then asked to lower his leg into the tank, and the amount of overflow of the upper tube was measured. The total volume was then calculated by measuring the fluid volume delivered to the upper tube and the volume difference between the lower and upper overflow tube, the so-called "reserve volume." The volume of ankle and foot was then measured by filling the tank to the upper overflow tube and measuring the amount of water in the cylinder when the foot and ankle were immersed in the water.

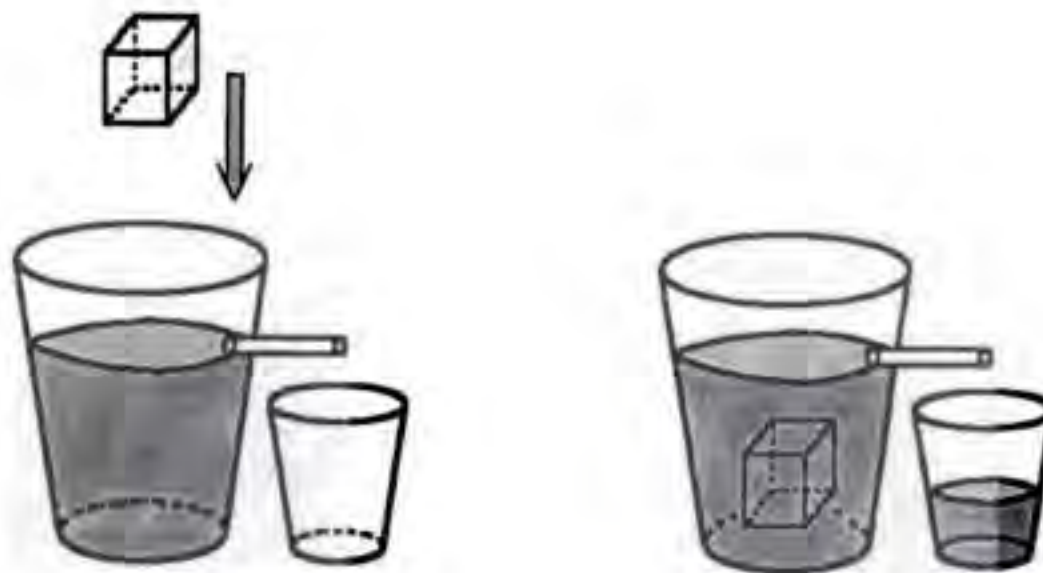


FIGURE 13.9 A schematic overview of a volume measuring method using water displacement: the desired volume of the object is represented by the overflow volume in the small tank after lowering the object into the large tank. This technique is especially useful for irregular-shaped objects.

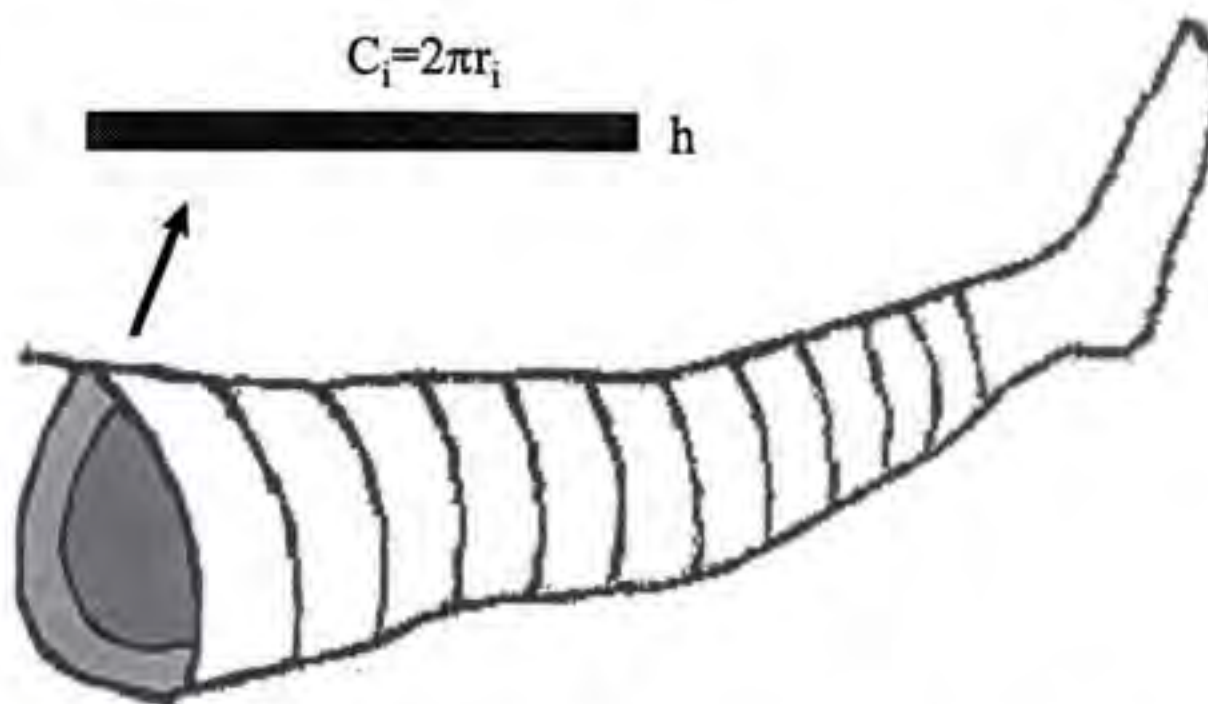


FIGURE 13.10 Illustration of the disk model to measure leg volume. The leg is divided in different sections with a fixed intersection distance the circumference of each section is obtained. This circumference is then used to describe the leg as perfect circles with an assumed radius calculated from the circumference. This radius is then used to calculate the contribution of that section to the volume and the total volume is obtained by summation.

Disadvantages of water displacement volumetry are hygiene problems, it is time-consuming, and not suitable for measurements of the volume of extremities of patients in the immediate postoperative period. Therefore, alternatives have been sought, with the disk model method as a promising one. The calculation of the volume of the leg is performed by dividing the leg into disks of thickness h (e.g., 3 cm) as illustrated in Figure 13.10. The total volume is equal to the sum of the individual disk volumes:

$$V = \sum_{i=1}^n \frac{C_i^2}{4\pi} h = \sum_{i=1}^n \pi r_i^2 h = h \sum_{i=1}^n \pi r_i^2 \quad (13.24)$$

where C_i is the circumference of the disk at position i with assumed radius of r_i .

A study to compare water displacement volumetry with the results of the disk model method indicated that both methods give similar results. Consequently, because of the ease of application, the disk model method is the method of choice to measure volumes of extremities. Assuming the length of a leg to be 75 cm, a ratio between length and step size (3 cm was proposed in [19]) of 25 is obtained. According to

6. W. G. Kubicek, F. J. Kottke, M. V. Ramos, R. P. Patterson, D. A. Witsoe, J. W. Labree, W. Remole, T. E. Layman, H. Schoening, and J. T. Garamala, The Minnesota impedance cardiograph — Theory and applications, *Biomed. Eng.*, 9, 410-416, 1974.
7. R. Shankar and J. G. Webster, Noninvasive measurement of compliance of human leg arteries, *IEEE Trans. Biomed. Eng.*, 38, 62-67, 1993.
8. N. Verschoor, H. H. Woltjer, B. J. M. van der Meer, and P. M. J. M. de Vries, The lowering of stroke volume measured by means of impedance cardiography during endexpiratory breath holding, *Physiol. Meas.*, 17, 29-35, 1996.
9. M. A. Cohn, H. Watson, R. Weisshaut, F. Stott, and M. A. Sackner, A transducer for non-invasive monitoring of respiration, in *ISAM 1977, Proc. Sec. Int. Symp. Ambulatory Monitoring*, London: Academic Press, 1978, 119-128.
10. J. A. Adams, Respiratory inductive plethysmography, in J. Stocks, P. D. Sly, R. S. Tepper, and W. J. Morgan (eds.), *Infant Respiratory Function Testing*, New York: Wiley-Liss, 1996, 139-164.
11. S. Webb, *The Physics of Medical Imaging*, Bristol, U.K.: IOP Publishing, 1988, 204-221.
12. M. K. Terris and T. A. Stamey, Determination of prostate volume by transrectal ultrasound, *J. Urol.*, 145, 984-987, 1991.
13. R. G. Aarnink, J. J. M. C. H. de la Rosette, F. M. J. Debruyne, and H. Wijkstra, Formula-derived prostate volume determination, *Eur. Urol.*, 29, 399-402, 1996.
14. P. J. Davis and P. Rabinowitz, *Methods of Numerical Integration*, San Diego: Academic Press, 1975, 40-43.
15. R. G. Aarnink, R. J. B. Giesen, J. J. M. C. H. de la Rosette, A. L. Huynen, F. M. J. Debruyne, and H. Wijkstra, Planimetric volumetry of the prostate: how accurate is it?, *Physiol. Meas.*, 16, 141-150, 1995.
16. E. D. Trautman and R. S. Newbower. The development of indicator-dilution techniques, *IEEE Trans. Biomed. Eng.*, 31, 800-807, 1984.
17. A. Rubini, D. Del Monte, V. Catena, I. Ittar, M. Cesaro, D. Soranzo, G. Rattazzi, and G. L. Alatti, Cardiac output measurement by the thermodilution method: an *in vitro* test of accuracy of three commercially available automatic cardiac output computers, *Intensive Care Med.*, 21, 154-158, 1995.
18. S. Jang, R. J. Jaszczak, F. Li, J. F. Debatin, S. N. Nadel, A. J. Evans, K. L. Greer, and R. E. Coleman, Cardiac ejection fraction and volume measurements using dynamic cardiac phantoms and radio-nuclide imaging, *IEEE Trans. Nucl. Sci.*, 41, 2845-2849, 1994.
19. D. M. K. S. Kaulesar Sukul, P. T. den Hoed, E. J. Johannes, R. van Dolder, and E. Benda, Direct and indirect methods for the quantification of leg volume: comparison between water displacement volumetry, the disk model method and the frustum sign model method, using the correlation coefficient and the limits of agreement, *J. Biomed. Eng.*, 15, 477-480, 1993.
20. P. Saucez, M. Remy, C. Renotte, and M. Mauroy, Thermal behavior of the constant volume body plethysmograph, *IEEE Trans. Biomed. Eng.*, 42, 269-277, 1995.
21. J. Rosell, K. P. Cohen, and J. G. Webster, Reduction of motion artifacts using a two-frequency impedance plethysmograph and adaptive filtering, *IEEE Trans Biomed. Eng.*, 42, 1044-1048, 1995.
22. K. P. Cohen, D. Panescu, J. H. Booske, J. G. Webster, and W. L. Tompkins, Design of an inductive plethysmograph for ventilation measurement, *Physiol. Meas.*, 15, 217-229, 1994.
23. O. S. Valerio Jimenez, J. K. Moon, C. L. Jensen, F. A. Vohra, and H. P. Sheng, Pre-term infant volume measurements by acoustic plethysmography, *J. Biomed. Eng.*, 15, 91-98, 1993.
24. R. G. Aarnink, R. J. B. Giesen, A. L. Huynen, J. J. M. C. H. de la Rosette, F. M. J. Debruyne, and H. Wijkstra, A practical clinical method for contour determination in ultrasonographic prostate images, *Ultrasound Med. Biol.*, 20, 705-717, 1994.
25. C. H. Chu, E. J. Delp, and A. J. Buda, Detecting left ventricular endocardial and epicardial boundaries by digital two-dimensional echocardiography, *IEEE Trans. Med. Im.*, 7, 81-90 1988.
26. J. Feng, W. C. Lin, and C. T. Chen, Epicardial boundary detection using fuzzy reasoning, *IEEE Trans. Med. Im.*, 10, 187-199, 1991.
27. S. Lobregt and M. A. Viergever, A discrete contour model, *IEEE Trans. Med. Im.*, 14, 12-24, 1995.

Further Information

- Anonymous, AARC Clinical Practice Guideline; Static lung volume, *Respir. Care*, 39, 830-836, 1994.
- Anonymous, AARC Clinical Practice Guideline; Body plethysmography, *Respir. Care*, 39, 1184-1190, 1994.
- E. F. Bernstein (ed.), *Noninvasive Diagnostic Techniques in Vascular Disease*, 3rd ed., St. Louis: Mosby, 1985.
- P. J. Davis and P. Rabinowitz, *Methods of Numerical Integration*, London: Academic Press, 1975.
- H. Feigenbaum, *Echocardiography*, 5th ed., Philadelphia: Lee & Febiger, 1993.
- W. N. McDicken, *Diagnostic Ultrasonics: Principles and Use of Instruments*, 3rd ed., London: Crosby Lockwood Staples, 1991.
- J. Nyboer, *Electrical Impedance Plethysmography*, Springfield, IL: Charles Thomas, 1959.
- S. Webb, *The Physics of Medical Imaging*, Bristol: IOP Publishing, 1988.
- J. B. West, *Respiratory Physiology — The Essentials*, Baltimore, MD: Williams and Wilkins, 1987.

TABLE 14.1 Defining Terms Relating to Angles

Term	Definition
Angle	A figure formed by two lines or planes that intersect one another.
Acute angle	An angle less than 90° .
Azimuth	The horizontal angle measured along the Earth's horizon, between a fixed reference (usually due south) and an object.
Bank	A lateral inclination.
Circle	A closed plane curve where all of its points are the same distance from its center point.
Declination = declivity	A negative slope.
Degree	Equal to $1/360$ of a circle.
Goniometer	An instrument for measuring angles (from the Greek word <i>gonio</i>).
Incline = Slope = Bias = Slant = Gradient = Grade	The deviation, plus or minus, from horizontal as defined by gravity.
Latitude	An angle measured north or south from the equator on a meridian to a point on the earth.
Lean = List = Tilt	The deviation from vertical as defined by gravity.
Longitude	The angle between the prime meridian (through Greenwich, England) and the meridian of a given place on Earth. This angle is defined as positive moving west.
Milliradian	An angle equal to $1/1000$ rad.
Minute	An angle equal to $1/60^\circ$.
Oblique angle	An obtuse or acute angle.
Obtuse angle	An angle greater than 90° .
Quadrant	One quarter of a circle (90°).
Radian	The angle subtended by an arc of a circle equal to the radius of that circle. One radian is equal to 57.29578° .
Rake	Equals the deviation in degrees from being perpendicular (90°) to a line or a plane.
Right angle	An angle of 90° .
Rise	A positive incline.
Second	An angle equal to $1/60'$ ($1/3600^\circ$).
Straight	An angle equal to 180° .
Taper	The change in diameter or thickness per unit length of axis.
Twist	The angle of turn per unit length of axis, as in a gun barrel or a screw thread.

14.2 Clinometers

A *clinometer* is an electronic device that measures vertical angle with respect to gravitational level. It is rectangular, with each side being a 90° to its adjacent sides. With a range of readings of at least $\pm 45^\circ$, this shape allows measurements up to a full 360° . Floating zero can be set anywhere and resolutions of $\pm 0.01^\circ$ are obtainable. Some models will convert readings to inches per foot, % of grade, and millimeters per meter (mm m^{-1}).

A clinometer can be used anywhere the angle of a surface with respect to gravity or another surface needs to be measured. High accuracy and resolution are obtainable, but calibration should be checked periodically with respect to a known level surface and a known angle. Surfaces that are remote to one another or have an intervening obstruction pose no problem for a clinometer.

14.3 Optical Comparator

An *optical comparator* measures angles, along with other dimensions or profiles, by referencing a visual image (usually magnified) of an object to a reticule that is calibrated in the measurement units desired. A hand-held optical comparator is placed directly over the object to be measured and the operator's eye is moved to the proper distance above the comparator for good focus of the magnified image. Some models contain a battery- or line-powered light source. Reticules for these hand-held devices are generally graduated in 1° increments.

TABLE 14.2 A Partial List of Manufacturers and Supplies of Angle Measurement Equipment

Company	Address
Flexbar Machine Corporation (Representative for Erich Preissr & Co., West Germany)	250 Gibbs Road Islandia, NY 11722-2697 Tel: (800) 879-7575
Fred V. Fowler Co., Inc.	66 Rowe Street P.O. Box 299 Newton, MA 02166 Tel: (617) 332-7004
L. S. Starrett Company	121 Crescent Street Athol, MA 01331 Tel: See local distributor
Brown & Sharpe Mfg. Co.	931 Oakton Street Elk Grove Village, IL 60007 Tel: (312) 593-5950
Swiss Precision Instruments, Inc. SPI	2206 Lively Blvd. Elk Grove Village, IL 60007 Tel: (708) 981-1300
Edmund Scientific Co.	101 East Gloucester Pike Barrington, NJ 08007-1380 Tel: (609) 573-6250

Projection-type optical comparators are available as bench or floor models and are made for either horizontal or vertical beam viewing. They use a high-intensity light source and magnifying optics to display an image of an object onto a rear-projection, frosted glass screen that is inscribed with angular as well as linear markings. The image displayed is the result of light being projected past the object, referred to as a shadow graph, or of light being reflected off the surface of the object. The method used is determined by the shape of the viewed object and its surface quality.

Magnification of the optical system in these devices can range from 10 \times by 100 \times , with screen diameters ranging from 0.3 m to 1 m.

These instruments are useful for measuring profiles of parts after final matching for the purpose of quality control or duplication.

As the name implies, the image that is projected can be superimposed on a mask or outline drawing placed directly on the view screen so that any deviations from the required shape can easily be determined. Optical comparators are heavy, nonportable devices that require a fairly high amount of maintenance and are best used in a fairly dark room.

14.4 Protractor

A *protractor* is an instrument used for measuring and constructing angles. A direct-reading protractor usually is graduated in 1 $^\circ$ increments and can be semicircular or circular in shape. The simplest models are of one-piece construction and made from metal or plastic. Other models include a blade or pointer pivoted in the center of the graduated circle.

More precise protractors are equipped with a vernier scale that allows an angle to be indicated to 5' of arc. See Figure 14.1 for an explanation of how to read such a vernier scale.

TABLE 14.3 Instruments and Devices Used to Measure or Indicate Angles

Type	Manufacturer	Model	Description	Approx. price
Sine bar	Flexbar	16292	5 in. × 15/16 in. wide	\$130.00
		16293	10 in. × 1 in. wide	
		16294	5 in. × 2 in. wide	
		12202	5 in. × 1 in. wide, economy	
	Fowler	52-455-010	5 in. center to center, 15/16 in. wide	\$30.00
		52-455-015	10 in. C. to C., 1 in. wide	
		52-455-030	2.5 in. C. to C., 1 in. wide	
	SPI	30-712-4	10 in. C. to C., universal bench center	\$3048.00
		98-379-1	5 in. C. to C., 1 in. wide, accuracy between rolls = 0.0003 in.	\$31.00
		30-091-3	10 in. C. to C., 1 in. wide, accuracy between rolls = 0.0001 in.	\$203.00
	Brown & Sharpe	598-291-121-1	5 in. C. to C., 1 in. wide	
		598-293-121-1	10 in. C. to C., 1 1/8 in. wide	
Sine plate	Flexbar	14612	5 in. C. to C., 6 in. × 3 in. × 2 in.	\$320.00
		14615	10 in. C. to C., 12 in. × 6 in. × 2 5/8 in.	\$1000.00
	Fowler	57-374-001	5 in. C. to C., 6 in. × 3 in. × 2 in.	
		57-374-004	10 in. C. to C., 12 in. × 6 in. × 2 5/8 in.	
Compound sine plate	SPI	77-026-3	10 in. C. to C., 12 in. × 6 in. × 2 5/8 in.	\$872.00
	Brown & Sharpe	599-925-10	10 in. C. to C., 12 in. × 6 in. × 2 3/8 in.	
	Flexbar	14616	5 in. C. to C., 6 in. × 6 in. × 3 1/8 in.	\$1100.00
	Fowler	57-375-001	5 in. C. to C., 6 in. × 6 in. × 3 1/8 in.	
Angle Computer Protractor-Direct	SPI	7-072-7	5 in. C. to C., 6 in. × 6 in. × 3 1/8 in.	\$926.00
	Brown & Sharpe	599-926-5	5 in. C. to C., 6 in. × 6 in. × 3 1/2 in.	
	Flexbar	19860	3-axis with vernier protractors	\$3750.00
	Flexbar	16337	Rectangular Head, 0–180°	\$25.00
Protractor-Vernier	Starrett	RP1224W	Head only, To fit 12 in., 18 in. & 24 in. blades	
		C183	Rectangular head, 0–180° 6 in. Blade	
	SPI	30-393-3	Rectangular head, 0–180°	\$23.00
		31-804-8	Head only. To fit 12 in., 18 in. & 24 in. Blades	\$39.00
	Flexbar	16339	360° range, 1' reading with magnifier, 12 in. & 6 in. blades incl.	\$400.00
		16338	360° range, 5' reading	\$75.00
	Starrett	C364DZ	12 in. Blade, 0–90° range thru 360°, 5' graduations	
		SPI	30-395-8	6 in. Blade, 0–90° range thru 360°, 5' graduations
Protractor, Digital (Inclinometer)	Flexbar	30-390-9	6 in. & 12 in. Blades, 0–90° range thru 360°, 1' graduations with magnifier	\$540.00
		599-490-8	8 in. Blade, 0–90° range thru 360°, Magnifier optional	
	Fowler	17557	±45° range, ±0.1° resolution	\$260.00
		17556	±60°, ±0.01° resolution, SPC output	\$450.00
Protractor, Dial Bevel	SPI	54-635-600	±45° range, ±0.01° resolution, RS232 output available	
		31-040-9	±45° range, resolution: ±0.01° (0 to ±10°), 0.1° (10° to 90°)	\$329.00
		30-150-7	8 in. Blade, 1 3/8 in. diameter dial, geared to direct read to 5'	\$527.30
Square-Reference Optical Comparator (Projector)	Brown & Sharpe	599-4977-8	8 in. Blade, dial read degrees and 5'	
		SPI	30-392-5	90° fixed angle
	Fowler	53-912-000	12 in. screen diameter, 10×, 20×, 25× lens available, horizontal beam, with separate light source for surface illumination	
Starrett		HB350	14 in. screen diameter, 10×, 20×, 25×, 31.25×, 50×, 100× lens available, horizontal beam	

TABLE 14.3 (continued) Instruments and Devices Used to Measure or Indicate Angles

Type	Manufacturer	Model	Description	Approx. price
		VB300	12 in. screen diameter, 10× through 100× lens available, Vertical beam	
		HS1000	40 in. screen diameter, 10× thru 100× lens available, Horizontal beam	
	SPI	40-350-1	14 in. screen diameter, 10×, 20×, 50× lens available, Horizontal beam	\$2995.00
Optical Comparator (hand-held)	SPI	40-145-3	10× Magnification, Pocket style	\$57.50
			Additional Reticles	\$11.00
	Edmund Scientific	40-140-6	7× Magnification, pocket style with illuminator	\$62.50
			Additional Reticles	\$10.50
	A2046	6× Magnification, pocket style, 360° protractor reticle, 1° increments	\$58.75	
Angle Plate	Fowler	52-456-000	Set of 2, 9/32 in. thick, steel, 30 × 60 × 90°, 45 × 45 × 90°	
	SPI	98-463	Set of 2, 5/16 in. thick, steel, 30 × 60 × 90°, 45 × 45 × 90°	\$32.00
Angle Positioning Block	SPI	70-997-2	0 to 60°, 10' vernier (for setting workpiece in a vice)	\$122.00
Angle Gage	Fowler	52-470-180	18 leaves, spring steel, 1 thru 10, 14, 14.5, 15, 20, 25, 30, 35, and 45°	
Angle Gage Blocks	SPI	31-375-9	18 gage set, 5° thru 90° in 5° steps, 5' accuracy	\$49.60
	Starrett	Ag18.TR	18 block set, use in combination for steps of 1", 1" accuracy	
		Starrett AG16.LM	16 block set, use in combination for steps of 1", 1/4" accuracy	
	SPI	30-140-8	10 block set, 1, 2, 3, 4, 5, 10, 15, 20, 25, 30°, Accuracy = ±0.0001" per inch. 1/4° and 1/2° blocks optional (each)	\$170.00 \$18.00

14.5 Sine Bar

A *sine bar* is a device used to accurately measure angles or to position work pieces prior to grinding and other machining procedures. It is constructed from a precisely hardened and ground steel rectangular bar to which are attached two hardened and ground steel cylindrical rods of the same diameter. The axis of each rod is very accurately positioned parallel to the other and to the top surface of the bar.

A sine bar is used in conjunction with precision gage blocks that are placed under one of the two cylindrical rods to raise that rod above the other rod a distance H (see Figure 14.2) equal to the sine of the angle desired, times the distance D between the two rods. The standard distance between the rods is 250 mm (5 in.) or 500 mm (10 in.). The governing equation in using a sine bar is $\sin A = H/D$.

A work piece positioned using a sine bar is usually secured with the use of a precision vice. The vice may clamp directly to the work piece or, when using a sine bar that has tapped holes on its top surface, to the sine bar sides with the work piece bolted to the top of the sine bar.

14.6 Sine Plate

A variation of the sine bar is the *sine plate*. A sine plate consists of the three elements of a sine bar plus a bottom plate and side straps used to lock the plate in the desired position. In addition, one of the ground steel rods is arranged to form a hinge between the top and bottom plates. When using a sine plate, a work piece is secured to the top plate using bolts or clamps and the bottom plate is secured to a machine tool table using clamps or a magnetic chuck.

Compound sine plates are bidirectional, allowing angles to be measured and set in each of two orthogonal planes (true compound angles).

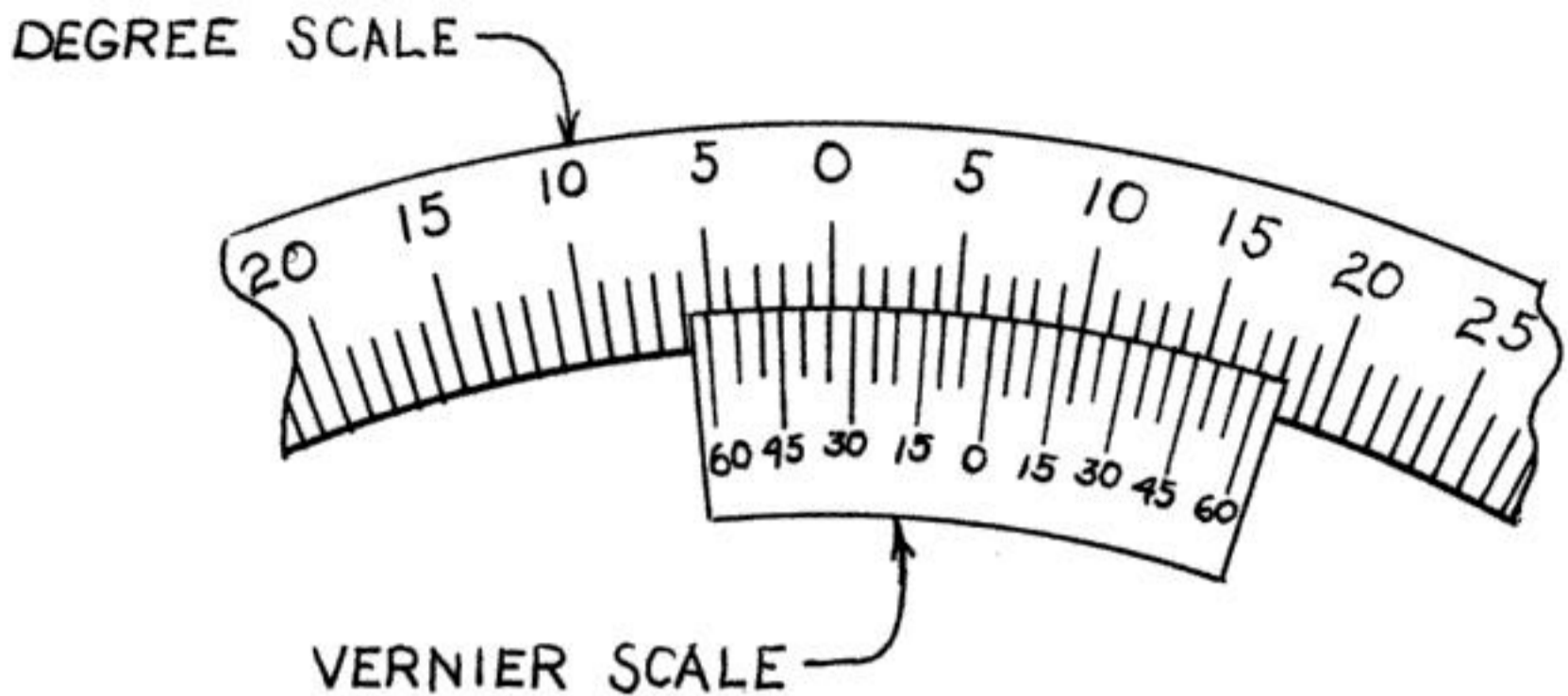


FIGURE 14.1 Vernier protractor. If the zero mark on the vernier scale is to the right of the zero mark on the main scale, as shown in this drawing, then the right side of the vernier scale must be used. Look for the mark on the vernier scale that best aligns with one of the marks of the protractor degree scale. Count the number of marks on the vernier scale from the zero mark to this mark. Each mark thus counted is, in this example, equal to 5' of arc and, therefore, the number of minutes to be added to the number of degrees indicated is 5 times the vernier marks counted. (In this example, the fourth mark aligns the best with the main scale indicating 20'). The number of degree indicated is the degree mark just to the left of the zero mark on the vernier scale. The left side of the vernier is similarly used when the indicated angle is to the left of the zero mark on the degree scale.

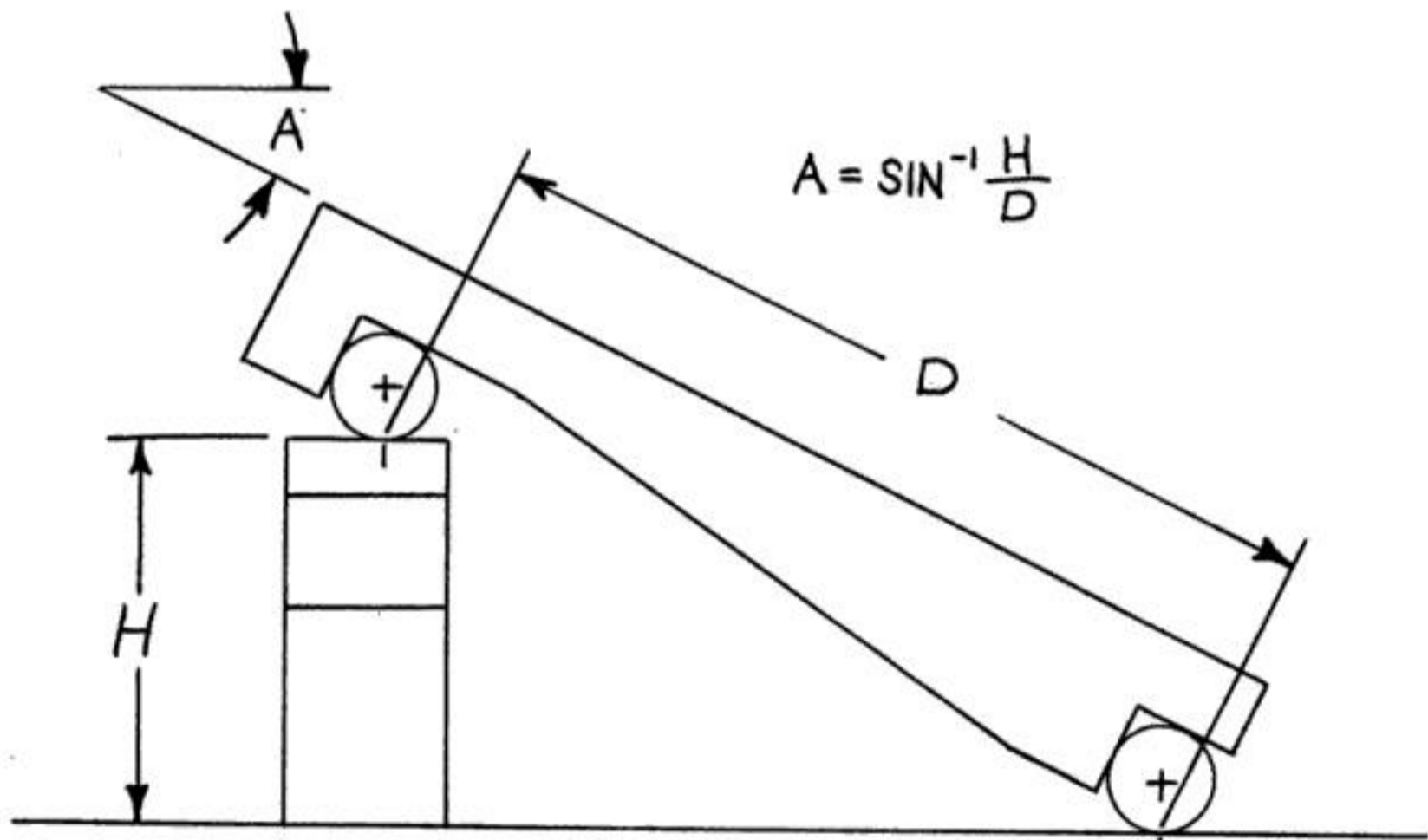


FIGURE 14.2 Sine bar. A sine bar is used in conjunction with precision gage blocks that are placed under one of the cylindrical rods, raising that end a distance, H , equal to the sine of the desired angle times the distance, D , between the two rods.

15

Tilt Measurement

Adam Chrzanowski
University of New Brunswick

James M. Secord
University of New Brunswick

15.1	Tiltmeters or Inclinometers.....	15-1
15.2	Geodetic Leveling.....	15-4
15.3	Hydrostatic Leveling.....	15-5
15.4	Suspended and Inverted Plumb Lines.....	15-7
15.5	Integration of Observations.....	15-9

If the relative position of two points is described by the three-dimensional orientation of a line joining them, then, in general, *tilt* is the angular amount that the orientation has changed, in a vertical plane, from a previous direction or from a reference direction. If the original or reference orientation is nearly horizontal, then the term "tilt" is usually used. If it is nearly vertical, then the change is often regarded as "inclination." Here, "tilt" will refer to either. The two points can be separated by a discernable amount, the base, or the tilt can be measured at a point with the reference orientation being defined by the direction of the force of gravity at that point. Thus, the same instrument that measures tilt at a point can be called either a *tiltmeter* or an *inclinometer* (or *clinometer*), depending on the interpretation of the results. The instrument used to measure a series of tilts along any vertical orientation is often called an *inclinometer* (e.g., Dunicliff [1]).

Angular tilt is directly related to the linear amount of change subtending the length of the base. Consequently, angular tilt does not have to be measured directly but can be derived from the mechanical or other measurement of this linear change if the length of the base is known.

Therefore, the following discussion has been subdivided into:

1. Tiltmeters or inclinometers (angular tilt at a point or over a limited, relatively small base length)
2. Geodetic leveling (tilt derived from a height difference over an extended base of virtually limitless length)
3. Hydrostatic leveling (tilt derived from a height difference over an extended base of limited length)
4. Suspended and inverted pendula, or plumb lines (inclination from a difference in horizontal relative position over a vertical base or height difference)

15.1 Tiltmeters or Inclinometers

Considering the basic principle of operation, tiltmeters may be divided into: liquid (including spirit bubble type), vertical pendulum, or horizontal pendulum. Dunicliff [1] provides a comprehensive review of tiltmeters and inclinometers according to the method in which the tilt is sensed (i.e., mechanical, with accelerometer transducer, with vibrating wire transducer, or with electrolytic transducer).

The sensitivity of tilt falls into two distinct groups: geodetic or geophysical special tiltmeters with a resolution of 10^{-8} rad (0.002") or even 10^{-9} rad; and engineering tiltmeters with resolutions from 0.1" to several seconds of arc, depending on the required range of tilt to be measured.

The first group includes instruments that are used mainly for geophysical studies of Earth tide phenomena and tectonic movements; for example, the Verbaander-Melchior [2] and the Zöllner [3, 4] horizontal pendulum tiltmeters, and the Rockwell Model TM-1 [5] liquid-bubble type. This category of instrument requires extremely stable mounting and a controlled environment. There are very few engineering projects where such sensitive instruments are required. However, deformation measurements of underground excavations for the storage of nuclear waste may be one of the few possible applications. An example is a mercury tiltmeter (Model 300) developed for that purpose by the Auckland Nuclear Accessory Co. Ltd. in New Zealand. In this instrument, the change in capacitance between electrodes and a mercury surface is proportional to the tilt. This tiltmeter, with a total range of 15", is claimed to give a resolution of 10^{-9} rad (0.0002"), which corresponds to a relative vertical displacement of only 6×10^{-7} mm over its base length of 587 mm.

In the second group, there are many models of liquid or pendulum tiltmeters of reasonable price (\$2000 to \$5000) that satisfy most needs in engineering studies. Apart from a spirit level or level vial by itself, the simplest form of tiltmeter is a base that is tens of centimeters long and leveled by centering the bubble in the vial by direct viewing or by an optical coincidence viewing of the two ends of the bubble. Direct viewing gives a resolution of 1/5 of a vial division (usually 2 mm), which typically has a sensitivity of 10" to 30" per division. Coincidence viewing increases the setting accuracy to 0.03 of the sensitivity of the vial. The discrepancy from horizontal between the two measurement points can be determined by a dial gage or micrometer that has a resolution of 0.0005 in. or 0.02 mm. Huggenberger AG claim a sensitivity of 0.3" (1×10^{-4} gon) over a range of $\pm 21'$ for their clinometer with a 100 mm base and coincidence centering of the bubble in the level vial. The clinometer can be attached to 1 m bases for either horizontal or vertical measurements.

If the vial is filled with an electrolytic liquid, the centering of the bubble can be done electrically. An example is the Electrolevel (by the British Aircraft Corp.), which uses the spirit bubble principle [6] and in which the movement of the bubble is sensed by three electrodes. A tilt produces a change in differential resistivity between the electrodes that is measured by means of a Wheatstone bridge. A resolution of 0.25" is obtained over a total range of a few minutes of arc. Many other liquid types of tiltmeters with various ranges (up to 30°) are available from various companies. Holzhausen [7] and Egan and Holzhausen [8] discuss the application of electrolytic tiltmeters (resolution of 2" over a range of $\pm 1^\circ$, manufactured by Applied Geomechanics) in the monitoring of various hydroelectric power dams in the U.S.

The Rank Organization in the United Kingdom [9] makes a liquid-dampened pendulum-type electronic level, the Talyvel, which gives an accuracy of $\pm 0.5''$ over a total range of $\pm 8'$. A similar transducer of the pendulum type is used in the Niveltronic tiltmeter (range of $\pm 150''$ with an accuracy of $\pm 0.2''$) produced by Tesa S.A. in Switzerland. Of particular popularity are servo-accelerometer tiltmeters with horizontal pendula. They offer ruggedness, durability, and can operate in low temperatures. The output voltage is proportional to the sine of the angle of tilt. Schaevitz Engineering produces such a servo-accelerometer that employs a small-mass horizontal paddle (pendulum) which tries to move in the direction of tilt, due to the force of gravity. Any resultant motion is converted by position sensors to a signal input to the electronic amplifier whose current output is applied to the torque motor. This develops a torque that is equal and opposite to the original. The torque motor current produces a voltage output that is proportional to the sine of the angle of tilt.

The typical output voltage range for tiltmeters is ± 5 V, which corresponds to the maximum range of tilt and readily allows for serial interfacing. The angular resolution of a tiltmeter depends on its range of tilt since a larger range would result in more angular tilt per unit voltage so a higher resolution tiltmeter would have a smaller range of measurable tilt. Typically, the resolution is 0.02% of the range (volts) [10].

There are many factors affecting the accuracy of tilt sensing, not just the resolution of the readout. A temperature change produces dimensional changes in the mechanical components and changes in the viscosity of the liquid in electrolytic tiltmeters and of the dampening oil in pendulum-type tiltmeters. Also, electric characteristics can alter with temperature changes. Drifts in tilt indications and fluctuations of the readout may also occur. Compensation for the effects of temperature changes can be incorporated

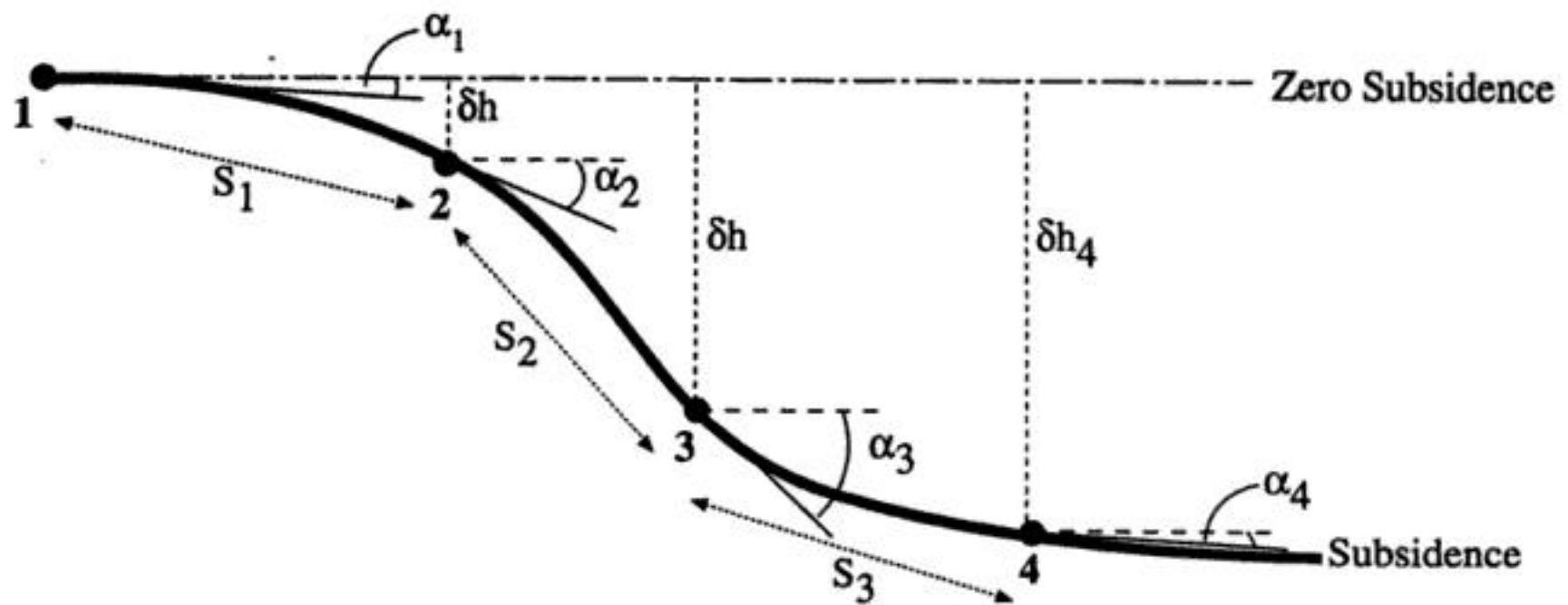


FIGURE 15.1 Ground subsidence derived from tilt measurements.

in the construction of an instrument, but at an increased cost. An alternative is to design a linear reaction by the instrument to the effects of temperature and to apply simple calibration corrections.

In less expensive models, compensation for the aforementioned sources of error is not very sophisticated, and such tiltmeters may show nonlinear output in reaction to changes in temperature and erratic drifts or other behavior that would be difficult to predict without testing. Consequently, very thorough testing and calibration are required even when the accuracy requirement is not very high [11]. Testing should investigate, at least, the linearity of output in reaction to induced tilts over the instrument's full range (\pm) and to changes in temperature. Some suggestions for testing and calibrating inclinometers, among other instruments, are given in Dunnycliff [1]. It is further emphasized that regular and up-to-date calibration is important in order to ensure continuity in the fidelity of the data being gathered. In most cases, the behavior being investigated changes with time and incorrect data cannot be recaptured.

Compensators for vertical circle readings in precision theodolites work on the same principle as some engineering tiltmeters. The liquid compensator of the Kern E2 electronic theodolite [12] gave a repeatability of better than 0.3" over a range of $\pm 150''$ and was incorporated in their tiltmeter, NIVEL 20, in 1989. The same compensation system has been used in the currently available Leica TC2002 precision electronic total station [13]. Consequently, the theodolite may also be used as a tiltmeter, in some applications, giving the same accuracy as the Electrolevel, for example.

Tiltmeters have a wide range of applications. A series of tiltmeters, if arranged along a terrain profile in a mining area, may replace geodetic leveling in the determination of ground subsidence [11] as shown in Figure 15.1. For example, the subsidence (i.e., the variation from the previous or original profile) of point 4 (δh_4) with respect to point 1 may be estimated from the observed changes in tilt, from a base or original position, (α_i in radians) and known distances between the points as:

$$\delta h_4 = s_1(\alpha_1 + \alpha_2)/2 + s_2(\alpha_2 + \alpha_3)/2 + s_3(\alpha_3 + \alpha_4)/2 \quad (15.1)$$

The fidelity of this method depends on the density of tilt measurements along the profile and the continuity of the profile (a constant slope of the terrain between measurement points is assumed). Similarly, deformation profiles of tall buildings may be determined by placing a series of tiltmeters at different levels of the structure [14]. Also, changes in borehole profiles can be created in a similar manner. The absolute profile of a borehole can be generated by considering the horizontal displacement in the direction of the orientation of the inclinometer (usually controlled by guide grooves in the borehole casing) for the i -th position, as it traverses a borehole with observation of α_i at a depth s_i . However, this would require calibration of the inclinometer to correct its output to show zero in its vertical position since the α_i are tilts from the vertical rather than angular changes from an original inclination.

In geomechanical engineering, the most popular application of tiltmeters and borehole inclinometers is in slope stability studies and in monitoring earth-fill dams. Torpedo-shaped biaxial inclinometers are used to scan boreholes drilled to the depth of an expected stable strata in the slope. By lowering the inclinometer on a cable with marked intervals and taking readings of the inclinometer at those intervals, a full profile of the borehole and its changes may be determined through repeated surveys, as mentioned above. SINCO and other producers of tiltmeters provide special borehole inclinometers (50 cm or 2 ft long) with guide wheels to control the orientation of the inclinometer. A special borehole casing (plastic or aluminum) with guiding grooves for the wheels is available. Usually, servo-accelerometer type inclinometers are used with various ranges of inclination measurements; for example, $\pm 6^\circ$, $\pm 53^\circ$, or even $\pm 90^\circ$. A 40 m deep borehole, if measured every 50 cm with an inclinometer having an accuracy of only $\pm 100''$, should allow for the determination of linear lateral displacement of the collar of the borehole with an accuracy of ± 2 mm.

In cases where there is difficult access to the monitored area or a need for continuous data acquisition or both, tiltmeters or borehole inclinometers can be left in place at the observing station with a telemetry monitoring system allowing for communication to the processing location. One example of a station setup of a telemetric monitoring of ground subsidence in a mining area near Sparwood, B.C. used a telemetry system developed for the Canadian Centre for Mining and Energy Technology (CANMET) by the University of New Brunswick [11, 15]. Terra Technology biaxial servo-accelerometer tiltmeters of $\pm 1^\circ$ range were used in the study. The telemetry system could work with up to 256 field stations. Each station accepted up to six sensors (not only tiltmeters but any type of instrument with electric output, e.g., temperature, power level or voltage). Another example is a fully automated borehole scanning system with a SINCO inclinometer and telemetric data acquisition that was also developed at the University of New Brunswick [16]. It has been used successfully in monitoring highwall stability at the Syncrude Canada Limited tarsands mining operation in northern Alberta.

15.2 Geodetic Leveling

Geodetic or differential leveling measures the height difference between two points using precision tilting levels, or precision automatic levels with compensators, with parallel plate micrometers and calibrated invar leveling staves or rods. Recent technology has provided digital automatic levels that use a CCD sensor in the instrument and bar codes, rather than linear graduations, on the staves [17]. Their ease of use and comparable precision have quickly made them rivals to the traditional optical instruments. In a single setup of the level, the height difference is the backsight rod reading minus the foresight rod reading. Any number of setup height differences can be combined to determine the height difference between two points of interest; however, the errors involved accumulate with the number of setups. With sight lengths limited to no more than 20 m, geodetic leveling can produce height differences with a precision of ± 0.1 mm per setup, which is equivalent to a precision of $\pm 0.5''$ in tilt. Although geodetic leveling is traditionally used to determine elevations, it is often used to monitor not only the settlement of sensitive structures but also to describe the tilt of components of a structure by determining the tilt between appropriate pairs of benchmarks (monumented in or on the structure) [18]. Since the level reference is created by an optical line of sight through a telescope (magnification up to $40\times$), a major source of systematic error is the effect of vertical atmospheric refraction. A vertical temperature gradient of even 1°C m^{-1} across the line of sight would bend the line of sight to be in error by 0.4 mm at 30 m. Gradients of this magnitude are commonly encountered in industrial settings and are usually even more evident outdoors. Less effect is realized if the sight lengths are kept short, but this must be weighted against the accumulation of random error with each additional setup (shorter sight lengths would require more setups to traverse the same height difference). The errors that seem to be insignificant in a single setup (or in a few setups) become magnified in height differences involving a large number of setups (e.g., rod scale error and settlement of the instrument or rods). Such errors have become quite noticeable in the monitoring of tectonic plate movement and undetected systematic effects can be misleading. Further discussion on precision leveling and sources of error is available in Vanicek et al. [19] and in Schomacker and Berry [20].

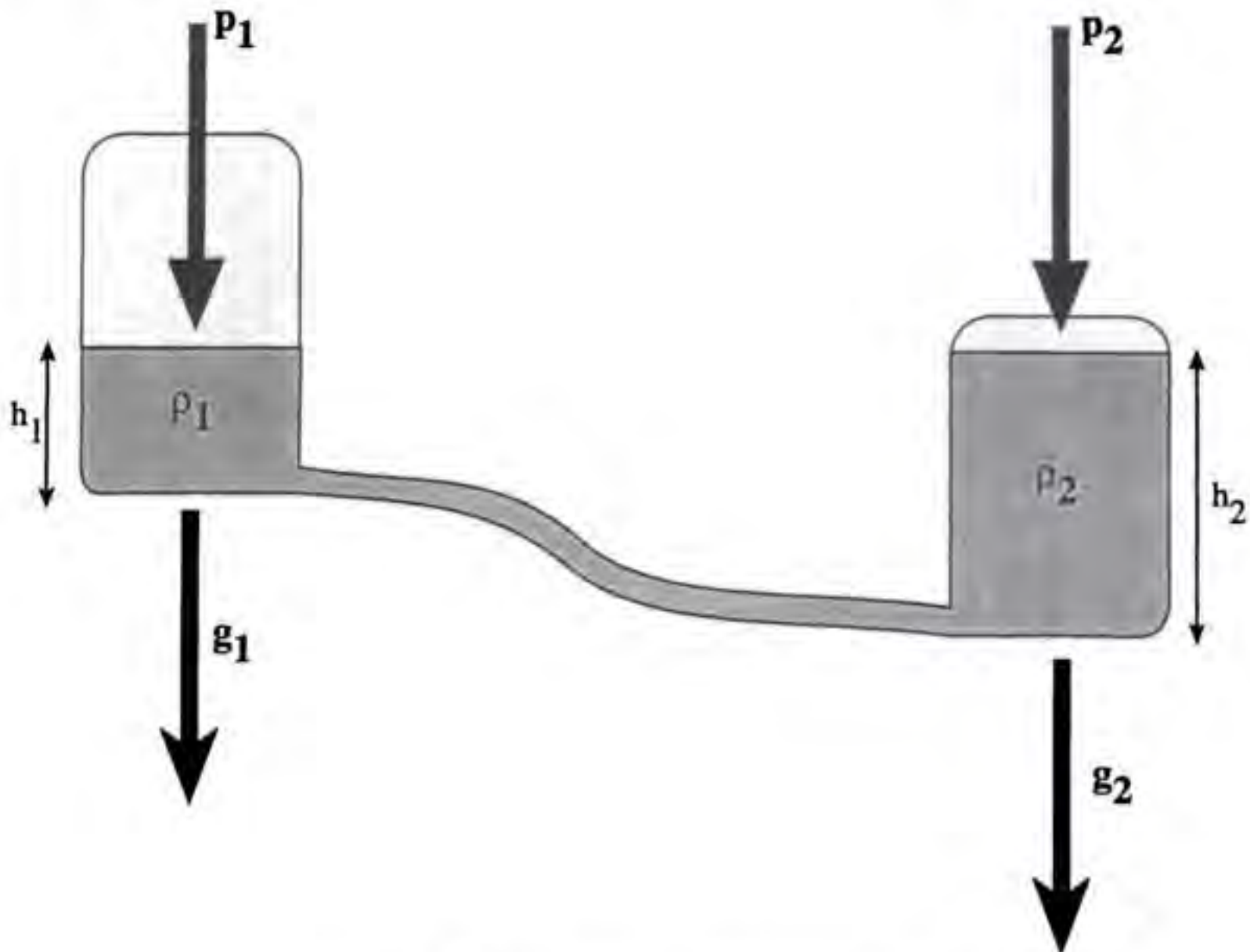


FIGURE 15.2 Hydrostatic equilibrium in connected vessels.

Having the elevations or heights, h_1 and h_2 , of two points or having measured, or determined, the height difference between them, $\Delta h_{12}^t = h_2 - h_1$, at a time t , means that, if $\delta\Delta h = \Delta h_{12}^{t2} - \Delta h_{12}^{t1}$, the tilt, T_{12} , can be calculated if the horizontal separation s_{12} is known, since $T_{12} = \delta\Delta h/s_{12}$. The separation does not have to be known as precisely as the height difference since the total random error is $\sigma_{T^2} = \sigma_{\delta\Delta h^2}/s^2 + \sigma_s(\delta\Delta h^2/s^4)$. As an example, for two points that are 60 m apart with a height difference of 0.5 m (extreme in most structural cases) with the height difference known to $\pm 50 \mu\text{m}$ ($\sigma_{\delta\Delta h}$) and the distance known to ± 0.01 m (σ_s), the tilt would have a precision (σ_T) of $\pm 0.3''$. Further, neither the measurement of the height difference nor the determination of the separation have to be done directly between the two points. The leveling can be done along whatever route is convenient and the separation can be obtained in a variety of ways, for example, inverting from coordinated values for the points [21].

15.3 Hydrostatic Leveling

If two connected containers (Figure 15.2) are partially filled with a liquid, then the heights h_1 and h_2 are related through the hydrostatic equation (Bernoulli's equation, as given in [22]):

$$h_1 + P_1/(g_1\rho_1) = h_2 + P_2/(g_2\rho_2) = c \quad (15.2)$$

where P is the barometric pressure, g is the force of gravity, ρ is the density of the liquid which is a function of temperature, and c is a constant.

The above relationship has been employed in hydrostatic leveling, as shown schematically in Figure 15.3. The air tube connecting the two containers eliminates possible error due to different air pressures at two stations. The temperature of the liquid should also be maintained constant because, for example, a difference of 1.2°C between two containers may cause an error of 0.05 mm in a Δh determination

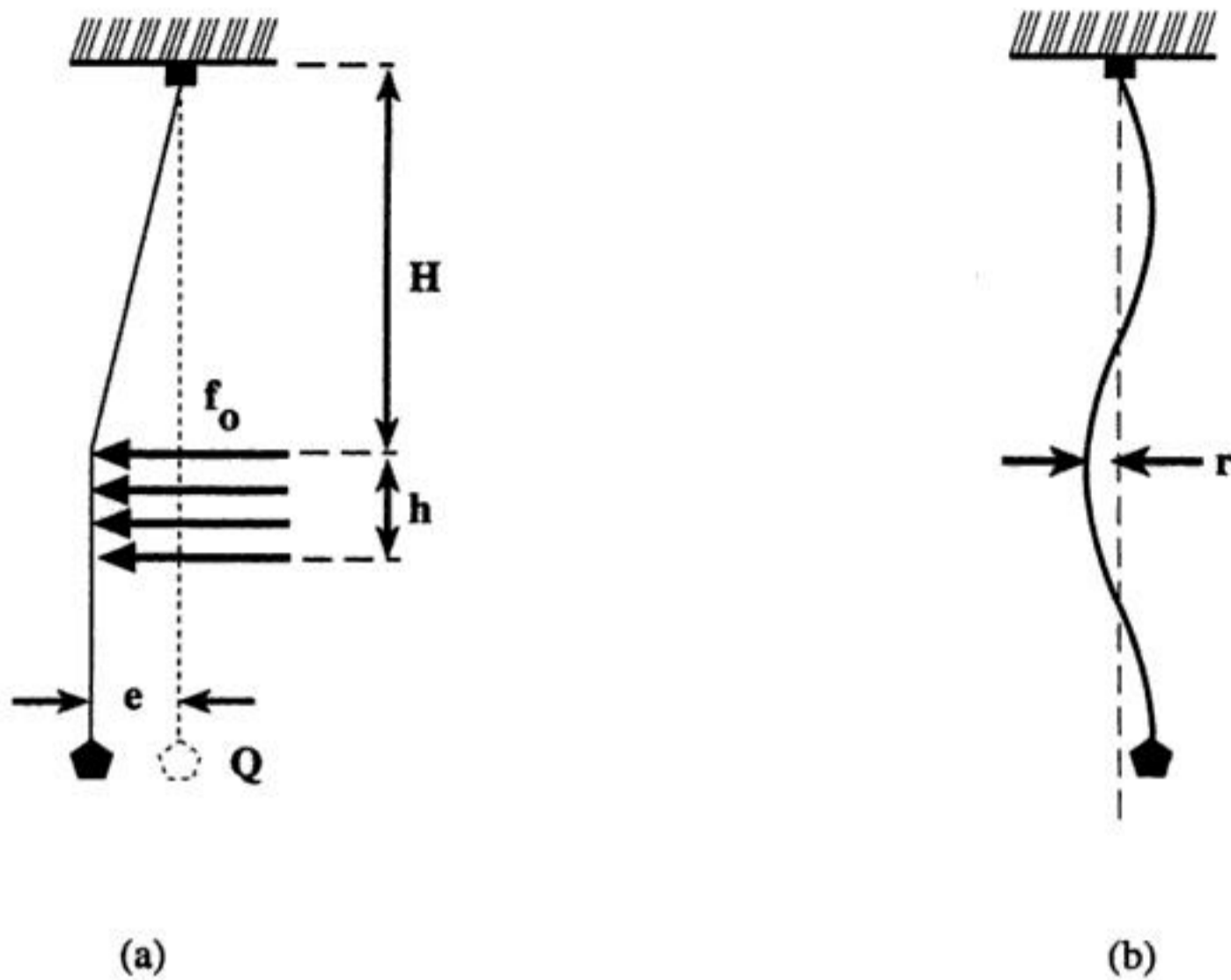


FIGURE 15.5 (a) Influence of air currents on a suspended plumbline. (b) Horizontal error due to the spiral shape of the wire.

AG in Switzerland. Telemac Co. (France) developed a system, Telependulum (marketed by Roctest), for continuous sensing of the position of the wire with remote reading and recording. A rigidly mounted reading table supports two pairs of induction type proximity sensors arranged on two mutually perpendicular axes. A hollow cylinder is fixed on the pendulum wire at the appropriate level, passing through the center of the table and between the sensors. Changes in the width of the gap between the target cylinder and the sensors are detected by the corresponding changes in the induction effect. The system has a resolution of ± 0.01 mm.

An interesting Automated Vision System has been developed by Spectron Engineering (Denver, Colorado). The system uses solid state electronic cameras to image the plumb line with a resolution of about $3 \mu\text{m}$ over a range of about 75 mm. Several plumb lines at Glen Canyon dam and at Monticello dam, near Sacramento, California, have been using the system since 1982 [29].

Two sources of error, which may be often underestimated by the user, may strongly affect plumb line measurements:

1. The influence of air currents
2. The spiral shape of the wire

If the wire of a plumb line (Figure 15.5(a)), with pendulum mass Q , is exposed along a length h to an air current of speed v at a distance H from the anchor point, then the plumb line is deflected by an amount [30]:

$$e = f_0 h H / Q \quad (15.3)$$

where f_0 is the acting force of air current per unit length of the wire. The value of f_0 may be calculated approximately from [30]

$$f_0 = 0.08 d v^2 / Q \quad (15.4)$$

where d is the diameter of the wire in millimeters, v is in meters per second, and Q is in kilograms. As an example, if $H = 50$ m, $h = 5$ m, $d = 1$ mm, $Q = 20$ kg, and $v = 1$ m s⁻¹ (only 3.6 km h⁻¹) then $e = 1$ mm.

The second source of error, which is usually underestimated in practice, is that the spiral shape (annealing) of the wire (Figure 15.5(b)) affects all wires unless they are specially straightened or suspended for a prolonged time (on the order of several months). If the wire changes its position (rotates) between two campaigns of measurements, then the recorded displacements could have a maximum error of $2r$. The value of r can be calculated from [30]:

$$r = \left(\pi d^4 E \right) / \left(64 R Q \right) \quad (15.5)$$

where E is Young's modulus of elasticity (about 2×10^{11} Pa for steel); R is the radius of the free spirals of the unloaded wire that, typically, is about 15 cm for wires up to 1.5 mm diameter; and d and Q are as above. For a plumb wire with $d = 1$ mm, $R = 15$ cm, and $Q = 196$ N (i.e., 20 kg), $r = 0.3$ mm.

If one plumb line cannot be established through all levels of a monitored structure, then a combination of suspended and inverted plumb lines may be used as long as they overlap at least at one level of the structure. At Hydro Quebec, the drill holes of the plumb lines are also used for monitoring height changes (vertical extension) by installing tensioned invar wires [31].

15.5 Integration of Observations

The above discussion has considered using individual instruments. Because many investigations, using other instrumentation as well as the measurement of tilt, involve the repetition of measurements, often over a long term, the fidelity of the measurements and their being referenced to the original or initial observation is vital to the investigation. It is risky to expect consistent behavior of instrumentation, particularly in environments with dramatic variations (especially in temperature), and over a long period of time. Any conclusion relating to the behavior of a structure is only as good as the data used in the analysis of the behavior. Two ways to ensure reliability are possible. One is to make regular testing and calibration a component of the observation regimen. The other is to analyze the observations as they are accumulated, either each observable as a temporal series of repeated measurements or observations of different locations or types together in an integrated analysis. The analytical tools for integrated analyses, as well as for calibration testing and temporal series analysis, have been developed [21, 32] and successfully implemented in several projects [15, 18, 33]. Proper calibration testing and correction, rigorous statistical analysis of trend in temporal series, and integrated analysis have proven to be valuable tools in the analysis of deformations and serve to enhance the monitoring of sensitive structures.

References

1. J. Dunnycliff, *Geotechnical Instrumentation for Monitoring Field Performance*, New York: John Wiley & Sons, 1988.
2. P. Melchior, *The Tides of the Planet Earth*, 2nd ed., Oxford, U.K.: Pergamon Press, 1983.
3. W. Torge, *Geodesy*, New York: Walter de Gruyter, 1991.
4. M. van Ruymbeke, Sur un pendule horizontal équipé d'un capteur de déplacement a capacité variable, *Bull. Géodésique*, 50, 281-290, 1976.
5. G. L. Cooper and W. T. Schmars, Selected applications of a biaxial tiltmeter in the ground motion environment, *J. Spacecraft Rockets*, 11, 530-535, 1974.
6. M. A. R. Cooper, *Modern Theodolites and Levels*, London: Crosby Lockwood, 1971.
7. G. R. Holzhausen, Low-cost, automated detection of precursors to dam failure: Coolidge Dam, Arizona, *Association of State Dam Safety Officials, 8th Annu. Conf.*, San Diego, CA, 1991.
8. N. H. Egan and G. R. Holzhausen, Evaluating structures using tilt (rotation) measurements, *Sensors Expo West Proc.*, 1991.

9. W. Caspary and A. Geiger, Untersuchungen zur Leistungsfähigkeit elektronischer Neigungsmesser, *Schriftenreihe, Vermessungswesen HSBW*, 3, München, Germany, 1978.
10. A. Chrzanowski, Geotechnical and other non-geodetic methods in deformation measurements, in Y. Bock (ed.), *Proc. Deformation Measurements Workshop*, Massachusetts Institute of Technology, Boston, 1986, 112-153.
11. A. Chrzanowski, W. Faig, B. Kurz, and A. Makosinski, Development, installation and operation of a subsidence monitoring and telemetry system, Contract report submitted to CANMET, Calgary, Canada, November 1980.
12. Kern & Co. Ltd., E2 Instruction Manual, Kern & Co. Ltd., Aarau, Germany, 1984.
13. Leica AG, Wild TC2002 User Manual, Leica AG, Heerbrugg, Germany, 1993.
14. H. Kahmen, *Elektronische Messverfahren in der Geodäsie*, Karlsruhe: Herbert Wichmann, 1978.
15. A. Chrzanowski and M. Y. Fisekci, Application of tiltmeters with a remote data acquisition in monitoring mining subsidence, *Proc. 4th Canadian Symp. Mine Surveying*, Banff, Alberta, Canada, June 1982.
16. A. Chrzanowski, A. Makosinski, A. Zielinski, and W. Faig, Highwall Monitoring System, Contract report to Syncrude Canada Limited, April 1988.
17. H. Ingensand, Wild NA2002, NA3000, Technical paper digital levels, Leica AG, Heerbrugg, Germany, 1993.
18. A. Chrzanowski and J. M. Secord, The 1989 Integrated Analysis of Deformation Measurements at the Mactaquac Generating Station, Contract Report to N.B. Power, May 1990.
19. P. Vanicek, R. O. Castle, and E. I. Balazs, Geodetic leveling and its applications, *Rev. Geophys. Space Phys.*, 18(2), 505-524, 1980.
20. M. C. Schomaker and R. M. Berry, *Geodetic Levelling*, NOAA Manual NOS NGS 3, U.S. Department of Commerce, National Oceanic and Atmospheric Administration, National Ocean Survey, Rockville, MD, 1981.
21. J. M. Secord, Development of the automatic data management and the analysis of integrated deformation measurements, Ph.D. dissertation, Department of Geodesy and Geomatics Engineering Technical Report 176, University of New Brunswick, Fredericton, 1995.
22. K. Schnädelbach, Neuere Verfahren zur präzisen Längen-und Höhenmessung, *Allgemeine Vermessungs-Nachrichten*, 1/1980.
23. G. R. Huggett, L. E. Slater, and G. Pavlis, Precision leveling with a two-fluid tiltmeter, *Geophys. Res. Lett.*, 3(12), 754-756, 1976.
24. H. Thierbach and W. Barth, Eine neue automatische Präzisionsschlauchwaage, *Z. Vermessungswesen*, 100, 470-478, 1976.
25. F. Robotti and T. Rossini, Analysis of differential settlements on monumental structures by means of the DAG automatic measuring device of levels and inclinations, in *Land Subsidence*, IAHA Publication No. 151, 1984.
26. T. H. Hanna, *Field Instrumentation in Geotechnical Engineering*, Clausthal-Zellerfeld, Germany, Trans Tech Publications, 1985.
27. E. Meier, A differential pressure tiltmeter for large-scale ground monitoring, *Water Power & Dam Construction*, 43(1), 38-40, 1991.
28. L. Dubreuil and R. Hamelin, Le forage vertical des trous de pendules inversés, *2nd Canadian Symp. Mining Surveying Rock Deformation Measurements*, Kingston, Ontario, Canada, 1974.
29. G. Kanegis, Automated vision system installed at Glen Canyon Dam, Spectron Engineering, n.d. [circa 1983].
30. A. Chrzanowski, E. Derenyi, and P. Wilson, Underground survey measurements — Research for progress, *The Canadian Mining and Metallurgical Bulletin*, June 1967.
31. B. Boyer and R. Hamelin, Monitoring Survey: Recent Developments in the Use of Inverted Pendula, Report No. 4, Question 56, *XV Congress Int. Commission Large Dams, Lausanne*, 1985.
32. A. Chrzanowski, Y. Q. Chen, P. Romero, and J. Secord, Integration of geodetic and geotechnical deformation surveys in the geosciences, *Tectonophysics*, 130, 1986.
33. A. Chrzanowski, Y. Q. Chen, J. Secord, and A. Szostak-Chrzanowski, Problems and solutions in the integrated monitoring and analysis of dam deformations, *CISM J. ACSGC*, 45(4), 547-560, 1991.

Appendix

A Sampling of Possible Suppliers of Tilt Measuring Instrumentation

Applied Geomechanics Inc. 1336 Brommer Street Santa Cruz, CA 95062	Auckland Nuclear Accessory Company Ltd. P.O. Box 16066 Auckland, 3. New Zealand
Eastman Whipstock GmbH Gutenbergstrasse 3 3005 Hannover-Westerfeld West Germany	Geotechnical Instruments Ltd. Station House, Old Warwick Road Leamington Spa, Warwickshire CV31 3HR England
Huggenberger AG Holstrasse 176 CH-8040 Zürich Switzerland	IRAD GAGE Etna Road Lebanon, NH 03766
Leica AG CH-9435 Heerbrugg Switzerland	Measurement Devices Limited 11211 Richmond Avenue, Suite 106, Building B Houston, TX 77082
Maihak AG Semperstrasse 38 D-2000 Hamburg 60 West Germany	Roctest Ltée Ltd. 665 Pine Montreal, P.Q. Canada J4P 2P4
RST Instruments Ltd. 1780 McLean Avenue Port Coquitlam, B.C. Canada V3C 4K9	Schaevitz Engineering P.O. Box 505 Camden, NJ 08101
Serata Geomechanics, Inc. 4124 Lakeside Drive Richmond, CA 94806	SINCO (Slope Indicator Co.) 3668 Albion Place N. Seattle, WA 98103
Soil Instruments Ltd. Bell Lane, Uckfield East Sussex TN22 1QI England	Solexperts AG Postfach 274 CH-8034 Zürich Switzerland
Solinst Canada Ltd. 2440 Industrial St. Burlington, Ontario Canada L7P 1A5	SIS Geotecnica Via Roma, 15 20090 Segrate (Mi) Italy
Spectron Engineering 800 West 9th Avenue Denver, CO 80204	Spectron Glass and Electronics Inc. 595 Old Willets Path Hauppauge, NY 11788
Telemac 2 Rue Auguste Thomas 92600 Asnieres France	Terrametrics 16027 West 5th Avenue Golden, CO 80401
P & S Enterprises, Ltd. 240 South Monaco Pkwy, # 302 Denver, CO 80224	Edi Meier & Partner 8408 Winterthur Switzerland

16

Velocity Measurement

Charles P. Pinney
Pinney Technologies, Inc.

William E. Baker
University of New Mexico

16.1	Introduction.....	16-1
16.2	Measurement of Linear Velocity	16-2
	Reference-Based Measurement • Seismic Devices	
16.3	Velocity: Angular	16-9
	Relative: Tachometer • Absolute: Angular Rate Sensors	
16.4	Conclusion	16-15

16.1 Introduction

The *linear velocity* of an object, or more correctly a particle, is defined as the time rate of change of position of the object. It is a vector quantity, meaning it has a direction as well as a magnitude, and the direction is associated with the direction of the change in position. The magnitude of velocity is called the speed (or pace), and it quantifies how fast an object is moving. This is what the speedometer in a car tells you; thus, the speedometer is well named. Linear velocity is always measured in terms of, or from, some reference object. Thus, the speedometer of a car tells how fast one is moving relative to the earth. Usually, linear velocity is identified using only the term “velocity.” Common units for velocity include meters per second and miles per hour, but any similar combination of units of length per unit of time is correct.

The *rotational velocity* (or angular velocity) of an object is defined as the time rate of change of angular position, and it is a measure of how fast an object is turning. It is completely analogous to linear velocity, but for angular motion. Common units are revolutions per minute, but any angular unit of measurement per unit of time can be used. Rotational velocity is a vector quantity also, with the direction of the vector being the same as the direction of the axis about which object is turning. For example, with a car stopped at a stop light with the motor running, the rotational velocity of the crankshaft of the motor is given by a magnitude (rotational speed), say 800 rpm (rev/min), and a direction associated with the direction in which the crankshaft is pointing. The axis of rotation of the object may be moving, rather than fixed as when the car is turning a corner. The roll, yaw, or pitch velocity of an airplane would be given in terms of rotational speeds about each of the turning axes in the same manner as for a crankshaft.

Usually, the reference from which linear or rotational velocity is given is understood from the context of the problem. It is often not stated explicitly. The measurement method used defines the reference.

Applications for velocity measurement include:

1. Controlling the speed at which metal stock is fed into a machine tool. If the metal is fed too quickly the result could be premature tool wear or it could even lead to machine failure. Feeding the material too slowly will reduce the yield of the machine tool.
2. Measuring the approach speed of a robotic tool onto its target.
3. Monitoring the speed of a generator in an electric power station.
4. An airport radar system measuring the speed of an approaching aircraft using the Doppler effect.
5. Measuring an automobile’s wheel speed in order to provide feedback to an antilock brake system.

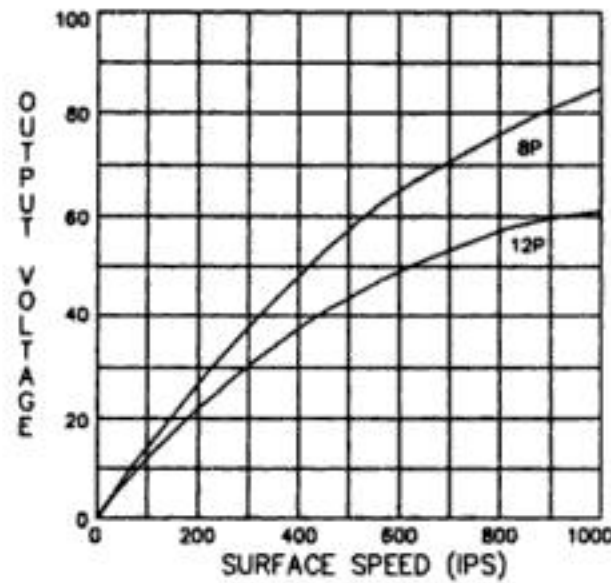


FIGURE 16.9 Magnetic speed sensor output voltage vs. speed. (Courtesy: Smith Systems, Inc., Brevard, NC.)

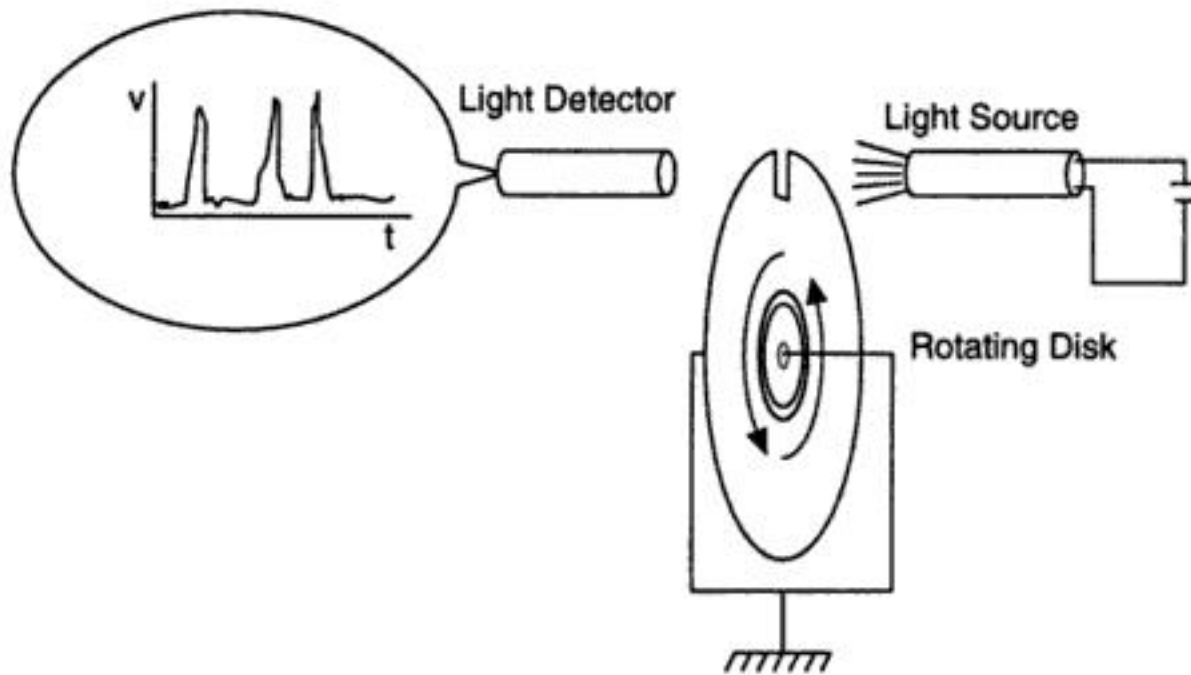


FIGURE 16.10 A slotted disk provides one pulse output for each rotation.

- Transmission speed
- Engine rpm
- Over/under speed
- Wheel speed
- Pump shaft speed
- Multiple engine synchronization
- Feedback for speed control
- Crankshaft position/engine timing
- Computer peripheral speeds

The typical specifications for magnetic speed sensors are given by a graph of output voltage versus surface speed in inches per second, as in Figure 16.9.

Sources for magnetic sensors include Smith Systems of Brevard, NC; Optek Technology of Carrollton, TX; Allied-Signal of Morristown, NJ; and Baluff of Florence, KY.

Optical Sensors

Optical methods of angular velocity detection employ a light emitter and a light detector. A light-emitting diode (LED) paired with a light-sensitive diode is the most common arrangement.

A slotted disk is placed in the axis of a rotating shaft. Each slot or slit will allow the light to pass through the disk. Figure 16.10 shows a typical arrangement. The detector will generate a pulse train with a rate proportional to the angular velocity.

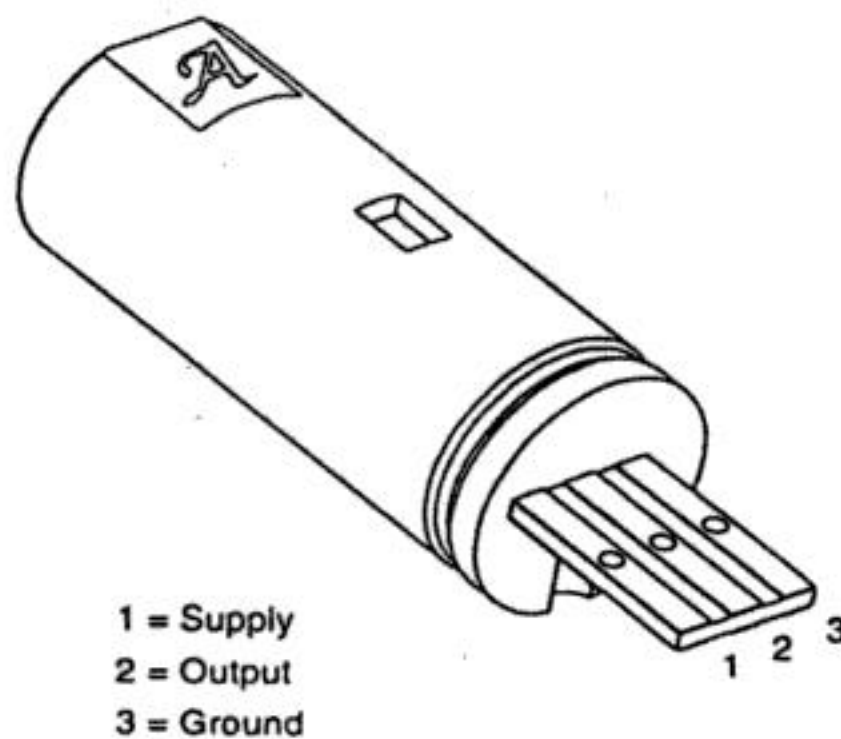


FIGURE 16.11 Hall-effect gear tooth sensor. (Courtesy: Allegro Microsystems, Inc., Worcester, MA.)

The effects of external light sources must be considered in the application of optical sensors.

Sources of optical sensor systems include Scientific Technologies of Fremont, CA; Banner Engineering Corp. of Minneapolis, MN; and Aromat Corp. of New Providence, NJ.

Hall Effect

The Hall effect describes the potential difference that develops across the width of a current-carrying conductor. E.H. Hall first used this effect in 1879 to determine the sign of current carriers in conductors. Hall effect devices are finding their way into many sensing applications. A typical Hall effect sensor application is the wheel speed sensor for antilock braking systems in automobiles. The Allegro ATS632LSC gear-tooth sensor, shown in Figure 16.11, is an optimized Hall-effect IC/magnet combination. The sensor consists of a high-temperature plastic shell that holds together a compound samarium-cobalt magnet, a single-element self-calibrating Hall effect IC, and a voltage regulator. The operation of this circuit is shown in Figure 16.12.

Wiegand Effect

The Wiegand effect is useful for proximity sensing, tachometry, rotary shaft encoding, and speed sensing in applications such as:

- Electronic indexing for water, gas, and electric meters and remote metering systems
- Measuring shaft speed in engines and other machinery
- Tachometers, speedometers, and other rotational counting devices

Wiegand effect technology employs unique magnetic properties of specially processed, small-diameter ferromagnetic wire. By causing the magnetic field of this wire to suddenly reverse, a sharp, uniform voltage pulse is generated. This pulse is referred to as a Wiegand pulse. Sensors utilizing this effect require only a few simple components to produce sharply defined voltage pulses in response to changes in the applied magnetic field. These sensors consist of a short length of Wiegand wire, a sensing coil, and alternating magnetic fields that generally are derived from small permanent magnets.

The major advantages of the Wiegand effect based sensors are:

- No external power requirement
- Two-wire operation
- Noncontact with no wear
- 20 kHz pulse rate
- High-level voltage output pulse
- Wide operating temperature range (e.g., -40°C to $+125^{\circ}\text{C}$)

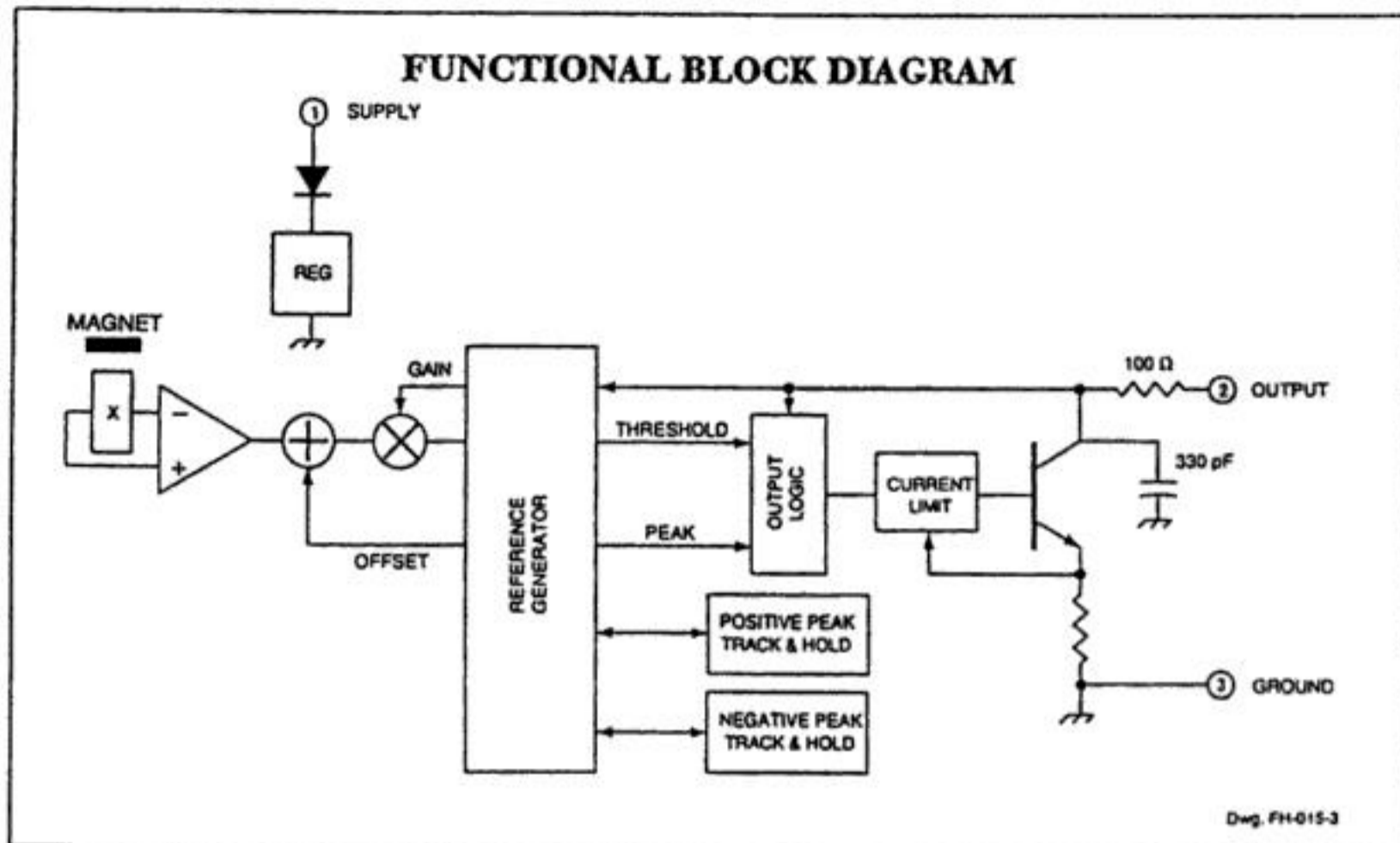


FIGURE 16.12 Hall-effect gear tooth sensor circuit. (Courtesy: Allegro Microsystems, Inc., Worcester, MA.)

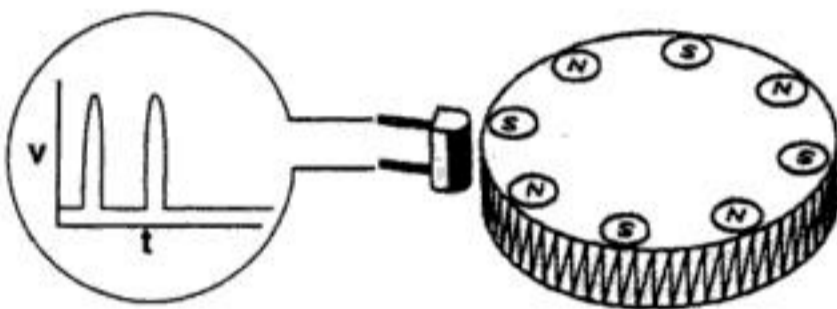


FIGURE 16.13 Small magnets cause sudden reversal in the ferromagnetic wire in a Wiegand sensor. (Courtesy: HID Corporation, North Haven, CT.)

When an alternating magnetic field of proper strength is applied to the Wiegand wire, the magnetic field of the core switches polarity and then reverses, causing the Wiegand pulse to be generated, as shown in Figure 16.13. The magnetic switching action of the Wiegand wire induces a voltage across the pickup coil of approximately $10 \mu\text{s}$ duration. These alternating magnetic fields are typically produced by magnets that are affixed to the rotating or moving equipment, by a stationary read head and moving Wiegand wires, or by an alternating current generated field.

Absolute: Angular Rate Sensors

Gyroscopes

Many absolute angular rate-measuring devices fall under the designation of gyroscope. A mechanical gyroscope is a device consisting of a spinning mass, typically a disk or wheel, mounted on a base so that its axis can turn freely in one or more directions and thereby maintain its orientation regardless of any movement of the base. It is important to make an initial distinction between angular velocity gyros and rate-integrating gyros. Angular velocity gyros are used to measure motion and as signal inputs to stabilization systems. Rate-integrating gyros are used as the basis for highly accurate inertial navigation systems. They allow a stable platform to maintain a fixed attitude with reference. These devices can be very complex. Three gyros are often teamed with three double-integrated accelerometers to provide an accurate measurement of absolute vehicle motion.

Ricardo Dao of Humphrey Inc. provided an excellent comparison of angular rate sensors in an article in *Measurements & Control* [14]. The five different technologies are summarized below.

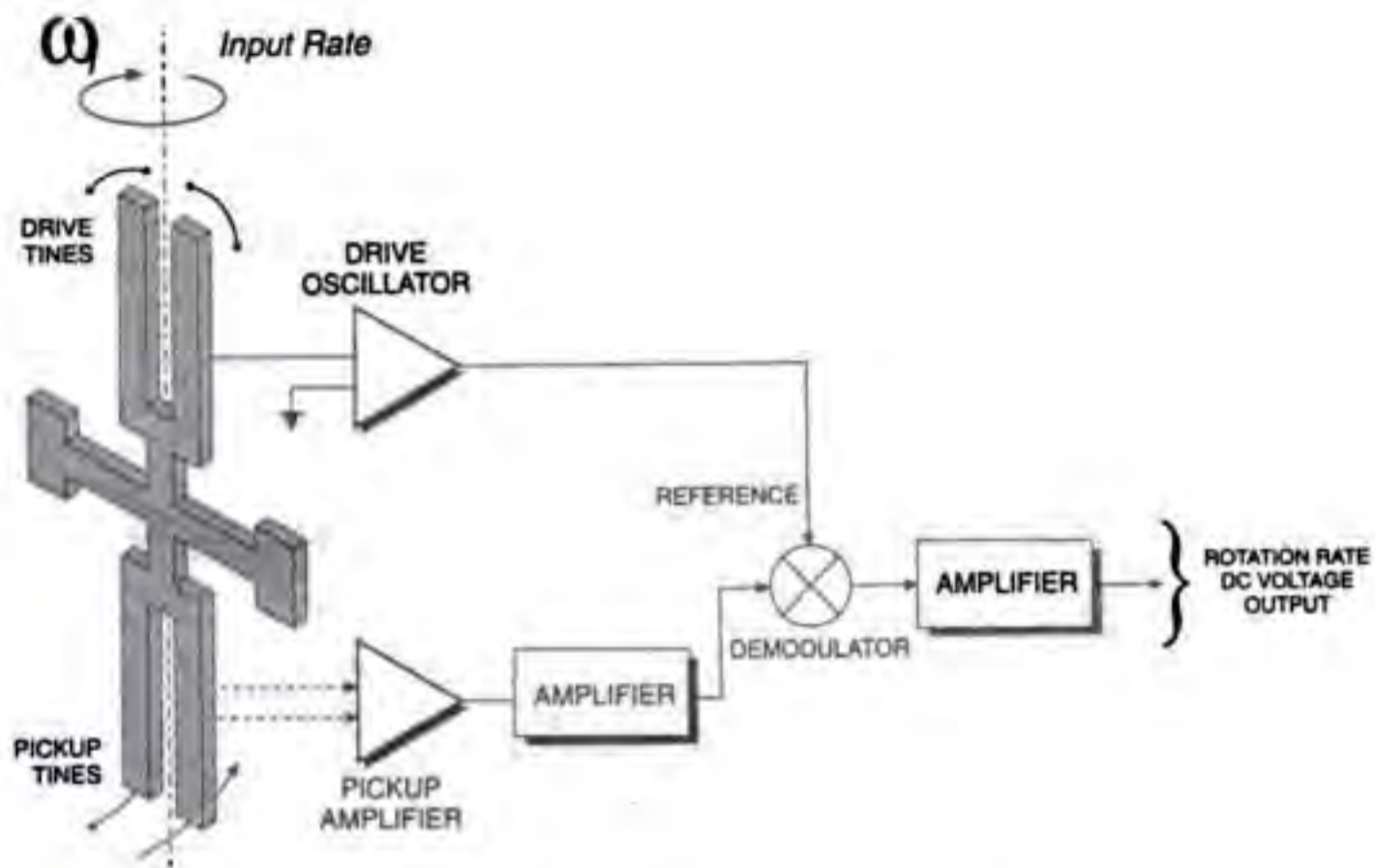


FIGURE 16.14 A vibrating quartz tuning fork uses the Coriolis effect to sense angular velocity. (Courtesy: BEI Sensors and Systems Co., Concord, CA.)

Spinning mass: The traditional gyro consists of a spinning wheel in a gimbaled frame. The principle of conservation of angular momentum provides the measurement tool.

Fluidic: A stream of helium gas flows past two thin tungsten wires [14]. The tungsten wires act as two arms of a Wheatstone bridge. At rest, the gas flow cools the sensing wires equally and the transducer bridge is balanced with zero output. When angular motion is applied to the sensor, one sensor wire will be subjected to increased flow while the other will see less flow. The resistance of the two wires will change and the bridge will be unbalanced. The sensor will produce a voltage output proportional to the angular velocity.

A pump is used to circulate the helium gas. This pump is a piezoelectric crystal circular disk that is excited with an external circuit. The pump produces a laminar flow of relatively high-velocity gas across the two parallel sensing wires.

Piezoelectric vibration: A number of angular velocity sensors have been developed that use micromachined quartz elements. A number of shapes are used, but the operating principle is similar for each. The quartz element vibrates at its natural frequency. Angular motion causes a secondary vibration that, when demodulated, is proportional to angular vibration. A description of one design follows.

The QRS and GyroChip™ family of products uses a vibrating quartz tuning fork to sense angular velocity [15, 16]. Using the Coriolis effect, a rotational motion about the sensor's longitudinal axis produces a dc voltage proportional to the rate of rotation. Figure 16.14 shows that the sensor consists of a microminiature double-ended quartz tuning fork and supporting structure, all fabricated chemically from a single wafer of monocrystalline piezoelectric quartz (similar to quartz watch crystals).

Use of piezoelectric quartz material simplifies the active element, resulting in exceptional stability over temperature and time. The drive tines, being the active portion of the sensor, are driven by an oscillator circuit at a precise amplitude, causing the tines to move toward and away from each another at a high frequency.

Each tine will have a Coriolis force acting on it of: $\{F = 2m W_i \times V_r\}$ where the tine mass is m , the instantaneous radial velocity is V_r , and the input rate is W_i . This force is perpendicular to both the input rate and the instantaneous radial velocity.

The two drive tines move in opposite directions, and the resultant forces are perpendicular to the plane of the fork assembly and in opposite directions. This produces a torque that is proportional to the

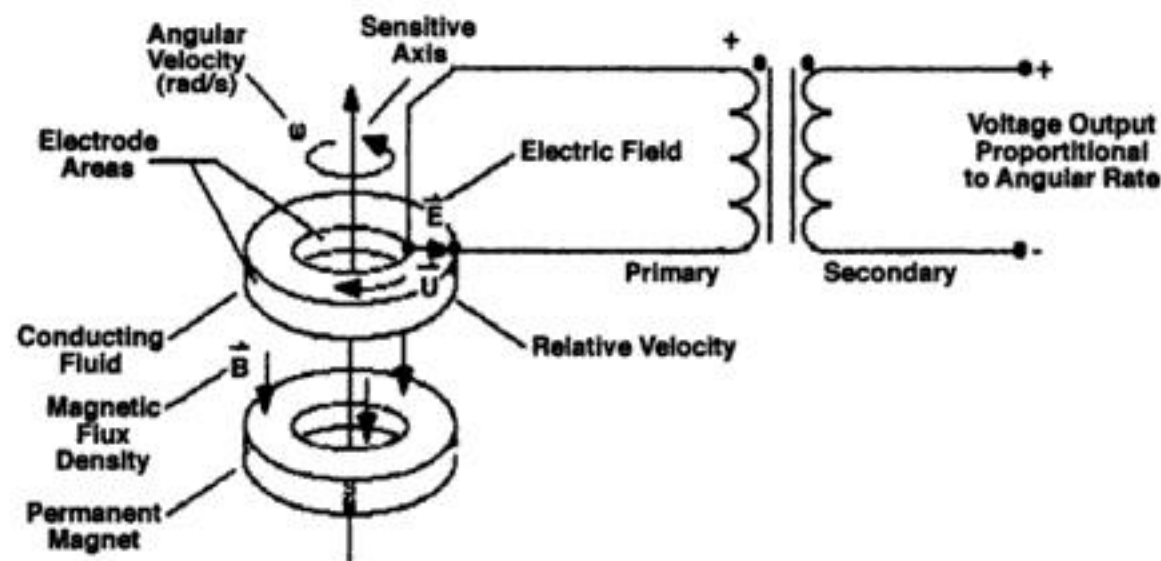


FIGURE 16.15 Magnetohydrodynamic angular rate sensor. (Courtesy: ATA Sensors, Albuquerque, NM.)

input rotational rate. Since the radial velocity is sinusoidal, the torque produced is also sinusoidal at the same frequency of the drive tines, and in-phase with the radial velocity of the tine.

The pickup tines, being the sensing portion of the sensor, respond to the oscillating torque by moving in and out of plane, producing a signal at the pickup amplifier. After amplification, those signals are demodulated into a dc signal that is proportional to the rotation of the sensor.

The output signal of the GyroChip™ reverses sign with the reversal of the input rate since the oscillating torque produced by the Coriolis effect reverses phase when the direction of rotation reverses. The GyroChip™ will generate a signal only with rotation about the axis of symmetry of the fork; that is, the only motion that will, by Coriolis sensing, produce an oscillating torque at the frequency of the drive tines. This also means that the GyroChip™ can truly sense a zero rate input.

MHD effect: The magnetohydrodynamic angular rate sensor is used to measure angular vibrations in the frequency range of 1 Hz to 1000 Hz. It is used where there is a high shock environment and a high rate of angular motion such as 10 to 250 rad s^{-1} . It does not measure a constant or dc velocity. It is used to measure impacts shorter than 1 s duration and vibrations between 1 Hz and 1000 Hz.

The principle of operation is illustrated in Figure 16.15 [17, 18]. A permanent magnet is attached to the outer case of the sensor. When the case turns, a moving magnetic field is produced (B). There is also a conductive fluid inside the sensor. When the sensor case turns, the fluid tends to stay in one place, according to Newton's first law. This produces a relative motion (U) between a magnetic field and conductor. This motion will produce a voltage (E) across the conductor proportional to relative velocity according to Faraday's law.

Since the fluid is constrained to move in an angular path, the voltage signal will be proportional to angular velocity about the center axis of the sensor. Due to this constraint, the sensor is insensitive to linear motion. The voltage signal is amplified through a transformer or an amplifier for output to a measuring device.

Fiber optic/laser: A beam of light is directed around the axis of rotation. A phase shift of the optical or laser beam is detected to measure angular velocity. The principle of operation is similar to the Doppler shift.

Differenced and integrated accelerometers: An array of accelerometers can be used to measure angular motion. The output of the accelerometers is differenced when they are aligned, or summed when they are mounted in opposite directions. This differencing will eliminate the linear component of motion. As shown in Figure 16.16, the magnitude of the differenced signals, a_1 and a_2 , is divided by the distance between the two sensors, l . This gives a measure of angular acceleration. The angular acceleration is integrated over time to give angular velocity. It is important to address the same concerns in this process as when integration was discussed in the linear section. It is assumed that there is a rigid mounting structure between the two accelerometers.

This technique is commonly applied to crash testing of anthropomorphic test devices (ATDs). The ATDs are used in automotive crash testing and aerospace egress system testing.

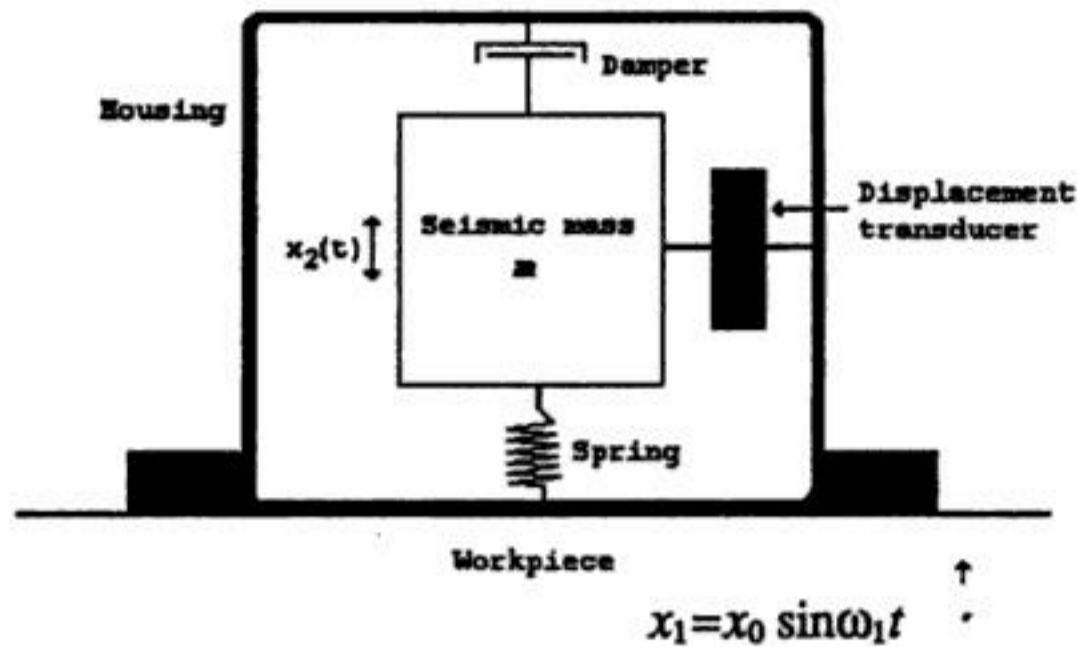


FIGURE 17.1 A typical deflection-type seismic accelerometer. In this basic accelerometer, the seismic mass is suspended by a spring or cantilever inside a rigid frame. The frame is connected to the vibrating structure; as vibrations take place, the mass tends to remain fixed so that relative displacements can be picked up. They are manufactured in many different types and sizes and they exhibit diverse characteristics.

true static measurements, passive sensors must be used. In general, accelerometers are preferred over displacement and velocity sensors for the following reasons:

1. They have a wide frequency range from zero to very high values. Steady accelerations can easily be measured.
2. Acceleration is more frequently needed since destructive forces are often related to acceleration rather than to velocity or displacement.
3. Measurement of transients and shocks can readily be made, more easily than displacement or velocity sensing.
4. Displacement and velocity can be obtained by simple integration of acceleration by electronic circuitry. Integration is preferred over differentiation.

Accelerometers can be classified in a number of ways, such as *deflection* or *null-balance* types, mechanical or electrical types, dynamic or kinematic types. The majority of industrial accelerometers can be classified as either deflection type or null-balance type. Those used in vibration and shock measurements are usually the deflection types, whereas those used for measurements of motions of vehicles, aircraft, etc. for navigation purposes may be either type. In general, null-balance types are used when extreme accuracy is needed.

A large number of practical accelerometers are of the deflection type; the general configuration is shown in Figure 17.1. There are many different deflection-type accelerometers. Although their principles of operation are similar, they only differ in minor details, such as the spring elements used, types of damping provided, and types of relative motion transducers employed. These types of accelerometers behave as second-order systems; the detailed mathematical analysis will be given in later sections.

Accelerometers can be classified as *dynamic*, meaning that the operation is based on measuring the force required to constrain a seismic mass to track the motion of the accelerated base, such as spring-constrained-slug types. Another type is the *kinematic* accelerometer, which is based on timing the passage of an unconstrained proof mass from spaced points marked on the accelerated base; this type is found in highly specific applications such as interspace spacecraft and in gravimetry type measurements.

For practical purposes, accelerometers can also be classified as *mechanical* or *electrical*, depending on whether the restoring force or other measuring mechanism is based on mechanical properties, (e.g., the law of motion, distortion of a spring or fluid dynamics, etc.) or on electrical or magnetic forces.

Calibrations of accelerometers are necessary in acceleration, vibration, and shock sensing. The calibration methods can broadly be classified to be *static* or *dynamic*. Static calibration is conducted at one

or several levels of constant acceleration. For example, if a tilting table calibration method is selected, the vertical component of the free fall is used without a choice of magnitude. On the other hand, if a centrifuge is selected, it produces a constant acceleration as a function of the speed of rotation, and the magnitudes can be chosen in a wide range from 0 to well over 50,000 g. The dynamic calibration is usually done using an electrodynamic shaker. The electrodynamic shaker is designed to oscillate in a sinusoidal motion with variable frequencies and amplitudes. They are stabilized at selected levels of calibration. This is an absolute method that consists of measuring the displacement with a laser interferometer and a precise frequency meter for accurate frequency measurements. The shaker must be driven by a power amplifier, thus giving a sinusoidal output with minimal distortion. The National Bureau of Standards uses this method as a reference standard. Precision accelerometers, mostly of the piezoelectric type, are calibrated by the absolute method and then used as the working standard. A preferred method is back-to-back calibration, where the test specimen is directly mounted on the working standard that, in turn, is mounted on an electrodynamic shaker.

Before providing details of different type of accelerometers, the common features such as accelerometer dynamics, velocity, distance, shock frequency responses, etc. will be introduced in the next section.

17.1 Accelerometer Dynamics: Frequency Response, Damping, Damping Ratio, and Linearity

This section concerns the physical properties of acceleration, vibration, and shock measurements in which accelerometers are commonly used. A full understanding of accelerometer dynamics is necessary in relation to characteristics of acceleration, vibration, and shock. The vibrations can be periodic, stationary random, nonstationary random, or transient.

Periodic Vibrations

In periodic vibrations, the motion of an object repeats itself in an oscillatory manner. This can be represented by a sinusoidal waveform:

$$x(t) = X_{\text{peak}} \sin(\omega t) \quad (17.1)$$

where $x(t)$ = time-dependent displacement

ω = $2\pi f$ = angular frequency

X_{peak} = maximum displacement from a reference point

The velocity of the object is the time rate of change of displacement:

$$u(t) = dx/dt = \omega X_{\text{peak}} \cos(\omega t) = U_{\text{peak}} \sin(\omega t + \pi/2) \quad (17.2)$$

where $u(t)$ = time-dependent velocity

$U_{\text{peak}} = \omega X_{\text{peak}}$ = maximum velocity

The acceleration of the object is the time rate change of velocity:

$$a(t) = du/dt = d^2 x/dt^2 = -\omega^2 X_{\text{peak}} \sin(\omega t) = A_{\text{peak}} \sin(\omega t + \pi) \quad (17.3)$$

where $a(t)$ = time-dependent acceleration

$A_{\text{peak}} = \omega^2 X_{\text{peak}} = \omega U_{\text{peak}}$ = maximum acceleration

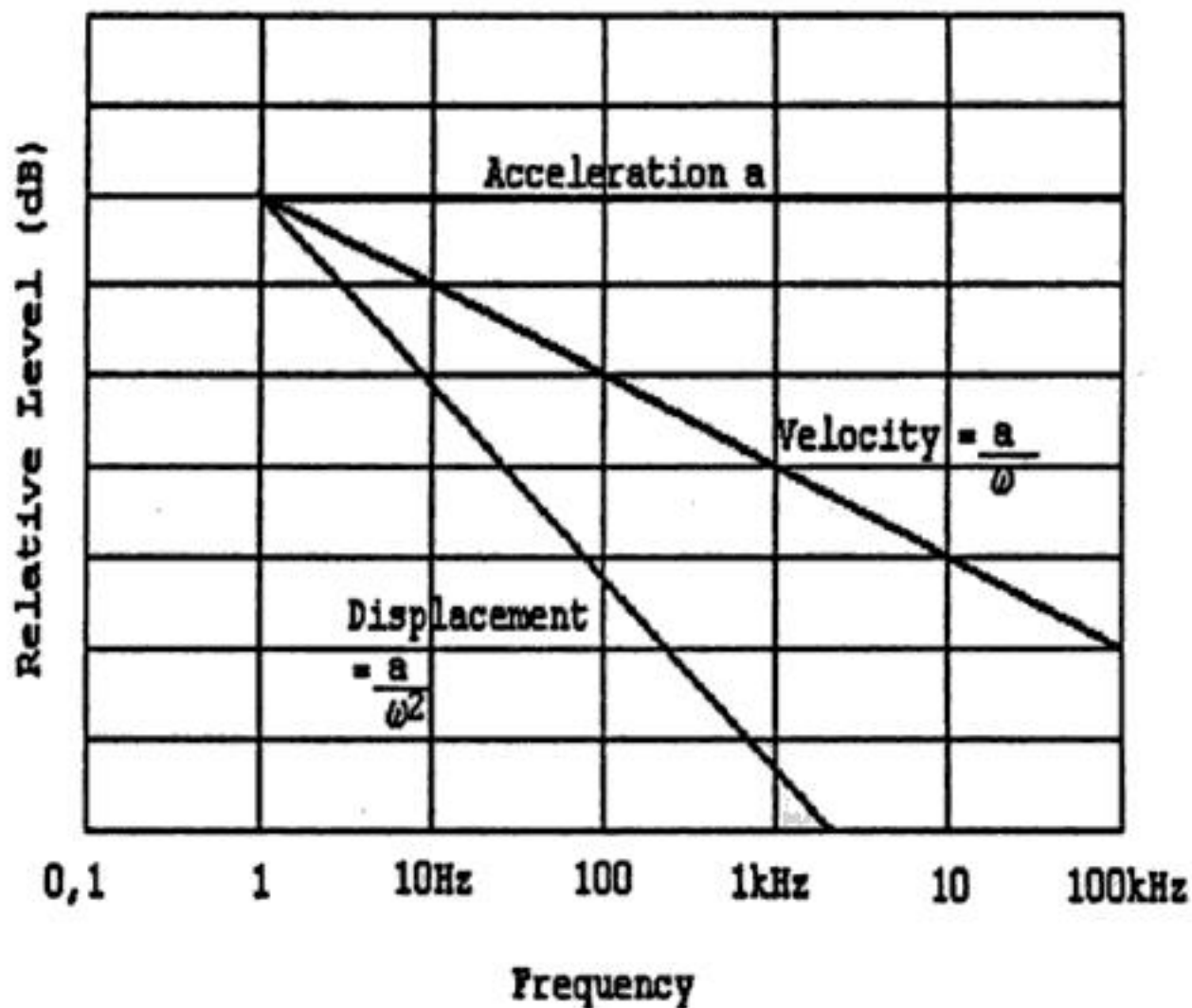


FIGURE 17.2 Logarithmic relationship between acceleration, velocity, and displacement. Velocity at a particular frequency can be obtained by dividing acceleration by a factor proportional to frequency. For displacement, acceleration must be divided by a factor proportional to the square of the frequency. Phase angles need to be determined separately, but they can be neglected in time-averaged measurements.

From the above equations, it can be seen that the basic form and the period of vibration remains the same in acceleration, velocity, and displacement. But velocity leads displacement by a phase angle of 90° and acceleration leads velocity by another 90° . The amplitudes of the three quantities are related as a function of frequency, as shown in Figure 17.2.

In nature, vibrations can be periodic but not necessarily sinusoidal. If they are periodic but nonsinusoidal, they can be expressed as a combination of a number of pure sinusoidal curves, described by Fourier analysis as:

$$x(t) = X_0 + X_1 \sin(\omega_1 t + \phi_1) + X_2 \sin(\omega_2 t + \phi_2) + \dots + X_n \sin(\omega_n t + \phi_n) \quad (17.4)$$

where $\omega_1, \omega_2, \dots, \omega_n$ = frequencies (rad s^{-1})

X_0, X_1, \dots, X_n = maximum amplitudes of respective frequencies

$\phi_1, \phi_2, \dots, \phi_n$ = phase angles

The number of terms may be infinite: the higher the number of elements, the better the approximation. These elements constitute the *frequency spectrum*. The vibrations can be represented in time domain or frequency domain, both of which are extremely useful in the analysis. As an example, in Figure 17.3, the time response of the seismic mass of an accelerometer is given against a rectangular pattern of excitation of the base.

Stationary Random Vibrations

Random vibrations are often met in nature where they constitute irregular cycles of motion that never repeat themselves exactly. Theoretically, an infinitely long time record is necessary to obtain a complete description of these vibrations. However, statistical methods and probability theory can be used for the

where $z = x_2 - x_1 =$ the relative motion between the mass and the base

$x_1 =$ displacement of the base

$x_2 =$ displacement of the mass

$\theta =$ the angle between sense axis and gravity

In order to lay a background for further analysis, taking the simple case, the complete solution to Equation 17.5 can be obtained by applying the superposition principle. The superposition principle states that if there are simultaneously superimposed actions on a body, the total effect can be obtained by summing the effects of each individual action.

Equation 17.5 describes essentially a second-order system that can be expressed in Laplace transform as:

$$X(s)/F(s) = 1/ms^2 + cs + k \quad (17.7)$$

or

$$X(s)/F(s) = K/[s^2/\omega_n^2 + 2\zeta s/\omega_n + 1] \quad (17.8)$$

where $s =$ the Laplace operator

$K = 1/k =$ static sensitivity

$\omega_n = \sqrt{k/m} =$ undamped critical frequency, rad/s

$\zeta = c/2\sqrt{km} =$ damping ratio

As can be seen, in the performance of accelerometers, important parameters are the static sensitivity, the natural frequency, and the damping ratio, which are functions of mass, velocity, and spring constants. Accelerometers are designed to have different characteristics by suitable selection of these parameters.

Once the response is expressed in the form of Equations 17.7 and 17.8, analysis can be taken further, either in the time domain or in the frequency domain. The time response of a typical second-order system for a unit-step input is given in Figure 17.4. The Bode plot gain phase responses are depicted in Figure 17.5. Detailed discussions about frequency response, damping, damping ratio, and linearity are made in relevant sections, and further information can be obtained in the references.

Systems in which a single structure moves in more than one direction are termed *multi-degree-of-freedom systems*. In this case, the accelerations become functions of dimensions as d^2x/dt^2 , d^2y/dt^2 , and d^2z/dt^2 . Hence, in multichannel vibration tests, multiple transducers must be used to create uniaxial, biaxial, or triaxial sensing points for measurements. Mathematically, a linear multidegree-of-freedom system can be described by a set of coupled second-order linear differential equations; and when the frequency response is plotted, it normally shows one resonance peak per degree of freedom.

Frequently, acceleration and vibration measurements of thin plates or small masses are required. Attaching an accelerometer with a comparable mass onto a thin plate or a small test piece can cause "mass loading." Since acceleration is dependent on the mass, the vibration characteristics of the loaded test piece could be altered, thus yielding wrong measurements. In such cases, a correct interpretation of the results of the measuring instruments must be made. Some experimental techniques are also available for the correction of the test results in the form of performing repetitive tests conducted by sequentially adding small known masses and by observing the differences.

The following sections discuss different types of accelerometers.

17.2 Electromechanical Force-Balance (Servo) Accelerometers

Electromechanical accelerometers, essentially servo or null-balance types, rely on the principle of feedback. In these instruments, acceleration-sensitive mass is kept very close to a neutral position or zero displacement point by sensing the displacement and feeding back this displacement. A proportional

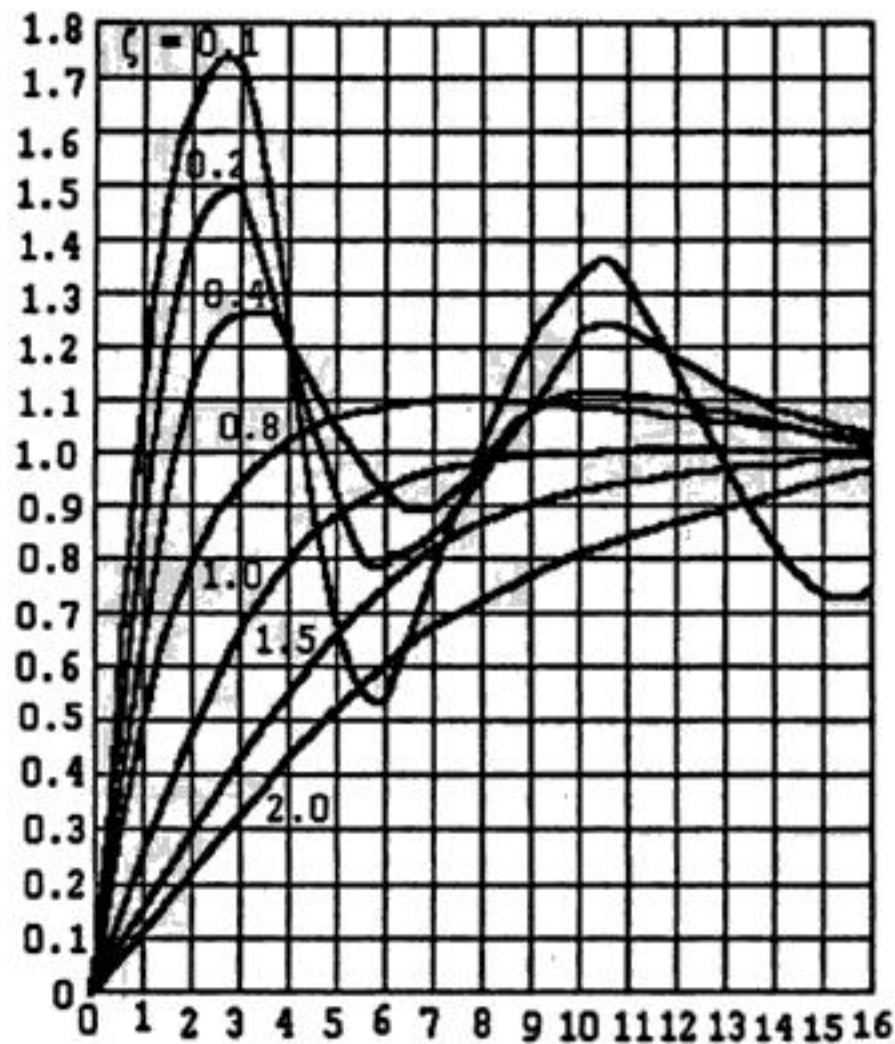


FIGURE 17.4 Unit step time responses of a second-order system with various damping ratios. The maximum overshoot, delay, rise, settling times, and frequency of oscillations depend on the damping ratio. Smaller damping ratios give faster response but larger overshoot. In many applications, a damping ratio of 0.707 is preferred.

magnetic force is generated to oppose the motion of the mass displaced from the neutral, thus restoring neutral position — just as a mechanical spring in a conventional accelerometer would do. The advantages of this approach are the better linearity and elimination of hysteresis effects as compared to mechanical springs. Also, in some cases, electric damping can be provided, which is much less sensitive to temperature variations.

One very important feature of null-balance type instruments is the capability of testing the static and dynamic performances of the devices by introducing electrically excited test forces into the system. This remote self-checking feature can be quite convenient in complex and expensive tests where it is extremely critical that the system operates correctly before the test commences. They are also useful in acceleration control systems, since the reference value of acceleration can be introduced by means of a proportional current from an external source. They are usually used for general-purpose motion measurements and monitoring low-frequency vibrations. They are specifically applied in measurements requiring better accuracy than achieved by those accelerometers based on mechanical springs as the force-to-displacement transducer.

There are a number of different types of electromechanical accelerometers: coil-and-magnetic types, induction types, etc.

Coil-and-Magnetic Type Accelerometers

These accelerometers are based on Ampere's law; that is: "a current carrying conductor disposed within a magnetic field experiences a force proportional to the current, the length of the conductor within the field, the magnetic field density, and the sine of the angle between the conductor and the field."

Figure 17.6 illustrates one form of accelerometer making use of the above principle. The coil is located within the cylindrical gap defined by a permanent magnet and a cylindrical soft iron flux return path. It is mounted by means of an arm situated on a minimum friction bearing so as to constitute an acceleration-sensitive seismic mass. A pick-off mechanism senses the displacement of the coil under

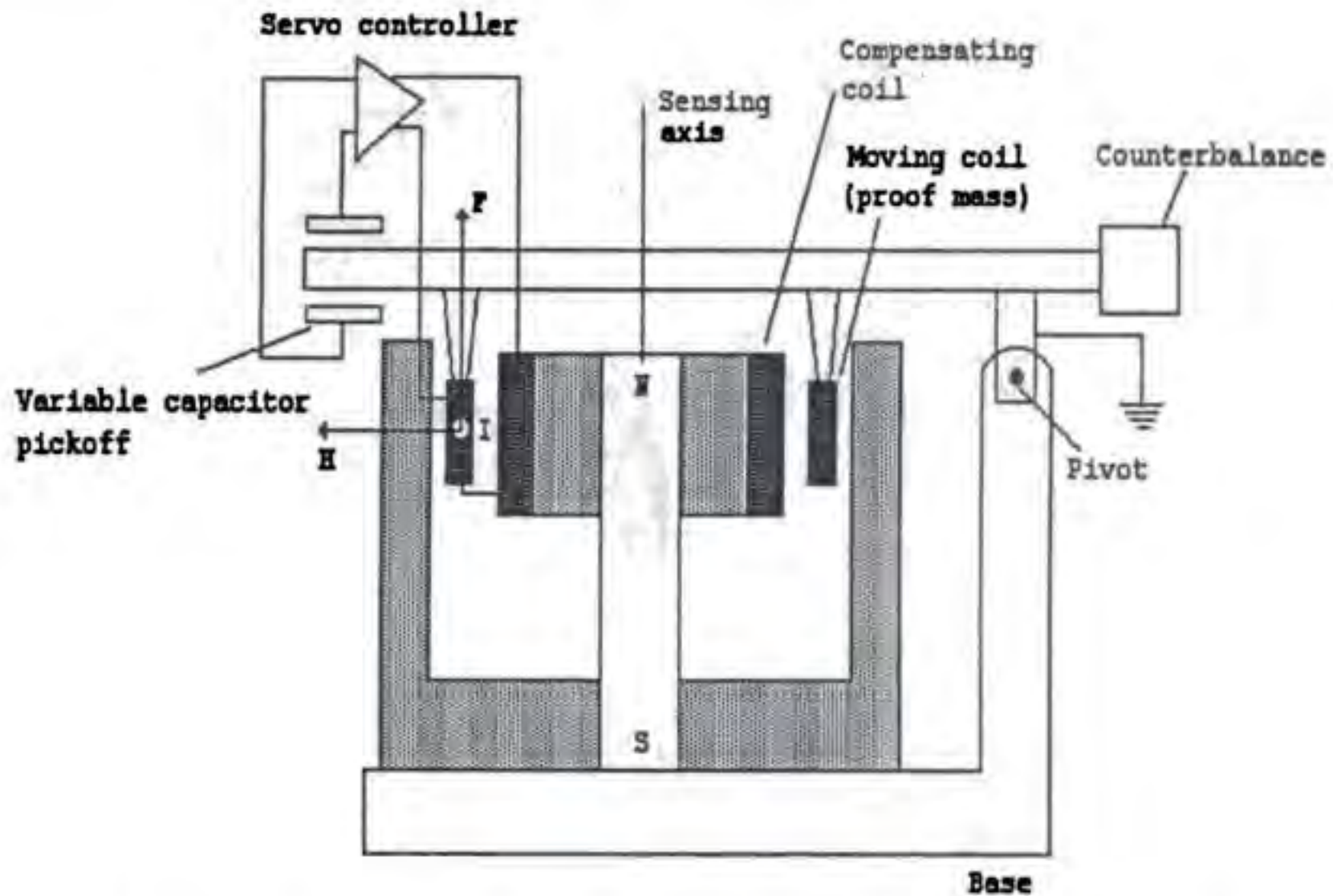


FIGURE 17.6 A basic coil and permanent magnet accelerometer. The coil is supported by an arm with minimum friction bearings to form a proof mass in a magnetic field. Displacement of the coil due to acceleration induces an electric potential in the coil to be sensed and processed. A servo system maintains the coil in a null position.

In these accelerometers, the magnetic structure must be shielded adequately to make the system insensitive to external disturbances or Earth’s magnetic field. Also, in the presence of acceleration, there will be a temperature rise due to i^2R losses. The effect of these i^2R losses on the performance is determined by the thermal design and heat transfer properties of the accelerometer. In many applications, special care must be exercised in choosing the appropriate accelerometer such that the temperature rises caused by unexpected accelerations cannot affect excessively the scale factors or the bias conditions.

A simplified version of another type of servo-accelerometer is given in Figure 17.7. The acceleration a of the instrument case causes an inertial force F on the sensitive mass m , tending to make it pivot in its bearings or flexure mount. The rotation θ from neutral is sensed by an inductive pickup and amplified, demodulated, and filtered to produce a current i_s directly proportional to the motion from the null. This current is passed through a precision stable resistor R to produce the output voltage signal and is applied to a coil suspended in a magnetic field. The current through the coil produces magnetic torque on the coil, which takes action to return the mass to neutral. The current required to produce magnetic torque that just balances the inertial torque due to acceleration is directly proportional to acceleration a . Therefore, the output voltage e_o becomes a measure of acceleration a . Since a nonzero displacement θ is necessary to produce the current i_s , the mass is not exactly returned to null, but becomes very close to zero because of the high gain amplifier. Analysis of the block diagram reveals that:

$$e_o/R = (mra - e_o K_c/R) \times (K_p K_a / K_s) / (s^2/\omega_{nl}^2 + 2\zeta_1 s/\omega_{nl} + 1) \tag{17.11}$$

Rearranging this expression gives:

$$mrRK_p K_a a / K_s = (s^2/\omega_{nl}^2 + 2\zeta_1 s/\omega_{nl} + 1 + K_c K_p K_a / K_s) e_o \tag{17.12}$$

THE
MEASUREMENT,
INSTRUMENTATION,
AND
SENSORS
HANDBOOK

Editor-in-Chief
John G. Webster

The **Measurement, Instrumentation, and Sensors Handbook** describes the use of instruments and techniques for practical measurements required in engineering, physics, chemistry, and the life sciences.

The book examines:

- Sensors
- Hardware
- Software
- Techniques
- Information processing systems
- Automatic data acquisition
- Reduction and analysis as well as their incorporation for control purposes

Organized according to the measurement problem, each section addresses the different ways of making a measurement for a given variable.

Chapters present three levels:

- Basic information without equations and a description of the subject that can be understood by the newcomer
- Detailed text and mathematical treatment essential for discovering applications and solving problems outside one's field of specialty
- Advanced applications of the subject, evaluative opinions, and areas for future study

The **Measurement, Instrumentation, and Sensors Handbook** provides a graded level of difficulty from start to finish, serving the reference needs of the broadest group of readers.

Edited by one of the more noted instrumentation experts in the world, the book contains nearly 150 contributions, covering all aspects on the design and implementation of various instrumentation.



Springer

ISBN 3-540-64830-5



9 783540 648307

<http://www.springer.de>