

FOUNDED 1925
INCORPORATED BY
ROYAL CHARTER 1961

*"To promote the advancement
of radio, electronics and kindred
subjects by the exchange of
information in these branches
of engineering."*

THE RADIO AND ELECTRONIC ENGINEER

The Journal of the Institution of Electronic and Radio Engineers

VOLUME 40 No. 3

SEPTEMBER 1970

Harmonized System for Quality Control

PROGRESS has been made in Great Britain in the past four years in building up a scheme for specifying electronic parts of assessed quality through the BS9000 series of specifications—these are prepared by Technical Committees of the British Standards Institution, supported by the Ministry of Technology. Probably one of the most significant advances of its kind in the electronics industry for some decades, the BS9000 concept provides for a single system of specifications in place of a variety of civil and military requirements and assured quality control for components graded for fitness-for-purpose. Over 100 British manufacturers, test houses and stockists have applied for approval under this scheme, generic and detailed specifications have been issued and components conforming to the requirements are now being put into production.

It has always been recognized since the inception of the BS9000 scheme that great mutual benefit would be gained if something similar could be promoted on a wider basis, multi-national or regional if not truly international. The persistence of British quality control engineers has now been justified by an agreed statement, issued simultaneously in the E.E.C. and E.F.T.A. countries, to announce the Harmonized System for Electronic Components which is being introduced in Western Europe. The System is intended to promote trade in electronic components between the participating countries by harmonizing their national systems of specifications and quality assurance for electronic components. In order to remove technical barriers to trade, the System will provide harmonized specifications which will be published individually by each country, but which will meet internationally agreed requirements. It will also provide for the multilateral recognition of approvals given by national inspectorates. It is important to note that the BS9000 Scheme will be fully compatible with the international Harmonized System.

The Comité Européen de Coordination des Normes Electrotechniques (CENEL) which comprises all E.E.C. and E.F.T.A. countries, has accepted overall responsibility for launching the Harmonized System for Electronic Components in these countries. CENEL has established a special committee, the CENEL Electronic Components Committee (CECC), which will be responsible for implementing the Harmonized System, except for certain quality assurance aspects. The responsibility for the quality assurance and inspection aspects of the Harmonized System is to be exercised by an independent Committee, known as the Electronic Components Quality Assurance Committee (ECQAC), and the second meeting of this Committee was held in London in June under the chairmanship of Mr. H. E. Drew, C.B., C.Eng., F.I.E.R.E., Director-General of Quality Assurance in the British Ministry of Technology.

The statement issued on behalf of CENEL envisages eventual widening of the Harmonized System to world-wide application as a desirable goal. It is not unduly optimistic to hope that just as the example of Great Britain has led to the European developments, so in turn will the International Electrotechnical Commission (IEC), who are regarding these with considerable interest, be lead to set up a comparable international scheme in due course.

Assured quality and performance for a given purpose and the simplicity of a single system of specifications in place of the numerous civil and military specifications clearly offer great advantages to manufacturer and user alike. It is now up to them to seize the opportunities offered.

F.W.S.

INSTITUTION NOTICES

Institution Premiums and Awards

The Council of the Institution announces that the following awards are to be made for outstanding papers published in *The Radio and Electronic Engineer* during 1969.

CLERK MAXWELL PREMIUM

'Optimum Transfer Functions for Feedback Control Systems with Plant Input Saturation' by D. R. Towill (October 1969).

REDIFFUSION TELEVISION PREMIUM

'Control of Gamma in C.R.T. Displays using Amplifiers with Exponential Negative Feedback' by S. L. Cachia (November).

P. PERRING THOMS PREMIUM

'Flashover in Picture Tubes and Methods of Protection' by A. Ciuciura (March)

J. LANGHAM THOMPSON PREMIUM

'Digital Computer Implementation of Bang-bang Process Control' by P. Atkinson and R. L. Davey (November).

HEINRICH HERTZ PREMIUM

'Transistor Abnormalities as Revealed by Current-Voltage Characteristics' by P. J. Holmes (November).

A. F. BULGIN PREMIUM

'A Radiometer for Measurement of the Noise Temperature of Low-noise Microwave Amplifiers' by J. W. Carter, H. N. Daghish and P. Moore (June).

LESLIE MCMICHAEL AWARD

'Global Communications: Current Techniques and Future Trends' by R. W. Cannon (May).

LORD BRABAZON AWARD

'Radar Polarization Comparisons in Sea-Clutter Suppression by Decorrelation and Constant False Alarm Rate Receivers' by J. Croney and A. Woroncow (October).

CHARLES BABBAGE AWARD

'A Stored Microprogram Control Unit using Tunnel Diodes' by N. E. Wiseman and P. C. Wright (March).

MARCONI AWARD

'A Wideband Amplitude Modulator as a Special Silicon Integrated Circuit' by A. Stewart and C. H. Jones (September).

The following Premiums and Awards are being withheld as papers published during the year which fall within their respective terms of reference were not of a sufficiently high standard:

Vladimir K. Zworykin Premium, Arthur Gay Premium, Dr. Norman Partridge Memorial Premium, Hugh Brennan Premium, Lord Rutherford Award.

The Premiums and Awards will be presented by the President of the Institution Mr. H. F. Schwarz, at the Annual General Meeting in London on Wednesday, 25th November, 1970.

Dinner of Council and Committees

A Dinner of the Council, its Committees and representatives will be held in the Lincoln and Manhattan Rooms, Savoy Hotel, London, on Thursday, 5th November, 1970. This will provide a rare occasion for all the members who assist the Institution to meet together socially and to be accompanied by their ladies and personal guests. It will also be an opportunity to express thanks to the Immediate Past President, Major-General Sir Leonard Atkinson, K.B.E., for his work for the Institution.

The charge for tickets, obtainable from 9 Bedford Square, is £4. 5s. each, which will include wines at table.

Thomson Lecture

This year's Thomson Lecture will be given by Sir Frederick Warner, B.Sc., Hon.D.Tech. C.Eng., Senior Partner of Cremer & Warner, whose subject will be 'Measurements, Models and Men'. The Lecture will be held at The Royal Institution of Great Britain, 21 Albemarle Street, London, W.1, at 6 p.m. on Thursday, 8th October, 1970.

Admission will be by ticket only and those wishing to attend should apply to:

The Secretary, The Institute of Measurement and Control, 20 Peel Street, London, W.8.

First British Quality and Reliability Convention

The National Council for Quality & Reliability is organising a Q & R Convention in London, 28th-29th October 1970, to promote the application of quality and reliability techniques by stressing profitability aspects. This, the first full-scale convention held by the National Council under its own auspices, is expected to play a vital role in helping more British products to become highly competitive in world markets. The Convention, entitled on 'Quality Rewards' is mainly directed at the numerically largest section of British industry—expanding and progressive companies in the small and medium range whose continued growth can be sustained by the use of Q & R measures tailored to their overall business situation.

After plenary opening sessions addressed by Dame Elizabeth Ackroyd, Director of the Consumer Council, Mr. J. R. F. Moss, Director of Naval Ship Production, Ministry of Defence (Navy), and Mr. E. L. G. Robbins, Member of the Bolton Committee of Inquiry on Small Firms, the Convention will divide up into a series of Management Seminars, Specialist Seminars for the control of Q & R, case study discussions and Q & R Clinics.

For further administrative details and application forms for the Convention and its Clinics, please write to: National Council for Quality & Reliability, Vintry House, Queen Street Place, London, E.C.4.

Adaptive Detection of Distorted Digital Signals

By

A. P. CLARK,

M.A., Ph.D., D.I.C., C.Eng., M.I.E.R.E.†

The paper describes a novel approach to the adaptive equalization of a channel, leading to an adaptive detector which promises to achieve in some applications a better performance over a slowly time-varying channel, for a given degree of equipment complexity, than is possible with the more conventional transversal-filter adaptive equalizer.

The transmitted signal contains a serial stream of binary data elements which are separated into orthogonal groups. Each group of elements is detected in a single detection process. The preferred method of detection is an iterative process, which is used first to determine the element binary values in the group and then to estimate the channel impulse response. By this means the detector is adjusted after each detection process so that it follows the variations in the channel transmission characteristics. No training signal is required, except at the start of transmission.

Mathematical Notation

$|x|$ magnitude of the real scalar x

$\{x_i\}$ the set x_1, x_2, \dots, x_k , where k is given in the text

(x_i) the row vector whose i th component is x_i

$[x_{ij}]$ the matrix whose component in the i th row and j th column is x_{ij}

X^T transpose of the matrix or vector X

1. Introduction

Much of the work so far carried out on the adaptive equalization of a time-varying channel has assumed the use of a transversal filter at the receiver.¹⁻¹¹ The input signal to the filter is a continuous stream of baseband data elements with an element rate of $1/T$ bauds. The filter is shown in Fig. 1 and contains a delay line tapped at T -second intervals. Each tap is connected through an amplifier or attenuator (and possibly an inverter) to the analogue adder. The tap gains $\{c_i\}$ are adjusted automatically to reduce the intersymbol interference in the output data signal. Most of the equipment in the adaptive equalizer is involved with the arrangements for adjusting the tap gains and is not shown in Fig. 1.

The adaptive equalizer acts essentially as a filter and is both simple and effective. However, the arrangement has three limitations. Firstly, for certain transmission characteristics, the channel cannot be equalized. Secondly, only approximate equalization can normally be achieved with a finite transversal filter, so that a reasonable number of taps, say around 20, is usually required. Thirdly, in order that the filter may adapt itself to variations in the channel transmission characteristics during the transmission of data, without the transmission of special training signals, the sequence of data element values must be reasonably random.

† Plessey Telecommunications Research Limited, Taplow Court, Taplow, Maidenhead, Berkshire.

This paper describes a different approach to adaptive equalization, which leads to an arrangement of adaptive detection. The system is only slightly more complex than the equivalent transversal equalizer but without any of the main disadvantages of the latter. The basic principles of the system are first described and then four suitable detection processes. Finally an outline is given of a particularly simple and effective adaptive detector.

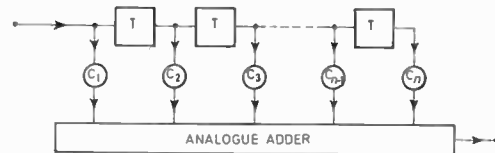


Fig. 1. Transversal filter.

The paper is aimed at those readers who have no previous knowledge of iterative methods and only a limited understanding of matrix algebra. The emphasis throughout is on basic concepts and practical usefulness. Details of the theoretical analysis and computer simulation tests, which have been carried out on the detection processes, are given in reference 12.

2. Principles of Detection Process

2.1. Basic Assumptions

Consider a synchronous serial data-transmission system as shown in Fig. 2. The input signal is a series of impulses $\{z_i \delta(t - iT)\}$ spaced at regular intervals of T seconds, where $\delta(t)$ is a unit impulse at time $t = 0$. Each impulse is a signal element. The elements are binary antipodal and have unit magnitude (area), so that for any integer i , z_i is 1 or -1 . It is assumed that the $\{z_i\}$ are statistically independent and equally likely to have either binary value.

The input to the modulator contains a low-pass filter to produce a baseband signal which is then used to modulate the carrier. The modulated-carrier signal may be a vestigial-sideband suppressed-carrier a.m. signal and the demodulator is a coherent detector whose output is a baseband signal. The modulator, transmission path and demodulator are assumed to be linear.

When no signal distortion is introduced in the transmission path, the nominal bandwidth of the baseband signal at the detector input is $1/2T$ Hz. The baseband signal is sampled at regular intervals of T seconds in the detector. The detector operates entirely on these sample values and its function is to obtain the best estimates $\{x_i\}$ of the transmitted $\{z_i\}$. The correct sampling instants are determined by a timing signal suitably synchronized to the received baseband signal.

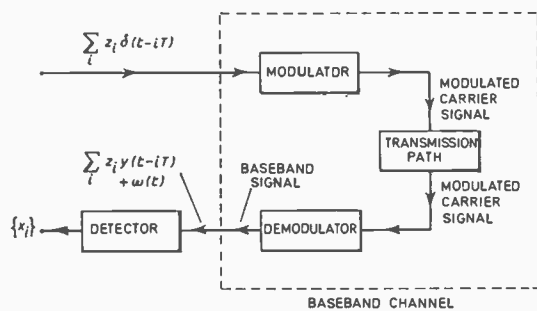


Fig. 2. Data transmission system.

It may readily be shown that the modulator, transmission path and demodulator are equivalent to a linear baseband channel. Suppose this has an impulse response $y(t)$. Where the transmission path is an h.f. radio link, $y(t)$ may vary slowly with time. Where it is a switched telephone circuit, $y(t)$ will in general differ from one transmission to the next but will not usually vary much during any one transmission. The data signal at the detector input is clearly

$$\sum_i z_i y(t-iT)$$

Suppose that additive white Gaussian noise is introduced at the output of the transmission path, giving the zero-mean Gaussian waveform $w(t)$ added to the data signal at the detector input. Although the additive noise over telephone circuits is not normally Gaussian, the assumption is adequate for our purposes. The sample values of $w(t)$ are taken to be statistically independent, which of course places a restriction on the overall response of the demodulator filters.¹

It is assumed that the various filters in the modulator and demodulator are designed so that, when there is

no signal distortion in the transmission path, there is no intersymbol interference at a sample value of the baseband signal in the receiver and the signal/noise ratio here is maximized. Reference 1 contains an excellent general treatment of this subject, which will not be considered further here.

The object of the present investigation is to determine the method in which the detector should use the sample values of the received baseband signal, in order to obtain the best compromise between performance and equipment economy, in those applications where the transmission path may introduce severe signal distortion. Thus in the arrangement of Fig. 2 we are concerned only with the design of the detector.

2.2. Detection of Orthogonal Groups of Signal Elements

Two groups of signal elements can be considered to be orthogonal when each of them produces no response in an optimum detection process on the other.

Where the baseband channel introduces negligible intersymbol interference in the sample values of the signal elements at the detector input, the optimum detector determines the binary value of each element from the sign of the corresponding sample value. Where there is appreciable intersymbol interference, such an arrangement is no longer optimum and may not even operate correctly in the absence of noise. A suitable transversal equalizer may, of course, be inserted at the input to the detector, in order to reduce the intersymbol interference to an acceptable level.¹⁻¹¹

An alternative approach to the optimum detection of the received signal elements is however suggested by studies into the optimum design and detection of a finite number of signals.¹²⁻¹⁸ This approach involves a simple modification to the data-transmission system of Fig. 2.

Suppose that with the most extreme time-dispersion of the received signal, a signal element can cause interference in the sample values of some or all of the p immediately preceding elements and in some or all of the q immediately following elements. The element stream is then divided into separate groups, each containing m consecutive elements which carry the transmitted data. Each group of m elements is separated from the following group by g elements, which are set to zero and act as a time guard band between the adjacent groups. Also

$$g = p + q \quad \dots\dots(1)$$

Associated with a group of m data elements there are

$$n = p + m + q = m + g \quad \dots\dots(2)$$

consecutive sample values which are dependent on

the particular group of data elements and independent of every other group. These n sample values are used for the detection of the m data elements. Clearly no data element in any one group can cause interference in the detection of an element in any other group, so that the different groups are orthogonal. n would be typically around 20 and g preferably less than $\frac{1}{2}n$, although in some applications a value as high as $\frac{1}{2}n$ may have to be used.

A special timing signal must of course now be transmitted to indicate to the detector the first sample value in each group of n .

Over any baseband channel likely to be used in practice, the impulse response decays sufficiently rapidly, away from its central peak, so that no serious error is introduced by assuming a finite time-dispersion of a received signal-element, even though the time dispersion may theoretically be infinite. Where the transmission path is an h.f. radio link or a switched telephone circuit, a time gap (time guard band) of around 3 or 4 ms between adjacent groups of elements, should be quite adequate.

Where the receiver has a fairly accurate prior knowledge of the n sample values corresponding to each of the two binary values, for every one of the m individual elements in a group, the optimum detection process determines which of the 2^m different combinations of the m binary values gives a resultant set of n sample values with the minimum mean square difference from the received n sample values. This detection process detects the sum of the m binary elements as a multi-level element having 2^m possible values. The detection process is optimum in the sense that it minimizes the probability of error (that is, the probability of one or more errors) in the detection of the m binary-elements. The disadvantage of the arrangement is that either the equipment complexity or the detection period is proportional to 2^m . In addition, for correct operation the detector may require a fairly accurate prior knowledge of the received signal level.¹²

An alternative and preferable approach to the detection of the m binary-elements in a group is through the solution of the appropriate set of simultaneous equations. The m elements are again detected simultaneously in a single detection process, and the operation of the arrangement can be explained as follows.

If there is no intersymbol interference, the n sample values of the first signal-element in a group of m , are

$$\underbrace{0 \dots 0}_p \underbrace{z_1 y_{p+1} \dots 0}_{m} \underbrace{0 \dots 0}_q$$

where z_1 is either 1 or -1 and carries the element binary value. y_{p+1} depends upon the channel and may vary slowly with time.

If there is intersymbol interference, the n sample values of the first element are

$$\underbrace{z_1 y_1 \ z_1 y_2 \dots z_1 y_{g+1}}_{g+1} \underbrace{0 \dots 0}_{m-1}$$

where y_j is the j th sample value of the first element when $z_1 = 1$. Clearly

$$g = n - m \quad \dots\dots(3)$$

and

$$y_j = 0 \quad \text{when } g+2 \leq j \leq n \quad \dots\dots(4)$$

y_j must be non-zero for at least one value of j in the range 1 to $g+1$, but it need not of course be non-zero for all j in this range.

The sample values of the i th signal-element in the group of m are

$$\underbrace{0 \dots 0}_{i-1} \underbrace{z_i y_1 \ z_i y_2 \dots z_i y_{g+1}}_{g+1} \underbrace{0 \dots 0}_{m-i}$$

and the n -component row-vector Y_i is now defined to be

$$\underbrace{0 \dots 0}_{i-1} \underbrace{y_1 y_2 \dots y_{g+1}}_{g+1} \underbrace{0 \dots 0}_{m-i}$$

for $i = 1, \dots, m$, so that the i th signal-element is given by $z_i Y_i$. z_i is of course a scalar and both z_i and Y_i are real. The binary value of the i th element is given by the sign of z_i , and

$$|z_i| = 1 \quad \text{for all } i \quad \dots\dots(5)$$

where $|z_i|$ is the magnitude of z_i .

The sum of the m signal-elements in a group is

$$\sum_{i=1}^m z_i Y_i = ZY \quad \dots\dots(6)$$

where Z is the row-vector (z_1, z_2, \dots, z_m) , written (z_i) , and Y is the $m \times n$ matrix whose i th row is Y_i . Clearly Y_i is the first row Y_1 shifted to the right by $i-1$ places, so that any one row of Y is a simple time-shift of any other.

It can readily be shown that the m vectors $\{z_i Y_i\}$ are always linearly independent. That is, no one of these can be expressed as a linear combination of some or all of the others. It can also be shown that ZY has a different value for every different combination of the m binary values and that the latter can always be uniquely determined from ZY , so long as the receiver has a prior knowledge of the $\{Y_i\}$ but not necessarily of the $\{z_i\}$.^{12, 13}

The determination of the vector Z from a given vector ZY , where the matrix Y is known, involves the solution of a set of simultaneous equations and these may be expressed in different ways. However, where the received signal contains additive Gaussian noise, there is one particular formulation of the simultaneous equations which minimizes the prob-

ability of error in a detection process. A detector which operates by solving this particular set of simultaneous equations will now be described.

2.3. Detection Process A

The n sample values of the received baseband signal, corresponding to a group of m signal-elements, are the n components of the row-vector $R = (r_i)$, where

$$R = ZY + W \quad \dots\dots(7)$$

W is an n -component row-vector whose components are the sample values of the additive noise in the baseband signal. It is assumed that the n components of W are sample values of statistically independent Gaussian random variables with zero mean and variance σ^2 .

The detector samples the received baseband signal at the correct n sampling points for a group of m signal-elements and it stores the corresponding vector R . Two stores are required, so that while one holds R for the detection process, the other is receiving the n sample values for the next vector R . The m signal-elements of the stored vector R are detected in the detection process A, shown in Fig. 3.

The correlation detector tuned to Y_i multiplies each component of R by the corresponding component of Y_i and adds the products to give the output signal

$$d_i = RY_i^T \quad \text{for all } i \quad \dots\dots(8)$$

where Y_i^T is the transpose of Y_i and RY_i^T is the inner product of R and Y_i .

The output signals from the m correlation detectors are the m components $\{d_i\}$ of the vector D . Thus

$$D = RY^T = (ZY + W)Y^T = ZA + WY^T \quad \dots\dots(9)$$

where

$$A = YY^T \quad \dots\dots(10)$$

$A = [a_{ij}]$ is an $m \times m$ real symmetric positive-definite matrix.

It can be shown that when the receiver knows the $\{Y_i\}$ but has no prior knowledge of the $\{z_i\}$ or σ^2 , then the best estimate it can make of the vector Z is the vector $X = (x_i)$, where

$$XA = D \quad \dots\dots(11)$$

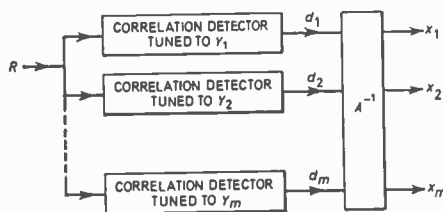


Fig. 3. Detection process A.

This minimizes the probability of error in a detection process, under the assumed conditions.¹² Since A is non-singular,

$$X = DA^{-1} \quad \dots\dots(12)$$

Thus to obtain the m estimates $\{x_i\}$ of the received signal values $\{z_i\}$, the output signals of the correlation detectors are fed through a network which performs the linear transformation A^{-1} on these signals, as shown in Fig. 3. It can be seen that the detection process A operates by solving the m linear simultaneous equations given by equation (11).

The wanted component in the output signal d_i from the i th correlation detector in Fig. 3 is

$$z_i Y_i Y_i^T = z_i a_{ii} \quad \dots\dots(13)$$

The correlation detector maximizes the ratio of the power level of this signal to the average power level of the noise component WY_i^T in d_i . However, d_i also contains $m-1$ components

$$z_j Y_j Y_i^T = z_j a_{ji} \quad \text{for } j \neq i \quad \dots\dots(14)$$

due to the other received signal-elements, so that there may be considerable intersymbol interference in d_i . The network A^{-1} processes the $\{d_i\}$ to eliminate all intersymbol interference and suitably adjusts the levels of the resultant signals to give the $\{x_i\}$ at its output terminals.

Since the signal elements $\{z_i Y_i\}$ are binary antipodal, the detection process is not affected by a constant attenuation applied to the $\{x_i\}$. Thus, under favourable conditions, each correlation detector can be a set of n attenuators, with inverters where necessary, and the network A^{-1} can be a set of m^2 attenuators together with arrangements for adding and subtracting the m signals at each of the m outputs.

In the final stage of the detection process, not shown in Fig. 3, the receiver examines the signs of the $\{x_i\}$ and allocates the appropriate binary values to the corresponding $\{z_i Y_i\}$. The detection process of course requires no prior knowledge of the received signal level.

The arrangement of Fig. 3 can be simplified to a set of m correlation detectors, where the i th detector is tuned to a vector whose inner product with Y_i is unity and whose inner product with Y_j , for $j \neq i$, is zero. These m vectors span the same m -dimensional subspace of the n -dimensional signal-space as is spanned by the $\{Y_i\}$. The output signals from the m correlation detectors in the simplified arrangement can be shown to be the required $\{x_i\}$.¹⁹ Where the channel impulse-response does not vary with time and is known at the receiver, this is clearly the preferred arrangement.

When the detector has no prior knowledge of the impulse response of the channel or when the impulse response varies with time, considerable equipment complexity may be involved with either of the arrangements just considered. The reason for this is that the component values of the network A^{-1} in Fig. 3 and of the correlation detectors in the simplified arrangement, cannot be derived directly or very simply from the received signal. The linear transformation A^{-1} in Fig. 3 may however be carried out in a basically different manner which does not suffer from this disadvantage. Four detection processes, each using such an arrangement, will now be described.

3. Detection Processes Suitable for Adaptive Working

3.1. Detection Process B

The generation of the matrix A^{-1} is the essential operation involved in the solution of the m linear simultaneous equations

$$XA = D \quad \text{.....(15)}$$

where A and D are given and X is to be determined. Analogue computer techniques provide a simple and effective means for solving these equations, where the matrix A is real, symmetric and positive definite as it is here.²² The application of such techniques to the detection of the m signal-elements, leads to the detection process B, shown in Fig. 4. This is in principle (but not in practice) the simplest of the four detection processes to be described here.

The output signal-vector from the m correlation detectors during a detection process is $E = (e_i)$, and the response of the integrators to the vector E is such that

$$\dot{X} = kE \quad \text{.....(16)}$$

where k is a positive constant. Both E and X vary continuously during a detection process.

At the start of a detection process the vector X is set to zero and the received vector R is fed to the input. Thus

$$E = RY^T = D \quad \text{.....(17)}$$

X is now permitted to vary freely, so that

$$E = (R - XY)Y^T = D - XA \quad \text{.....(18)}$$

and, from equation (16),

$$\dot{X} = k(D - XA) \quad \text{.....(19)}$$

It may be shown¹² that when A is real, symmetric and positive definite, the system is asymptotically stable and converges along the path of steepest descent with respect to the error function

$$(R - XY) \cdot (R - XY)^T \quad \text{.....(20)}$$

towards the single point of equilibrium, where

$$\dot{X} = 0 \quad \text{.....(21)}$$

and

$$D - XA = 0 \quad \text{.....(22)}$$

Thus at the end of the detection process,

$$E \approx 0 \quad \text{.....(23)}$$

and

$$X \approx DA^{-1} \quad \text{.....(24)}$$

Clearly the circuits associated with the correlation detectors in Fig. 4 perform the same function as the network A^{-1} in Fig. 3, so that the tolerance of the detection process B to the additive noise and signal distortion introduced in the channel should be the same as that of the detection process A.

3.2. Iterative Detection Processes

In an iterative process the solution vector X of the matrix equation

$$XA = D \quad \text{.....(25)}$$

is obtained as a result of a sequence of separate steps, giving successively closer approximations to the required solution vector.

A large number of different iterative processes are described in the published literature, but the majority of these require considerable equipment complexity and are not therefore suitable for our purposes.¹² There is however one iterative process which is ideally suited to the present application. This is the point Gauss-Seidel iterative process.^{20, 21, 23, 24}

The method of operation of this process and of two further developments of the process will first be described with reference to Fig. 5. In Section 5 it is shown how the practical implementation of these detection processes may be modified to give a much simpler arrangement than the detection process B.

3.3. Detection Process C

At the start of the detection process, the vector X in Fig. 5 is set to zero and the received vector R is fed to the input, so that

$$X = 0 \quad \text{.....(26)}$$

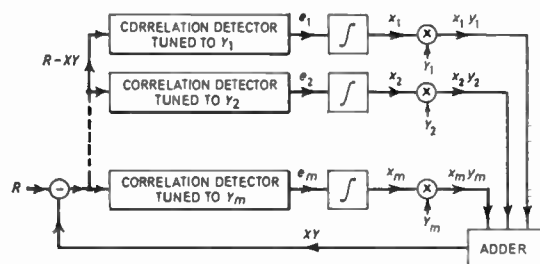


Fig. 4. Detection process B.

(Inputs to adder are $x_1 Y_1, x_2 Y_2, \dots, x_m Y_m$)

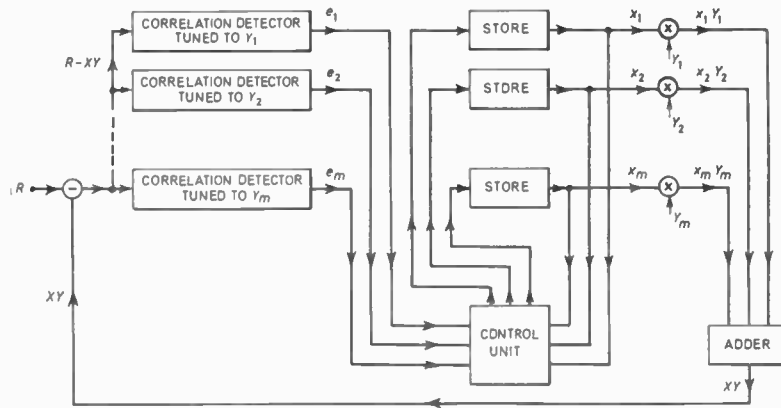


Fig. 5. Detection processes C, D and E.

and

$$E = D \quad \dots\dots(27)$$

x_1 is then adjusted so that the output signal from the first correlation detector is reduced to zero. This in general changes all m output signals $\{e_i\}$ from the correlation detectors. x_2 is now changed so that the output signal from the second correlation detector is reduced to zero, and so on sequentially to x_m , which completes the first cycle of the iterative process. The procedure is then repeated for the second cycle, the $\{x_i\}$ being changed sequentially and in the same order as before, and so on for as many cycles as required.

When x_i is adjusted to reduce to zero the output signal from the i th correlation detector, the change in x_i is

$$\Delta x_i = \frac{e_i}{v} \quad \dots\dots(28)$$

where e_i is the output signal from the i th correlation detector immediately preceding the change and

$$v = Y_i Y_i^T = a_{ii} \quad \text{for all } i \quad \dots\dots(29)$$

Clearly, X and E , instead of varying continuously as in process B, now vary in steps.

The above process can be modified so that the change in x_i becomes

$$\Delta x_i = h \frac{e_i}{v} \quad \dots\dots(30)$$

where h is a constant and

$$0 < h < 2 \quad \dots\dots(31)$$

This is known as overrelaxation,²⁴ and equation (28) is of course a special case of (30).

It can be shown theoretically that the detection process will converge to the required solution-vector X , so long as the matrix A is real, symmetric and positive definite and $0 < h < 2$, that is so long

as the m signal-elements $\{z_i Y_i\}$ are linearly independent and $0 < h < 2$.^{12,24} Thus at the end of the detection process

$$E \approx 0 \quad \dots\dots(32)$$

and

$$X \approx DA^{-1} \quad \dots\dots(33)$$

just as for the detection process B. To obtain the maximum rate of convergence of the iterative process, h should normally have a value equal to or a little greater than 1.

3.4. Detection Process D

A reduction in the equipment complexity of process C can be achieved by using a fixed magnitude for the change in x_i instead of the value given by equation (30). This leads to the process D which operates as follows, using the arrangement of Fig. 5.

The sequence of operations is exactly as described for process C, except that Δx_i is no longer given by equation (30). If, immediately preceding the change Δx_i in the stored value of x_i ,

$$|e_i| < fv \quad \dots\dots(34)$$

where f is a positive constant such that $f \ll 1$,

then

$$\Delta x_i = 0 \quad \dots\dots(35)$$

If

$$|e_i| \geq fv \quad \dots\dots(36)$$

then

$$\Delta x_i = \pm c |z_i| = \pm c \quad \dots\dots(37)$$

where c is a positive constant such that $f \ll c \ll 1$. The sign taken for Δx_i is the same as that of e_i .

Since $|z_i| = 1$, for all i , v is the magnitude of the output signal from a correlation detector if only the corresponding signal-element is received. fv is a

threshold level with which the correlation detector output signal is compared, before deciding whether or not to make a change in the corresponding x_i . c is typically less than 0.1 and f is typically 0.01.

The process D will converge so long as sufficiently small changes in x_i are used. It has an appreciably lower rate of convergence than process C, particularly with unfavourable signals requiring very small changes in x_i . This disadvantage is largely overcome in the process E, which avoids the need for very small changes in x_i to ensure convergence.

3.5. Detection Process E

This is a simple modification of the process D. As in D the values of the $\{x_i\}$ are changed sequentially and in a fixed cycle.

The change Δx_i in the stored value of x_i is now carried out in two steps. The first of these is exactly as for the process D, so that

$$\Delta x_i = 0 \text{ or } \pm c \quad \dots\dots(38)$$

depending on the value of e_i immediately preceding this change. In the second step, the sign of x_i is determined and the signal

$$\pm b|z_i| = \pm b \quad \dots\dots(39)$$

is added to it, the sign chosen for the added signal being the same as that of x_i . When x_i is zero the sign is chosen at random. b is a positive constant such that $b \ll c$. Typical values of f , b and c are 0.01, 0.01 and 0.1 respectively. Thus, immediately after each non-zero or zero change in x_i , determined as for the process D, the magnitude of x_i is slightly increased by a fixed amount b .

3.6. The Constraint on X

The tolerance to additive noise of the detection processes B to E can be improved by applying the following constraint to the vector X . The value of x_i is constrained to satisfy

$$|x_i| \leq 1 \text{ for all } i \quad \dots\dots(40)$$

throughout the detection process, where

$$|z_i| = 1 \text{ for all } i \quad \dots\dots(41)$$

In an iterative process the constraint overrides and so, if necessary, truncates the change in x_i dictated by the process.

The detection processes B to E, with and without the constraint on X , have been tested by means of computer simulation, for values of m up to 10 and for a general class of received signals probably less favourable to the detection process than the signals likely to be obtained here. In each trial of a computer simulation test, the m signal-elements $\{z_i Y_i\}$ were selected at random from the permitted signals, these being such as to ensure the linear independence of the

signal elements, while allowing the cross-correlation coefficient between any two of these to have one of the values

$$-0.9, -0.7, -0.5, \dots, +0.7, +0.9.$$

Each detection process, with values of m up to 6, was tested with and without additive Gaussian noise and with independent random variations in the individual levels of the m signal-elements over various ranges from 0 to about 10 dB.^{1,2}

The results of the tests show that the constraint on X , when correctly applied, does not prevent the correct convergence of any of the four detection processes. Furthermore, each of the detection processes B, C and D gains an advantage in tolerance to additive Gaussian noise of typically 1 or 2 dB when the constraint is applied. All processes with the constraint on X have the same tolerance to additive Gaussian noise and the same tolerance to inaccuracies in the setting of the constraint. In addition, when the levels $\{|z_i|\}$ of 6 signal-elements $\{z_i Y_i\}$ vary independently and at random over a range approaching 0.7 to 1.3, but with the vector X constrained as in equation (40), an advantage is still obtained in the average tolerance to additive Gaussian noise, together with correct detection in the absence of noise. The results clearly suggest that the operation of a detection process with the constraint on X is not critically dependent on the correct setting of the constraint. Further tests suggest that when there is any uncertainty in this setting, it should be adjusted to the highest of the range of likely values, say l , so that

$$|x_i| \leq l \text{ for all } i \quad \dots\dots(42)$$

Under these conditions the constraint on X should never prevent the correct convergence of the detection process and should always give some improvement in tolerance to additive Gaussian noise.

For correct operation of the process E, the constraint on X must be applied. This is to counteract the built-in tendency for the magnitudes of the $\{x_i\}$ to increase.

The constraint on X slightly reduces the rate of convergence of process B, it does not noticeably affect the rate of convergence of process D, and it greatly increases the rate of convergence of process C, which is now maximum when $1.25 \leq h \leq 1.5$. Except where otherwise stated, it will be assumed that the constraint on X as given by equation (40) is applied to each of the detection processes B to E and that $1.25 \leq h \leq 1.5$ for the process C. Equation (41) is assumed to hold, as before.

3.7. Assessment of Detection Processes

The detection processes B to E have been shown to operate correctly so long as they use suitable parameter

values, and they achieve a better tolerance to additive Gaussian noise than the process A. They have, however, a more important advantage over the process A in that all the stored values used in a detection process, that is R and Y , can be derived directly and very simply from the received signal, which is not so in the process A. Thus, with a time-varying channel, where the detector must be held correctly matched to the channel during the transmission of data, they should involve appreciably less complex equipment than an equivalent arrangement of the process A.

The important advantage of processes C, D and E over B is that they may be simplified to use only one correlation detector instead of the m needed by process B. The simplified arrangement of these processes is described in Section 5.

The processes D and E are entirely digital in the sense that they require only 'yes/no' decisions in determining the changes to the $\{x_i\}$. They should therefore be less complex than C.

The process E, although slightly more complex than D, has better convergence properties. Computer simulation tests have shown that with $m = 10$, $f = 0.01$, $b = 0.01$ and $c = 0.1$, the process E will converge in less than 50 iterative cycles, even under very unfavourable conditions. This is only a little slower than the process C under equivalent conditions.¹²

4. Adaptive Detection Techniques

4.1. Preset Detection

In Figs. 4 and 5 the m vectors $\{Y_i\}$, used both in the correlation detectors and to multiply the $\{x_i\}$, are the appropriate n -component segments of a single stored vector

$$L = \underbrace{0 \dots 0}_{m-1} \underbrace{y_1 y_2 \dots y_{g+1}}_{g+1} \underbrace{0 \dots 0}_{m-1} \dots (43)$$

where the $g+1$ potentially non-zero components of this vector appear in each of the m vectors $\{Y_i\}$. Clearly the knowledge of any one of these vectors, say Y_1 , completely determines the rest.

In an arrangement of preset detection, the start of a transmission contains a training signal in which each group of m elements has

$$z_1 = 1 \dots (44)$$

and

$$z_i = 0 \text{ for } i = 2, \dots, m \dots (45)$$

so that each group contains just the signal-vector Y_1 . To reduce the effects of noise, many vectors Y_1 are transmitted consecutively and the average of the received vectors is used to give the stored value of L .

With a knowledge of L the detector of Fig. 4 or 5 can now detect any received vector R in the following data signal.

Where the impulse response of the channel does not vary over any one transmission, the value of L determined at the start may be used throughout the transmission, giving a simple and effective arrangement of preset detection.

With a time-varying channel, the vectors Y_1 of the training signal may be interspersed between the data signals at regular known intervals and the value of L periodically changed so that the $\{Y_i\}$ in the detector follow the changes in impulse response of the channel. This, however, appreciably reduces the data transmission rate. A more effective arrangement for a time-varying channel will now be described.

4.2. Estimation of Y_1 by means of an Iterative Process

It may be shown that

$$ZY = Y'Z' \dots (46)$$

where Y' is the row vector $(y_1, y_2, \dots, y_{g+1})$ and Z' is the $(g+1) \times n$ matrix of rank $g+1$, whose i th row is

$$\underbrace{0 \dots 0}_{i-1} \underbrace{z_1 z_2 \dots z_m}_m \underbrace{0 \dots 0}_{g-i+1}$$

Thus

$$R = Y'Z' + W \dots (47)$$

Clearly, if Y' is determined, the $m \times n$ matrix Y is completely defined.

If the m correlation detectors tuned to the $\{Y_i\}$ in Fig. 3 are replaced by $g+1$ correlation detectors, where the i th correlation detector is tuned to Z'_i , the i th row of Z' , then the $(g+1)$ -component output signal-vector from the correlation detectors is

$$C = R(Z')^T = Y'B + W(Z')^T \dots (48)$$

where

$$B = Z'(Z')^T \dots (49)$$

The $(g+1) \times (g+1)$ matrix B is real, symmetric and positive definite.

It can be shown that when the receiver knows the $\{z_i\}$ but has no prior knowledge of the $\{y_i\}$ or σ^2 , then the best estimate it can make of Y' is the row-vector U , where

$$UB = C \dots (50)$$

U has $g+1$ components which are estimates respectively of the first $g+1$ components of Y_1 . From equation (50),

$$U = CB^{-1} \dots (51)$$

so that U can be determined by feeding the vector C

through a linear network B^{-1} or preferably by performing the equivalent iterative detection process, using the $g + 1$ correlation detectors tuned to the $\{Z_i\}$. No constraints are applied here to the components of U . Clearly there is a close parallel between equations (51) and (33).

4.3. Adaptive Detection

After the detection process for the $\{z_i\}$, the receiver can be assumed to know the $\{z_i\}$ with only a small probability of error, so that it can set up the $g + 1$ correlation detectors tuned to the $\{Z_i\}$ with a high probability of these being correct. Since the vector R is held at the input of the detector, the receiver now has all the necessary information to determine U . Each detection process for the $\{z_i\}$ is therefore followed by a detection process to determine U and so to estimate Y_1 .

It is obviously desirable that the same type of detection process should be used for both the $\{z_i\}$ and U , so that the maximum quantity of equipment may be shared between the two processes, thus reducing the total equipment complexity.

The vector U , determined at the end of each detection process, is used to adjust the stored value of L (equation (43)), so that the $\{Y_i\}$ in the detector follow the variations in the impulse response of the channel. Under normal conditions the variation of the impulse response should be slow enough to enable the m th to n th components of L to be given by the corresponding components of the running average of U over a number of these vectors, so that the $\{Y_i\}$ in the detector are only slightly affected by the additive noise in R . Since the element error rate is normally less than say 1 in 10^4 , errors in the detection of the $\{z_i\}$ should not be important. The effect of such errors may be reduced by limiting the maximum change between successive values of each component of U , or even by neglecting any U which shows an excessive change. An excessive change in U may also be used to indicate a probable error in the detection of one or more of the $\{z_i\}$.

The arrangement just described enables any of the detection processes B to E to be used in a fully adaptive system. The training signal is still required at the start of a transmission. The $\{z_i\}$ are here given a set of non-zero values (1 or -1) which are known at the receiver, and Y_1 is estimated in each detection process without first having to detect the $\{z_i\}$. When the detector has determined Y_1 , it is ready to receive data and is held correctly matched to the slowly time-varying channel during the following transmission. Any sequence of binary element values may be transmitted in the data signal without adversely affecting the operation of the adaptive detector.

4.4. Estimation of Y_1 by means of a Feedback Shift-Register

An alternative detection process, which may be used to estimate Y_1 , is shown in Fig. 6. A square here represents an element of a shift register. When triggered, an element transfers the signal at its input to its output, in the direction shown. The signal may be in analogue or digital form. A circle marked z_i represents a switched inverter, whose output signal is z_i times the input signal, where z_i may be 1, -1 or 0. It is assumed that $g < m$.

Before the start of a detection process, all signal voltages and all $\{z_i\}$ in Fig. 6 are set to zero and the vector R is fed into the feedback shift-register, to the position shown. The detection process for the $\{z_i\}$ is now carried out on R , after which the $\{z_i\}$ in Fig. 6 are set to their detected values (1 or -1). The two shift-registers are then triggered simultaneously $g + 1$ times. It may be shown that if there is no noise in R and the $\{z_i\}$ have been correctly detected, the output shift-register will now hold the vector Y' , that is the first $g + 1$ components of Y_1 .

The arrangement of Fig. 6 provides a very simple means of estimating Y_1 from R . It does however tend to accentuate the effects of additive noise in R and can be seriously affected when there is an error in the setting of one or more of the $\{z_i\}$. It would not therefore in general operate as well as an equivalent iterative detection process, over a noisy time-varying channel.

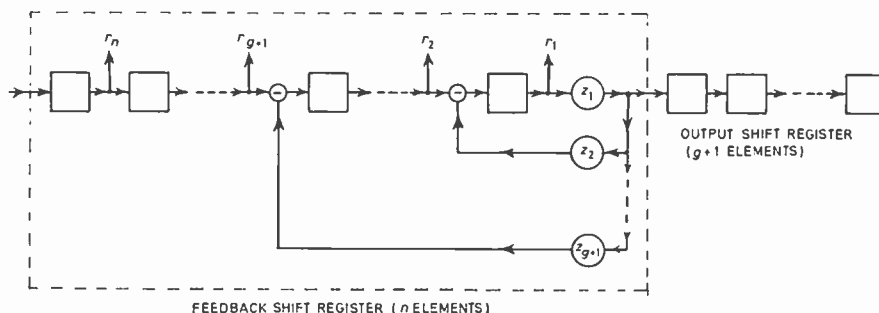


Fig. 6. Estimation of Y_1 from R .

4.5. Automatic Gain Control

The length of the vector Y_1 is held approximately constant at a nominal value of unity, as follows. The components of U , obtained after each detection process, are squared and added to give an estimate of the squared length of Y_1 . The running average of this estimate or of its square root, taken over several successive vectors U , is then used to control the gain of the a.g.c. amplifier at the input of the demodulator (Fig. 2). The latter is of course a part of the baseband channel.

Where the attenuation of the transmission path varies over a range of 40 or 50 dB, an efficient a.g.c. system should hold the length of Y_1 , and therefore the length of each Y_i , to within 1 or 2 dB of unity. This not only avoids the risk of overloading in the receiver but, by greatly reducing the variations in the channel impulse-response, it improves the operation of the adaptive detector. Under these conditions it is reasonable to assume (as has been done here) that the m vectors $\{Y_{ij}\}$, used in the detector, are the same as the corresponding vectors which make up the received vector R .

5. Simplified Adaptive Detector

5.1. Description of System

The simplified arrangement of the detection processes C, D and E is shown in Fig. 7. Each of the three shift registers stores a set of values in digital form, each value being represented by a binary number of 7 or more bits.

After the reception of the training signal, the vector Y_1 is stored in the shift-register F, as shown in Fig. 7. The detector is now ready to receive data and the vector R corresponding to the first group of m data-elements is fed into the shift-register G, to the position shown.

The $\{y_i\}$, for $i = g+2, \dots, n$, are always held at zero. The vector S , stored in the n summing circuits, is initially set to zero and so is the vector X , stored in the shift-register H. The output signal from the correction generator is also initially set to zero.

Each component of S is now subtracted from the corresponding component of R , and each component of the resultant vector $R-S$ is multiplied by the corresponding component of Y_1 . The n products are then added to give the correlation detector output signal

$$e_1 = (R-S)Y_1^T \quad \dots\dots(52)$$

where of course $S = 0$. The appropriate change Δx_1 in x_1 is now determined and held at the output of the correction generator. The value of Δx_1 is determined according to whichever of the detection processes C, D or E is being used.

Δx_1 is added to x_1 in the shift-register H, to give the signal $x_1 + \Delta x_1$ at the input to the following element. Δx_1 is also used to multiply each component of the vector Y_1 , and the product $\Delta x_1 Y_1$ is added to S . At this stage

$$x_i = 0 \quad \text{for all } i \quad \dots\dots(53)$$

and

$$S = \Delta x_1 Y_1 \quad \dots\dots(54)$$

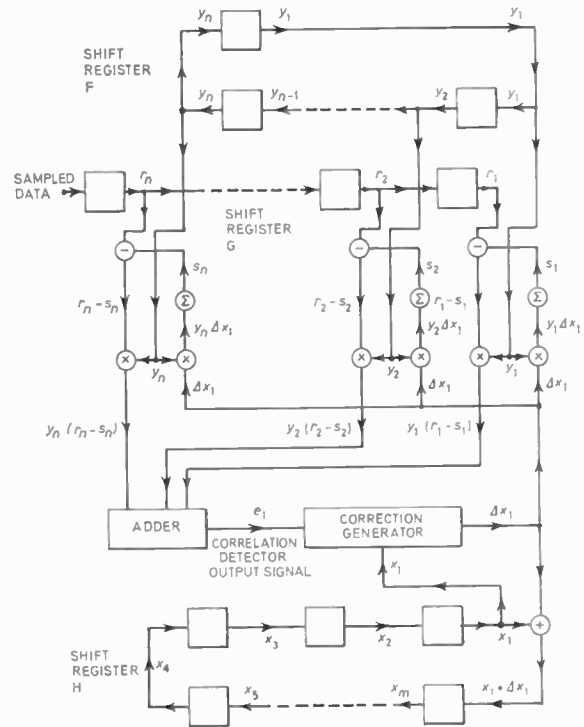


Fig. 7. Simplified iterative detector.

The shift-registers F and H are now triggered simultaneously so that each of the values stored in these shift registers is shifted one place in the direction shown. The shift-register F thus holds the vector Y_2 and the signal x_2 is fed to the correction generator. The output signal from the correlation detector now becomes

$$e_2 = (R-S)Y_2^T \quad \dots\dots(55)$$

where

$$S = \Delta x_1 Y_1 \quad \dots\dots(56)$$

The change Δx_2 in x_2 is determined and held at the output of the correction generator. Δx_2 is added to x_2 in the shift-register H and the vector $\Delta x_2 Y_2$ is added to S . At this stage

$$x_1 = \Delta x_1 \quad \dots\dots(57)$$

$$x_i = 0 \quad \text{for } i = 2, \dots, m \quad \dots\dots(58)$$

and

$$S = \Delta x_1 Y_1 + \Delta x_2 Y_2 \quad \dots\dots(59)$$

Clearly

$$S = XY \quad \dots\dots(60)$$

where the m components of X here are the input signals to the corresponding elements of the shift-register H.

The shift-registers F and H are now triggered again and the process continues exactly as described, until a change has been made to the stored value of each x_i and the first cycle of the iterative process has therefore been completed. At the start of the next iterative cycle, the signals stored in the shift-register H are automatically in the correct positions as in Fig. 7, but the signals in shift-register F require to be reset to the positions shown.

After the completion of a given number of iterative cycles, the receiver examines the signs of the $\{x_i\}$ stored in the shift-register H, and allocates the appropriate binary values to the corresponding signal-elements.

Having determined the $\{z_i\}$ the detector is now ready to estimate Y_1 . In the shift-register F, y_i is replaced by the detected value of z_i (1 or -1), for $i = 1, \dots, m$, and the remaining signals are set to zero. The shift-register G holds the vector R , as before. In the shift-register H, x_i is replaced by u_i , for $i = 1, \dots, g+1$, and the remaining signals are held at zero. It is assumed that $g < m$. At the start of this process the $\{u_i\}$ are set to the currently estimated values of the components of Y' and S to its corresponding value UZ' . The change Δu_i in u_i is determined either according to the detection process C, with $h = 1$, or else according to the detection process D, with $c = 0.01$ and $f = 0.005$. No constraints are applied to the $\{u_i\}$. With a slowly time-varying channel, where no u_i changes by more than say 5% from one detection process to the next, an adequate rate of convergence should be obtained with either arrangement.

Since the same piece of equipment is used first to detect the $\{z_i\}$ and then to estimate Y_1 , only a small increase in equipment complexity is involved in estimating the channel impulse response.

The arrangement of Fig. 7 uses only one correlation detector. It is therefore considerably less complex than the equivalent arrangement of Fig. 5. Where the circuits can handle an appreciable increase in speed of operation, the complexity of the system can again be greatly reduced by arranging the individual operations of multiplication and addition in the correlation detector to take place sequentially, instead of simultaneously as in Fig. 7. The resultant system should not be significantly more complex than the equivalent adaptive transversal equalizer.

5.2. Comparison with Transversal Equalizer

An adaptive detector overcomes the three main weaknesses of an adaptive transversal equalizer, mentioned in Section 1, and achieves the following advantages. Firstly, so long as adequate time-gaps are inserted between adjacent groups of signal elements, correct operation with near-optimum detection can be obtained over any channel likely to be used in practice, except of course when the channel introduces excessive attenuation. Secondly, a near-optimum detection process is used to estimate the channel impulse-response. Thus, the detector should be capable of following accurately any likely variation in the channel impulse-response which takes place over a time interval as short as $50nT$ seconds or typically $1000T$ seconds. This compares favourably with the equivalent adaptive transversal equalizer. Thirdly, the detector can adapt itself correctly to variations in the channel impulse-response, for any received sequence of data element values. Finally, only binary-coded elements need be used, even when the bandwidth of the baseband channel is appreciably less than $1/2T$ Hz,¹⁴ so that the equipment complexity involved with multi-level signals can be avoided. Where required, of course, multi-level elements may be used.

The adaptive detector has two main disadvantages. It is basically a little more complex than the adaptive transversal equalizer and it requires the signal elements to be separated into groups, which results in some reduction in tolerance to additive noise for a given transmission rate. However, this reduction is at least partly offset by the advantage in tolerance to additive noise gained through the constraint on X . A further disadvantage of the adaptive detector is the fact that when the channel introduces signal distortion, the received groups of elements will normally only become accurately orthogonal as $g \rightarrow \infty$. Thus the value of g should be kept as large as conveniently possible. If it is required to keep $g \ll n$ and so achieve the maximum tolerance to additive noise when the signal distortion is low, it may be necessary to accept significant intersymbol interference between adjacent groups of elements, when the signal distortion is high, with the result that the tolerance to additive noise will now be inferior to that obtained with a larger value of g .

The above considerations suggest that in some applications an adaptive detector should tolerate more severe and more rapidly time-varying signal-distortion than that tolerated by an adaptive transversal equalizer of similar complexity.

To avoid undue changes in the signal values stored in the shift-registers F and H in Fig. 7, during long iterative processes, all signals in the iterative detector must be stored in digital form. At the present time

this involves more complex equipment than if they were stored in analogue form, as in a typical transversal-equalizer. However, in a few years time it will probably be cheaper to handle the signals in digital rather than analogue form, so that the adaptive detector should then be at no disadvantage on account of its restriction to digital signals.

The adaptive detector appears likely to achieve the greatest advantage over the adaptive transversal equalizer in those applications where a received signal-element introduces severe intersymbol interference in the neighbouring elements but only over the two or three immediately adjacent elements on each side. It is unlikely to achieve any useful advantage in those applications where the intersymbol interference is spread over a large number of the neighbouring elements.

The adaptive detector does not, of course, equalize the channel but instead it accepts the distorted signals as they are and performs a near-optimum detection process on each group of m signal-elements, using the whole of each element in this process.

6. Conclusions

A relatively simple and most effective arrangement for an adaptive detector is obtained as follows.

The transmitted serial stream of binary data-elements is separated into groups, with adequate time-gaps between adjacent groups in order to eliminate intersymbol interference between the different groups in the received signal. Since the m received baseband data-elements in a group are linearly independent, they can be detected correctly in the absence of noise, for any channel impulse-response which does not introduce excessive attenuation.

The m elements are detected simultaneously in a single detection process. With a time-invariant channel this can be carried out by means of a set of m correlation detectors, each tuned to a different one of the received elements in a group, followed by a linear network which eliminates the intersymbol interference in the output signals from the correlation detectors. The system may be further simplified to only m correlation detectors which give the required output signals directly. With a time-varying channel, however, considerable equipment complexity is involved with either of these arrangements. The detection of the m signal-elements is now best carried out with the linear network in the original detector replaced by an iterative process which performs essentially the same function. The iterative process adjusts the estimates $\{x_i\}$ of the m element-values $\{z_i\}$, sequentially and in a fixed cycle, until no further significant changes in the $\{x_i\}$ are obtained. This not only leads to a much simpler system but it also

permits the introduction of constraints on the $\{x_i\}$, giving a useful increase in tolerance to additive noise.

Immediately following each iterative detection process, a second iterative process is used to obtain an estimate of the channel impulse-response and to adjust the detector appropriately, so that the detector follows the changes in the channel transmission-characteristics. No training signal is required, except at the start of transmission.

The arrangement promises to achieve in some applications a better performance over a time-varying channel, for a given degree of equipment complexity, than is possible with the more conventional transversal-filter adaptive equalizer.

7. Acknowledgments

Much of the work reported in this paper was carried out at the Imperial College of Science and Technology, University of London, as part of a course of research for the Ph.D. degree.¹² The work was sponsored jointly by the Science Research Council and Plessey Telecommunications Research Ltd., under the Industrial Fellowship Scheme and the author gratefully acknowledges their financial support. He would like to thank the Director of Research, Plessey Telecommunications Research Ltd., for permission to publish the paper.

8. References

8.1. Adaptive Equalization

1. Lucky, R.W., Salz, J. and Weldon, E.J., 'Principles of Data Communication', pp. 40-165 (McGraw-Hill, New York, 1968).
2. Rapoport, M. A., 'Automatic equalization of data transmission facility distortion using transversal equalizers', *Trans. Inst. Elect. Electronics Engrs on Communication Technology*, COM-12, No. 3, pp. 65-73, September 1964.
3. Becker, F. K., Holzman, L. N., Lucky, R. W. and Port, E., 'Automatic equalization for digital communication', *Proc. Inst. Elect. Electronics Engrs*, 53, No. 1, pp. 96-7, January 1965.
4. Lucky, R. W., 'Automatic equalization for digital communication', *Bell Syst. Tech. J.*, 44, No. 4, pp. 547-88, April 1965.
5. Gorog, E., 'A new approach to time-domain equalization', *I.B.M. J. Res. Devel.*, 9, No. 4, pp. 228-32, July 1965.
6. Lucky, R. W., 'Techniques for adaptive equalization of digital communication systems', *Bell Syst. Tech. J.*, 45, No. 2, pp. 255-86, February 1966.
7. Lucky, R. W. and Rudin, H. R., 'An automatic equalizer for general purpose communication channels', *Bell Syst. Tech. J.*, 46, No. 9, pp. 2179-208, November 1967.
8. Lytle, D. W., 'Convergence criteria for transversal equalizers', *Bell Syst. Tech. J.*, 47, No. 8, pp. 1775-800, October 1968.
9. Di Toro, M. J., 'Communication in time-frequency spread media using adaptive equalization', *Proc. I.E.E.E.*, 56, No. 10, pp. 1653-79, October 1968.

10. Gersho, A., 'Adaptive equalization of highly dispersive channels for data transmission', *Bell Syst. Tech. J.*, **48**, No. 1, pp. 55-70, January 1969.
 11. Proakis, J. G. and Miller, J. H., 'An adaptive receiver for digital signalling through channels with intersymbol interference', *Trans. I.E.E.E. on Information Theory*, IT-15, No. 4, pp. 484-97, July 1969.
- ### 8.2. Optimum Design and Detection of Signals
12. Clark, A. P., 'The transmission of digitally-coded-speech signals by means of random access discrete address systems', Ph.D. Thesis, Faculty of Engineering, University of London, approved July 1969.
 13. Zadeh, L. A. and Miller, K. S., 'Fundamental aspects of linear multiplexing', *Proc. Inst. Radio Engrs*, **40**, No. 9, pp. 1091-7, September 1952.
 14. Tufts, D. W., 'Nyquist's problem—the joint optimization of transmitter and receiver in pulse amplitude modulation', *Proc. I.E.E.E.*, **53**, No. 3, pp. 248-59, March 1965.
 15. Aaron, M. R. and Tufts, D. W., 'Intersymbol interference and error probability', *Trans. I.E.E.E. on Information Theory*, IT-12, No. 1, pp. 26-34, January 1966.
 16. Deighton, P. D., 'Optimisation of realisable receiver in pulse-amplitude modulation', *Electronics Letters*, **3**, No. 3, pp. 129-30, March 1967.
 17. Deighton, P. D., 'Joint optimisation of realisable receiver and transmitter in data-transmission systems', *Electronics Letters*, **3**, No. 7, pp. 342-4, July 1967.
 18. Shnidman, D. A., 'A generalized Nyquist criterion and an optimum linear receiver for a pulse modulation system', *Bell Syst. Tech. J.*, **46**, No. 9, pp. 2163-77, November 1967.
 19. Franks, L. E., 'Signal Theory', pp. 1-65 (Prentice-Hall, Englewood Cliffs, N.J., 1969).
- ### 8.3. Techniques for Solving Linear Simultaneous Equations
20. Beckenbach, E. F. (editor), 'Modern Mathematics for the Engineer', 1st series, pp. 428-79 (McGraw-Hill, New York, 1956).
 21. Faddeeva, V. N., 'Computational Methods of Linear Algebra', pp. 117-145 (Dover Publications, New York, 1959).
 22. Horn, R. E. and Honnel, P. M., 'Electronic network synthesis of linear algebraic matrix equations', *Trans. Amer. Inst. Elect. Engrs*, **78**, Part 1, pp. 1028-32, 1960. (*Communication and Electronics*, No. 46, January 1960.)
 23. Martin, D. W. and Tee, G. J., 'Iterative methods for linear equations with symmetric positive-definite matrix', *Computer J.*, **4**, No. 3, pp. 242-54, October 1961.
 24. Varga, R. S., 'Matrix Iterative Analysis' (Prentice-Hall, Englewood Cliffs, N.J., 1962).

Manuscript first received by the Institution on 29th December 1969 and in revised form on 29th May 1970. (Paper No. 1339/Com. 32)

© The Institution of Electronic and Radio Engineers, 1970

Conference on 'Laboratory Automation'

Middlesex Hospital Medical School, London W.1, 10th to 12th November, 1970

This Conference is organized by The Institution of Electronic and Radio Engineers with the association of The Institution of Electrical Engineers, The Institute of Physics and The Physical Society, The Royal Institute of Chemistry, and The Institute of Measurement and Control.

In recent years, automatic techniques have been introduced into many laboratories to speed up experimental and analytical procedures, and to reduce time spent by staff on repetitive work. On-line computers have been installed, facilitating the handling of large quantities of data, with immediate processing and presentation to the experimenter. In some cases, the computers have been used to control the actual experiments, for example the movement of a set of detectors to new, accurately controlled positions when sufficient data have been accumulated. Automatic equipment has been designed for carrying out chemical and biological analysis on a large number of

samples simultaneously, such as blood and urine in a pathological laboratory. Similar applications have occurred in other research laboratories.

The aim of the Conference will be to bring together workers who are already applying automatic techniques in their laboratories, or who may be interested in so doing, and designers and manufacturers of such equipment.

The term 'Laboratory' is intended to cover scientific observatories and routine testing laboratories as well as research and development laboratories. Papers are being contributed by workers in many fields, as will be seen from the provisional programme below.

Further information and registration forms for the Conference will be available in due course from the I.E.R.E., 9 Bedford Square, London, WC1B 3RG. Telephone: 01-637 2771 (Ext. 8).

PROVISIONAL PROGRAMME

Tuesday, 10th November

AUTOMATIC METHODS IN SPECTROMETRY (including electro-magnetic, nuclear and mass spectrometry and gas chromatography)

Data Processing at Harwell

Exploiting Small Computers for On-line Applications

The Use of On-line Computers for Instrument Control and Simple Data Processing

Mass Spectrometric Data Acquisition and Handling—
a Modular Low Cost High Performance System

An On-line Digital Computer for Enhancement and
Integration of Nuclear Magnetic Resonance Spectra

A Digital Control System for an Emission Spectrograph

A Computer Orientated Microdensitometer for
Spectrographic Plates

Laboratory Analysis with a Sealed-Tube Neutron
Generator

A Manual and Punched Tape Programmer for an X-Ray
Spectrometer

A Data Handling System for Use in Time-dependent
Fourier Transform Spectrometry

A Computer-Aided Gas Chromatography Laboratory
Continuous Gas Chromatography

Wednesday, 11th November

AUTOMATION IN CHEMICAL, BIO-CHEMICAL AND NUCLEAR LABORATORIES

Some Applications of Automation to Analytical Chemistry

Application of Ion-Selective Sensors to Automatic
Analysis

Autolab—A System for Automatic Chemical Analysis
Using Discrete Samples

Productivity in the Analytical Laboratory—A Rational
Approach to Automation

Interface Equipment for the Simultaneous Control of
Three Neutron Experiments by a Small Computer

Control and Monitoring of Neutron Beam Experiments
Using Data Processors

The Development of Automatic Analysis Equipment at
the Coal Research Establishment

Thursday, 12th November

AUTOMATION IN MECHANICAL, ELECTRONIC, ELECTRICAL AND OTHER LABORATORIES AND OBSERVATORIES

Use of a Process Control Computer in the Automation
of Materials Testing and Analytical Laboratories

Random Loading Fatigue Machines On-line to a
Digital Computer

The Use of a Laboratory Computer as Signal
Generator and Correlator

Automatic Temperature Monitoring in a Large Creep
Laboratory

Automated Underwater Ultrasonic Measurements

The Automatic Recording and Analysis of Magnetic
Fields

Low Level Language Programming for the On-line
Correction of Microwave Measurements

A Multi-Channel Analogue Recording System with
Computer Interfaced Data Processing Facilities

Automated Methods in Optical Astronomy

Computer Control of Optical Telescopes

Acquisition and Reduction of Acoustic Noise Data

Transmission Factors of Microwave Filters with Prescribed Attenuation and Group Delay

By

B. D. RAKOVICH,

Dip.Eng., Ph.D.†

and

A. D. JOVANOVICH,

Dip.Eng.†

Synthesis procedure is presented for determining transmission factors of a broad class of microwave filters using commensurate line lengths which satisfy prescribed attenuation and group delay characteristics. This technique, which lends itself to automatic programming, derives its origin from a previously described method of increasing selectivity in linear phase microwave filters and depends on the introduction of a second frequency transformation which is well known in the theory of lumped element filters. The organization of the computer program is discussed and a few numerical examples are worked out to facilitate the comparison of the results with those previously obtained.

1. Introduction

It is well established that Richards' frequency transformation¹ enables many problems of microwave networks using commensurate line lengths to be solved by methods developed for lumped element networks. Quite a number of studies have been reported involving 'low-pass' and 'band-pass' configuration which can be synthesized to yield Butterworth and Chebyshev type of magnitude response without need for new approximating methods. By contrast, although several papers have more recently appeared²⁻⁶, the corresponding problem of obtaining constant delay approximation in microwave networks has not been solved to the same extent. This may be attributed to the fact that, due to the nonlinear transformation, this technique cannot be directly applied to an approximation of the delay characteristic.

The procedure proposed by Carlin and Zysman² is based on the approximation of magnitude and not the delay characteristic about $\omega = 0$ so that the accuracy of delay response cannot be readily determined beforehand. In addition, the resulting magnitude response is not particularly suitable for filter applications. Scanlan and Rhodes⁴ have presented the method of synthesis of a class of microwave functions which provides a constant delay at all frequencies. However, these functions cannot be realized as interdigital or cascaded transmission line filters but require the use of C- and D-sections. Besides, the narrow-band filter cannot be constructed with these functions unless a large number of elements is used.⁶

An exact solution for the maximally-flat type of approximation of delay response in microwave filters has first been obtained by Abele³ from the use of finite convergents of the continued fraction expansion

† Faculty of Electrical Engineering, University of Belgrade, Yugoslavia.

for $\tan(\arctan \Omega)$. It has been demonstrated⁵ that the same result can be obtained by direct application of a method of synthesis of low-pass phase equalizers.⁷ These transmission factors can be realized as a transmission line or conventional interdigital network, but, since the resulting magnitude deviates considerably from that of an ideal filter they have limited use in the synthesis of band-pass filters.

Subsequently, an interesting method has been developed by Rhodes⁶ for simultaneous approximation of both group delay and magnitude responses of the filter in a maximally-flat sense about band centre. In order to obtain a maximally-flat delay the poles of the transmission factor are chosen to be the zeros of the symmetrical Jacobi polynomials and then the zeros of an odd polynomial in the numerator of the transmission factor are selected so that the maximally-flat magnitude criteria are met. This procedure, which is reminiscent of the technique used by Allemandou⁸ in lumped element network synthesis, yields an improved passband magnitude response when compared to that of the Abele's solution. Unfortunately, the stopband attenuation is deteriorated so that the overall magnitude performance is rather unsatisfactory. The realization of these transmission factors in the form of a generalized interdigital networks has also been given.^{6,9}

More recently, a method has been presented¹⁰ to find transmission factors of a general class of transmission line filters which yield an increased stopband attenuation of the filter while still retaining an excellent delay characteristic in the passband. These transmission factors are derived from the maximally-flat type of approximation of the delay response by reducing the number of flatness conditions by one and using the remaining parameter to adjust the magnitude response.

As with the previous paper¹⁰, the main purpose of the work described herein is to present a new method of determining transmission factors of linear phase transmission line filters which leads to further considerable improvement of the stopband performance of the magnitude response. On the other side, any particular set of passband specifications can be met so that these transmission factors may be regarded as approximating a constant delay and a constant magnitude simultaneously. They are capable of being realized by a transmission line or conventional interdigital network and are equally suitable for narrow-band and wide-band filters. Although the following discussion will be focused for the most part on band-pass filter functions, this method can also be employed in the synthesis of low-pass configurations.

The paper is organized as follows. In the first part the new transmission factor is derived from the results obtained for the maximally-flat type of approximation of delay response using a frequency transformation which is well-known in the synthesis of lumped element networks. Two variable parameters are then introduced and their influence on the magnitude and group delay characteristics is discussed. The second part is mainly concerned with the organization of a digital computer program for automatic computation of a transmission factor of minimum complexity that will satisfy the prescribed filter specifications. Numerical examples are worked out and the results are compared with those for other known methods.

2. Theoretical Background

2.1. Approximation Problem

Networks considered here consist of a finite number of series and shunt transmission line lossless elements having the same phase constant $\theta = \tau_0\omega$ and the characteristic impedances which may be different but are neither zero nor infinite. The generator and load impedances, between which the network is operated, are both positive real and need not be of equal value.

Let

$$S_{12}\left(j\frac{\pi\omega}{2\omega_0}\right) = S_{12}(j\tau_0\omega),$$

where τ_0 is the one-way delay of the line element and ω_0 is the radian frequency for which the line is a quarter-wavelength,

be a scattering transmission factor such that under the transformation

$$\Omega = \tan \tau_0\omega \quad \dots\dots(1)$$

or

$$\Omega = -\cot \tau_0\omega \quad \dots\dots(1')$$

the magnitude squared function is rational, then $|S_{12}(j\Omega)|^2$ must be of the form

$$|S_{12}(j\Omega)|^2 = \frac{\Omega^{2m}(1+\Omega^2)^r}{Q_n(\Omega^2)} \quad \dots\dots(2)$$

$$0 \leq |S_{12}(j\Omega)|^2 \leq 1 \quad \text{for} \quad -\infty \leq \Omega \leq +\infty \quad \dots\dots(3)$$

where $Q_n(\Omega^2)$ is a real polynomial of degree n in Ω^2 and $m = 0, 1, 2, \dots$ and $r = 0, 1, 2, \dots$ are constants depending on filter structure.

For band-pass filters the transformation (1') is pertinent and the function $S_{12}(p)$, as determined from (2), is

$$S_{12}(p) = \frac{p^m(p^2-1)^{r/2}}{H_n(p)} \quad \dots\dots(4)$$

where $p = \Sigma + j\Omega$ is the complex frequency variable in the transformed plane. If, on the other hand, the passband of the filter is located around $\tau_0\omega = 0, \pm\pi, \dots$ (low-pass filters) the transformation (1) is more appropriate and $S_{12}(p)$ becomes

$$S_{12}(p) = \frac{p^m(1-p^2)^{r/2}}{H_n(p)} \quad \dots\dots(5)$$

In the following analysis we shall assume $m = 0$ which corresponds to most practical cases under consideration. This does not represent a real restriction on the generality of the method since the numerator in $S_{12}(p)$, being an even or odd polynomial in p , does not contribute to the phase.

The approximation problem to be solved can now be stated in the following form: Subject to the constraint (3) determine an expression for $H_n(p)$ which must be a Hurwitz polynomial so that the scattering transmission factor (4), when evaluated in the original frequency plane $s = \sigma + j\omega$, approximates to the prescribed magnitude and group delay characteristics of the filter.

2.2. Determination of $H_n(p)$

It has been shown¹⁰ that starting with the ideal delay function¹¹

$$S_k(p) = \frac{1}{(1+p)^k} \quad \dots\dots(6)$$

where k is a positive real number and not necessarily an integer, the following recurrence formula can be derived for determining the polynomial $H_n(p)$ in the denominator of (5) corresponding to the maximally-flat type of delay approximation

$$H_n(p) = H_{n-1}(p) + \frac{k^2 p^2}{a_{n-2} a_{n-1}} H_{n-2}(p) \quad \dots\dots(7)$$

where $H_0(p) = 1, H_1(p) = 1 + kp$ and

i even:

$$a_i = \frac{\prod_{v=1}^{i/2} [k^2 - (2v-1)^2]}{\prod_{v=1}^{i/2} [k^2 - (2v)^2]} \quad \dots\dots(8)$$

odd:

$$a_i = \frac{(2i+1) \prod_{v=0}^{(i-1)/2} [k^2 - (2v)^2]}{\prod_{v=1}^{(i+1)/2} [k^2 - (2v-1)^2]} \dots\dots(9)$$

The realizability conditions restrict the lower limit of the parameter k to $k \geq n-1$. It has also been found¹⁰ that by substituting a new variable parameter ξ , lying in the range $n \leq \xi \leq a_{n-1}$, for a_{n-1} in (7) the selectivity of the magnitude response of the filter can be considerably increased while a very good delay characteristic is still retained.

These results will now be used to derive new transmission factors by means of some mapping properties of the function

$$p = \frac{\Omega_c}{2} \left(z - \frac{1}{z} \right) \dots\dots(10)$$

where Ω_c is the useful frequency band in the p domain. It is well known from the synthesis of lumped element networks that by this mapping function a maximally-flat type of approximation can be converted into a nearly equal-ripple type of approximation in the p plane.¹²

With this in mind we shall map the function $S_k(p)$ (equation (6)) in the z -plane retaining only those z_i which fall outside the unit circle

$$S_k(z) = \frac{1}{(z+z_0)^k} \dots\dots(11)$$

where

$$z_0 = \frac{1}{\Omega_c} + \sqrt{\frac{1}{\Omega_c^2} + 1}$$

or introducing a new variable $y = z/z_0$

$$S_k(y) = \frac{1}{z_0^k(1+y)^k} \dots\dots(12)$$

Since the function retains the same form after transformation, equations (7)–(9) can be directly applied to find the auxiliary approximating polynomial $H_n(\xi, z)$ in the z -plane. When this is accomplished, and the auxiliary polynomial solved for its roots, the mapping function (10) is used again to convert the results back into the p -plane. A similar technique was employed in the synthesis of lumped element networks leading to the so-called quasi-Chebyshev-type of delay response.^{13–16}

This procedure is straightforward but some important points must be taken into account before applying this method. Since, by means of the function (10), the left half of the z -plane outside the unit circle maps on the entire left half of the p -plane, the moduli of all

zeros of the auxiliary polynomial must be greater than unity; otherwise its p -plane transformation will not be a Hurwitz polynomial. Moreover, none of these zeros must be close to the unit circle if a satisfactory delay approximation is to be obtained in the original frequency plane. For any particular n and Ω_c (i.e., the fractional bandwidth of the filter $w = (f_2 - f_1)/f_0$) this imposes the upper limit of k since increasing k decreases the moduli of the z -plane zeros. For example, it has been found by numerical computation that for $n = 5$, $\xi = a_{n-1}$ and for large fractional bandwidth $w = 1$, corresponding to 3 : 1 bandwidth filter, one of the z -plane zeros lies inside the unit circle if $k = 9$ is chosen. On the other hand, for moderate values of the bandwidth factor, say $w = 0.5$ or less, higher values of k may be used.

The magnitude response of the filters is greatly affected by the values of k , which is very approximately equal to the normalized midband delay. As may be inferred from (6) the stopband performance is improved for higher values of k . However, increasing k increases the passband attenuation so that a higher-order network must be used than would otherwise be necessary from the point of view of the requirements imposed on delay response. This situation is similar to that encountered with the maximally-flat type of approximation, but simple means of improving the passband magnitude response, particularly in wideband cases have been found.

Suppose we know the required values of n and k and let for the moment $\xi = a_{n-1}$, corresponding to the maximally flat type of approximation in the z -plane. If all z -plane zeros of $H_n(\xi, z)$ are multiplied by a positive real factor $\lambda > 1$, the type of approximation remains unchanged but the passband and stopband attenuation are both decreased because of the frequency scaling. The zero frequency delay in the p -plane is decreased by the same factor and the general shape of the delay characteristic of the filter tends towards that of the maximally-flat type. If, on the other hand, the z -plane zeros and the parameter k are multiplied by λ simultaneously, the stopband attenuation increases very slowly, while the passband attenuation is decreased at the expense of an increased delay distortion. These results are illustrated in Figs. 1–3 for the third, fourth and fifth order approximating functions with $r = n-1$, $\Omega_c = 0.414$ (the fractional bandwidth $w = 0.5$) and the stopband beginning at $\Omega_s = 1$. As may be seen from Fig. 1, the passband attenuation decreases very slowly for higher values of the product $k\lambda$ so that little can be gained from increasing $k\lambda$ above, say, $(n+1)^2$, where n is the order of the network. Since in the design of the filters with small fractional bandwidth, $w = (f_2 - f_1)/f_0 = 0.1$ or less, a large value of k may be required to meet the stopband specifications,

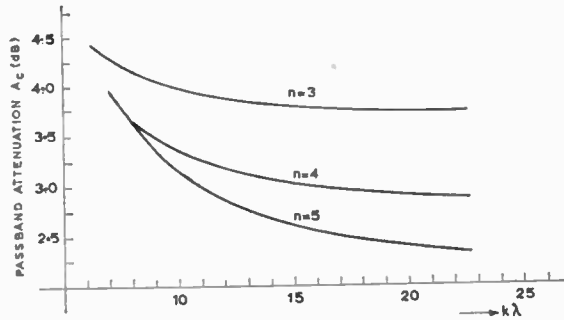


Fig. 1. Variation of passband attenuation with $k\lambda$ for $n = 3$ to 5 , $w = 0.5$, $r = n - 1$, $k = n + 1$.

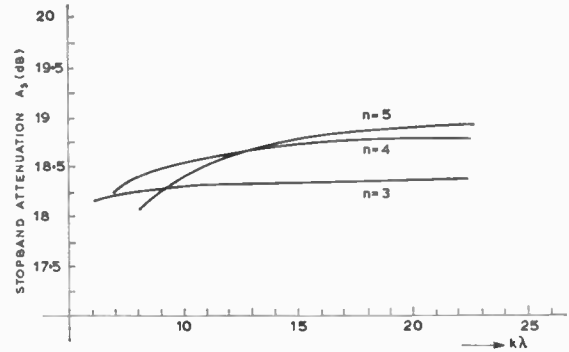


Fig. 2. Variation of stopband attenuation with $k\lambda$ for $n = 3$ to 5 , $w = 0.5$, $r = n - 1$, $k = n + 1$.

this technique becomes less efficient in these cases and, hence, $\lambda = 1$ should be chosen. It is fortunate, however, that in narrow-band filters the usual passband specifications can be easily fulfilled in almost all practical situations due to the fact that the useful band is narrow.

The diagrams shown in Figs. 1-3 have been obtained first by converting the z -plane zeros of $H_n(\xi, z)$ into the p -plane to form the approximating polynomial

$$h_n(p) = \sum_{i=0}^n B_i p^i \quad \dots\dots(13)$$

which has been normalized so that

$$h_n(0) = 1 \quad \dots\dots(14)$$

Then,

$$S_{12}(p) = \frac{(p^2 - 1)^{r/2}}{h_n(p)} \quad \dots\dots(15)$$

from which it follows for $r = n - 1$

$$|S_{12}(j\Omega)|^2 = \frac{(1 + \Omega^2)^{n-1}}{\sum_{i=0}^n C_i \Omega^{2i}} \quad \dots\dots(16)$$

where

$$\sum_{i=0}^n C_i \Omega^{2i} = h_n(p)h_n(-p)|_{p=j\Omega}$$

Denoting by $\tau(\omega)$ the group delay responses in the original s -plane, we find

$$\tau(\omega) = -\tau_0 \frac{d\Omega}{d\omega} \frac{d}{d\Omega} [\arg S_{12}(j\Omega)] \quad \dots\dots(17)$$

from which the normalized delay $T_n(\Omega)$ is obtained

$$T_n(\Omega) = (1 + \Omega^2) \frac{d}{d\Omega} [\arg S_{12}(j\Omega)] \quad \dots\dots(18)$$

The curves on Figs. 1-3 have been plotted by assuming $\xi = a_{n-1}$ but the same general relationships are preserved for other values of ξ .

2.3. Further Restrictions on ξ

The parameter ξ plays an important role in adjusting the shape of the magnitude response of the filter. Suppose first that the maximum number of flatness conditions is imposed on the delay response in the z -plane, i.e., $\xi = a_{n-1}$, leading to a quasi-Chebyshev type of delay approximation in the p -plane. In this case the magnitude response of the filter is found to be a monotonic function of frequency provided that $k \geq n - 1$ for $r \leq n$ and that all zeros of $H_n(\xi, z)$ are in the left half or the z -plane outside the unit circle. The magnitude performance is superior when compared with that of the maximally-flat type of delay approximation studied by Abele³ but the general shape of the magnitude response remains the same. Since $|S_{12}(j\Omega)|_{\max} = S_{12}(0)$ the realizability condition (3) is automatically fulfilled.

Now, if, for any particular n , ξ is decreased below the value $\xi = a_{n-1}$, while k and λ remain unchanged, the monotonic character of the magnitude response tends to be lost and a minimum appears in the attenuation characteristic just beyond the useful

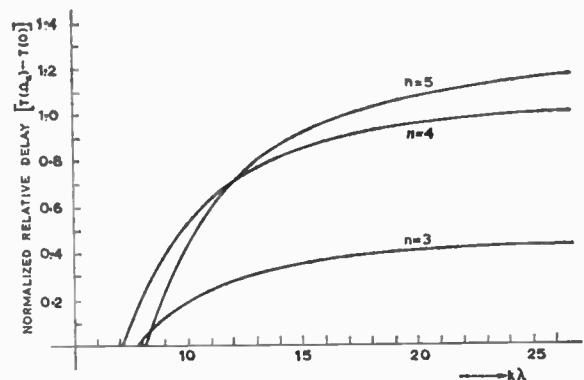


Fig. 3. Variation of normalized relative delay with $k\lambda$ for $n = 3$ to 5 , $w = 0.5$, $r = n - 1$, $k = n + 1$.

frequency band. The negative peak corresponding to this minimum increases and becomes sharper with decreasing ξ , as shown in Fig. 4. Since $S_{12}(p)$ must be bounded by unity on $\text{Re}(p) = 0$, the transmission factor must be multiplied by a constant $C < 1$, so that we have

$$S_{12}(p) = \frac{C(p^2 - 1)^{r/2}}{h_n(p)} \quad \dots\dots(19)$$

In order to find the upper limit of the constant C the frequency at which the attenuation minimum occurs should first be determined. Equating to zero the derivative $d/d\Omega |S_{12}(j\Omega)|$ we get after some manipulation

$$\sum_{i=0}^n [(r-i)C_i - (i+1)C_{i+1}] \Omega^{2i} = 0 \quad \dots\dots(20)$$

where, as before,

$$\sum_{i=0}^n C_i \Omega^{2i} = h_n(p)h(-p)|_{p=j\Omega}$$

Positive real zeros of (20) represent the frequencies at which analytical maxima and minima of the attenuation response occur. These zeros may be simple or multiple and the largest one always corresponds to a minimum. Moreover, it has been found that for $r \leq n-1$, the equation (20) has only two positive real zeros, if any, corresponding to the maximum and minimum of the attenuation characteristic. On the other hand, for $r = n$ and lower values of ξ the number of positive real zeros of (20) may become larger than two which depend on the specific values of the other two variable parameters k and λ . But again, with increasing ξ the attenuation response takes the same shape as for $r < n$ and it seems that these types of attenuation responses (Fig. 4) provide the best compromise between the magnitude and group-delay characteristics in all practical cases.

The upper limit for C in (19) is equal to unity only if

$$\left| \frac{(p^2 - 1)^{r/2}}{h_n(p)} \right|_{p=j\Omega_{min}} \leq 1 \quad \dots\dots(21)$$

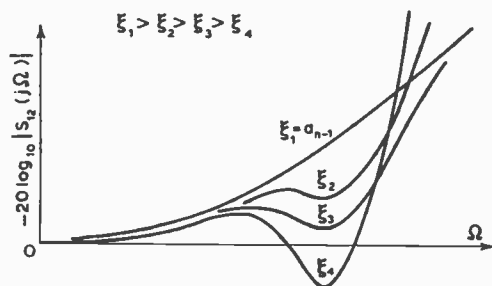


Fig. 4. Typical attenuation characteristics for different values of ξ .

In some cases, depending on filter specifications $\Omega_{max} < \Omega_c$, where Ω_c is the limit of the passband meaning that the maximum passband attenuation of the filter need not be at Ω_c . This point will be referred again in the following section where the organization of the computer program is described.

3. Computer Program

A computer program has been written enabling the automatic computation of the transmission factor described in this paper. The program has been organized to accept the following set of filter specifications:

- midband frequency f_0 in MHz
- fractional bandwidth $w = (f_2 - f_1)/f_0$
- maximum passband attenuation A_{ic} in dB
- stopband in MHz
- minimum stopband attenuation A_{is} in dB
- maximum passband delay distortion $\Delta\tau_0$ in ns
- values of the constants r and C (equation (19))
- starting values for n (the order of the approximating function), and k .

If for given filter specifications, the required order of the polynomial $H_n(\xi, p)$ cannot readily be estimated beforehand, the computation may start with $n = 3$. The realizability conditions imply that $k \geq n-1$, but in order to save the computational time a much larger initial value of k should be chosen, say $k = 3n - 4n$ if the fractional bandwidth of the filter is small ($w = 0.1$, or less).

A simplified flow chart is shown in Fig. 5. The main steps of the design procedure are as follows:

- (i) Using equation (1) compute

$$\tau_0 = \frac{1}{4f_0} \quad \dots\dots(22)$$

and the p -plane frequencies Ω_c and Ω_s corresponding to the limits of the passband and the stopband of the filter in the original plane respectively.

- (ii) Compute the coefficients a_1, a_2, \dots, a_{n-1} from (8) and (9). Since the proper sign for the increment $\Delta\xi$ is automatically adjusted during the computation the starting value of ξ may be conveniently chosen as $\xi = a_{n-1}/2$.

- (iii) Compute the coefficients of

$$H_n(\xi, y) = H_{n-1}(y) + \frac{k^2 y^2}{a_{n-2}\xi} H_{n-2}(y) \quad \dots\dots(23)$$

by using the recurrence relation (7), solve $H_n(\xi, y)$ for its roots $y_v = a_v + jb_v (v = 1, 2, \dots)$, compute $z_v = z_0 y_v = z_0(a_v + jb_v)$ and then transform the z -plane zeros into the p -plane to form the polynomial $h_n(p)$ equation (13).

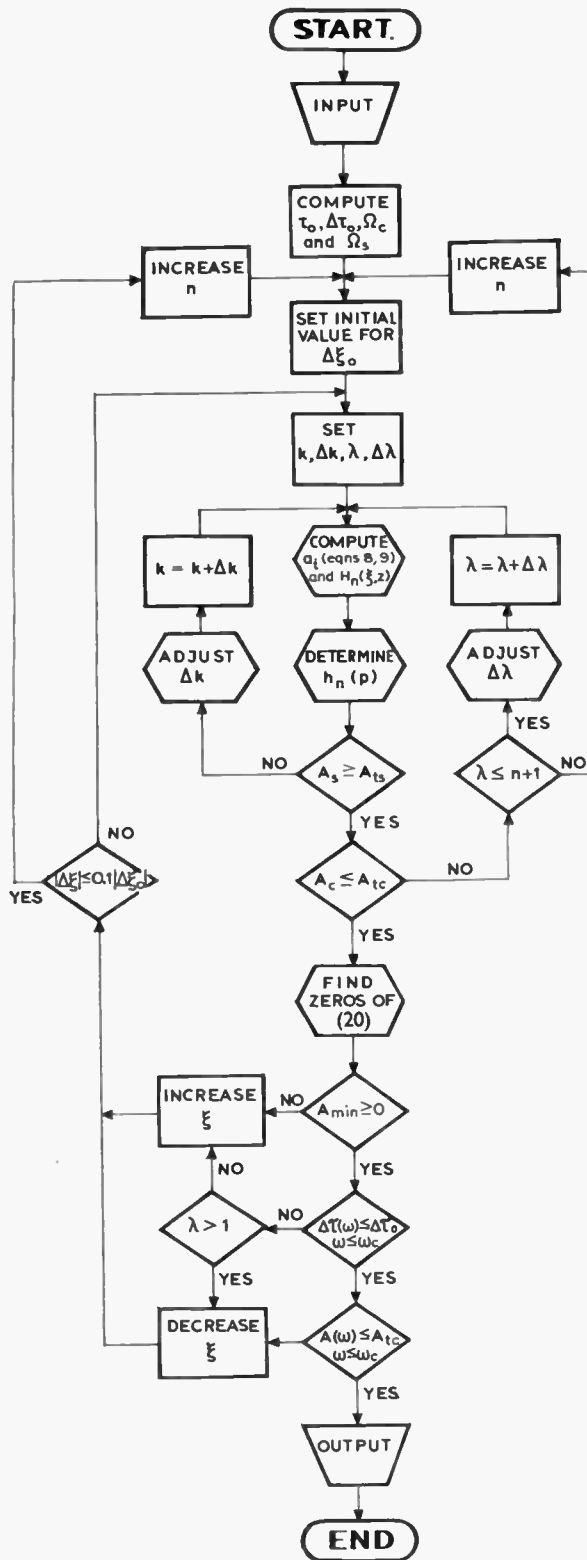


Fig. 5. Simplified flow chart for the program organization.

- (iv) Compute the minimum stopband attenuation $A_s = -20 \log_{10} |S_{12}(p)|$, where $S_{12}(p)$ is defined by (15). If for Ω_s , $A_s < A_{ts}$, increase k and repeat computation until the required attenuation A_{ts} is met.
- (v) Find the attenuation A_c at the limit of the passband. If $A_c > A_{tc}$ multiply k , obtained in step (iv), by $\lambda > 1$, evaluate the z -plane zeros and multiply them by the same factor λ . Convert the results into the p -plane and find the attenuation A_c at the limit of the passband. This operation is repeated until $A_c \leq A_{tc}$ or $k\lambda = (n+1)^2$ is reached. If $A_c > A_{tc}$ for $k\lambda = (n+1)^2$, decrease ξ and repeat computation.
- (vi) Solve (20) for its positive real roots Ω_{max} and Ω_{min} and check whether the condition (21) is fulfilled. If not, the instruction is given to increase ξ and repeat the computation. If for $\Omega = \Omega_{min}$ the attenuation is positive and $\Omega_{max} < \Omega_c$, meaning that the maximum passband attenuation is not at Ω_c , adjust k and λ in fine steps to meet the required specifications.
- (vii) For any n , the passband attenuation is minimum if the parameter ξ is adjusted so that the attenuation is equal to zero at $\Omega = \Omega_{min}$ (Fig. 4). Therefore, if the specified passband attenuation is not met in steps (v) and (vi), this value of ξ should first be determined before directions are given to increase the order of the function. This value of ξ can be obtained first by determining two values of ξ which correspond to positive and negative attenuations at $\Omega = \Omega_{min}$ respectively and then using a simple interpolation formula.
- (viii) Compute the delay response in the passband using (17) and, if necessary, decrease ξ and repeat computation to fulfil delay requirements. If the passband delay requirements are not met and the parameter ξ was adjusted in the preceding step to give zero attenuation at $\Omega = \Omega_{min}$, instructions are given to increase the order of the network and to repeat the computation. Otherwise, the value of ξ corresponding to zero attenuation at Ω_{min} should be determined and the passband delay response checked again before instructions are given to increase n .

In this connection one additional remark seems to be appropriate. It might seem, at first glance, that if the delay specification is not met in step (viii), instead of decreasing ξ , this parameter should be increased since with increasing ξ towards $\xi = a_{n-1}$, the auxiliary polynomial approaches the maximally-flat type of delay approximation in the z -plane. While it is generally true that for fixed k and λ , the passband delay distortion decreases with increasing ξ , the delay

response is much more affected by the parameter λ . For a specified maximum passband attenuation the necessary value of λ decreases with decreasing ξ and the passband delay characteristic is improved (see Fig. 4). Only in those instances where the passband attenuation requirements are met with $\lambda = 1$, the delay response is improved with increasing ξ . Hence, an instruction is introduced in step (viii) to change the sign of $\Delta\xi$ if $\lambda = 1$.

4. Numerical Examples and Comparison with Other Design Methods

The program just described was employed to find transmission factors of transmission line filters satisfying various specifications both for narrow-band and wide-band filters. It has been found that, for any prescribed maximum delay distortion and maximum attenuation in the useful band, these transmission factors provide a considerable improvement in the stopband performance of the filter when compared with the results obtained by any other method so far described. On the other hand, if the minimum stopband attenuation is also given, the transmission factor determined by the present method has a lower degree and hence the resulting filter is of smaller complexity than in any other comparable case. This will be illustrated by several examples which were

deliberately chosen to facilitate the comparison of the proposed procedure with the methods previously reported.

As a first example, the narrow-band case is considered and we propose to determine the transmission factor of an interdigital filter with short circuited input lines with the following specification:

midband frequency $f_0 = 1980$ MHz; passband 20 MHz (the fractional bandwidth factor $w \approx 0.01$); maximum passband attenuation 0.15 dB; minimum stopband attenuation of 45 dB at $f_0 \pm \Delta f_0 = 1980 \pm 140$ MHz; maximum passband delay distortion $\Delta\tau_0 = 0.2$ ns.

It is known^{4,17} that in order for the transmission factor to be realizable as the aforementioned network it must be of the form

$$S_{12}(p) = \frac{(p^2 - 1)^{(n-1)/2}}{h_n(p)} \dots\dots(24)$$

Hence, $r = (n-1)/2$ and we find that the third-order network with $k = 105.12$, $\lambda = 1$, $\xi = 2.34$, yields the maximum passband delay distortion $\Delta\tau_0 = 0.19$ ns, the maximum passband attenuation $A_c = 0.15$ dB and the minimum stopband attenuation $A_s = 47$ dB. All specifications are fulfilled and hence the third-order network can be chosen. The attenuation and group delay of the filter are shown in Fig. 6.

Now, using the method described by Rhodes⁶, the third-order function, realizable as a generalized interdigital network, with $\alpha = 150$ (α is the notation in Rhodes' paper for the normalized midband delay) satisfies the same passband requirements but provides only 20.7 dB attenuation at the limits of the stopband. For comparison the attenuation and delay characteristics of this filter are also shown in Fig. 6.

Another design procedure leading to the so-called transitional maximally-flat linear phase filter has also been described by Rhodes⁶. Again, a symmetrical Jacobi polynomial is used as the denominator in the transfer function in order to ensure a maximally-flat type of delay response. However, the number of constraints imposed upon the magnitude response in the passband is reduced to $q < (n-1)/2$ and $q < (n-2)/2$ for n odd and even respectively. The additional parameters are used to provide transmission zeros at the point $p = \pm 1$. In this way the stopband performance of the filter is improved at the expense of the quality of the passband magnitude response. This form of transfer functions can be realized by combining conventional and generalized interdigital sections as discussed by Rhodes.

In accordance with the second example from Rhodes' paper suppose the following design specifications are to be fulfilled:

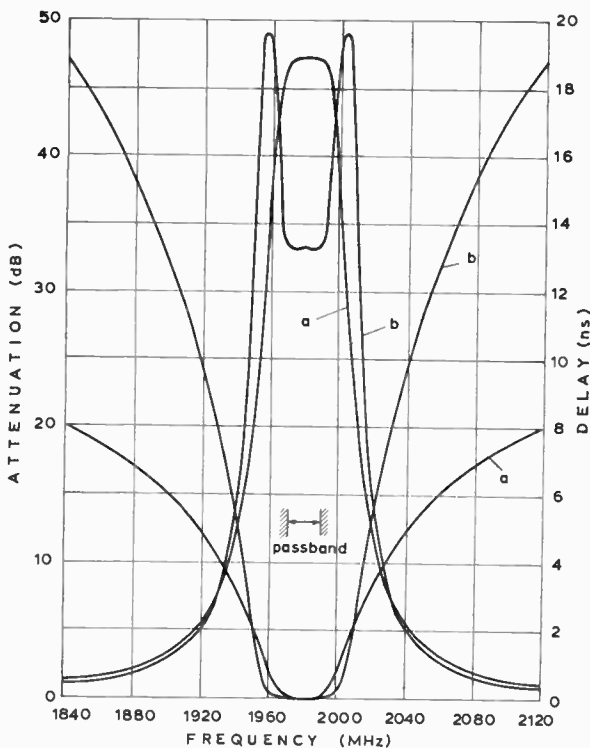


Fig. 6. Attenuation and group delay characteristics of the third-order narrow-band filters: (a) Rhodes' solution ($\alpha = 150$); (b) described system ($k = 105.12$, $\lambda = 1$, $\xi = 2.34$).

fractional bandwidth $w = 0.11$; maximum passband attenuation 2 dB, minimum stopband attenuation of 45 dB at $f_0 \pm 0.155 f_0$; maximum passband delay distortion $\Delta\tau_0 = 0.15$ ns.

Using the proposed method and choosing the transmission factor of the form (24) we find that the sixth-order network with $k = 39.86$, $\lambda = 8$ and $\xi = 3.25$ provides maximum passband attenuation 1.9 dB, minimum stopband attenuation 45.3 dB and maximum passband delay distortion 0.11 ns. Thus, $n = 6$ is necessary to meet all requirements. This represents a significant improvement over the twelfth-order transitional maximally-flat linear phase filter which is required to fulfil the same specifications.⁶ The latter was designed to have two transmission zeros at $p = \pm 1$, three flatness constraints on the passband magnitude response and the normalized midband delay $\alpha = 100$.

The wide-band case will be illustrated by the following example from Abele's paper. The transmission factor of an interdigital filter with short-circuited input lines is to be determined with the following specifications:

midband frequency $f_0 = 1000$ MHz; passband 850–1150 MHz (the fractional bandwidth factor $w = 0.3$); maximum passband attenuation 1.1 dB; minimum stopband attenuation 15 dB; maximum passband delay distortion $\Delta\tau_0 = 0.1$ ns.

As has been shown by Abele, the fifth-order network of the maximally-flat type is needed to meet these specifications. If the method described in reference 10 is applied, the same specifications are fulfilled with the third-order network. On the other side, if the procedure described in the present paper is used all passband requirements are met with the third-order network ($k = 6.37$; $\lambda = 1$; $\xi = 3.30$) but the stopband attenuation is increased to 22 dB.

If the comparison is made on the basis of equal n , then the fifth-order network of the present type ($k = 11.47$, $\lambda = 1.21$, $\xi = 3.21$), satisfying all passband requirements yields the minimum stopband attenuation of 42 dB, which is to be compared with 29 dB and 15.3 dB for the method described in Reference 10 and

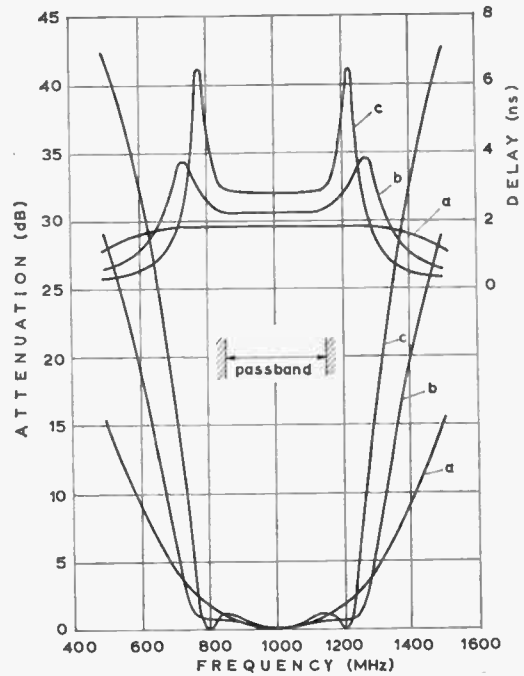


Fig. 7. Attenuation and group delay characteristics of the fifth-order wide-band filter ($w = 0.3$): (a) Abele's solution; (b) function described in Reference 10; (c) described system ($k = 11.47$; $\lambda = 1.21$; $\xi = 3.21$).

the maximally-flat type of approximation respectively (Fig. 7). The results of comparison are summarized in Table 1.

As a further illustration of the superiority of the proposed procedure, we shall assume that the useful band of the filter is increased to 500 MHz (750–1250 MHz) corresponding to the fractional bandwidth factor $w = 0.5$, while all other requirements are left unchanged as in the last example. Now, we find that the fifth-order network with $k = 7.63$, $\lambda = 1.09$, $\xi = 3.64$, yields the maximum passband attenuation 1.09 dB, the minimum stopband attenuation 20.2 dB and the maximum passband delay distortion $\Delta\tau_0 = 0.08$ ns so that all requirements are met. On the other hand, the fifth-order network determined by the method described in Reference 10 which satisfies the passband

Table 1. Comparison of bandpass filters with fractional bandwidth $w = 0.3$

	$n = 5$ Maximally flat	$n = 3$ Reference 10	$n = 5$ Reference 10	$n = 3$ Described system	$n = 5$ Described system
Maximum passband attenuation (dB)	1.15	0.49	0.71	1.07	1.07
Minimum stopband attenuation (dB)	15.29	15.04	29.03	22.08	42.46
Maximum passband delay distortions (ns)	0.000	0.097	0.057	0.092	0.086

delay specification provides the minimum stopband attenuation of 18.66 dB, but the maximum passband attenuation is increased to 3 dB. If, for the sake of comparison, the maximum passband attenuation that can be tolerated is increased to 3 dB, we find by applying the present technique that the transmission factor with $k = 9.04$, $\lambda = 1$ and $\xi = 4.23$ satisfies all passband specifications but provides a minimum stopband attenuation at $f_s = 1000 \pm 250$ MHz of 28.02 dB. This represents an improvement of almost 10 dB over the comparable figure for the method described in Reference 10, and an improvement of 15.5 dB when compared with the maximally-flat type of delay approximation. The results of this comparison are presented in Table 2.

As the last example of bandpass filters, we shall consider the case where more stringent requirements are imposed on the passband magnitude response. Suppose a bandpass filter is required with the following specifications:

midband frequency $f_0 = 2500$ MHz; passband 2000–3000 MHz (fractional bandwidth factor $w = 0.4$), maximum passband attenuation 0.15 dB; minimum stopband attenuation at 1000 and 4000 MHz 20 dB; maximum passband delay distortion $\Delta\tau_0 = 0.05$ ns.

Using (24) and retaining the same filter structure as before, i.e. $r = (n-1)/2$, we find that the third-order function satisfying all passband requirements cannot provide more than 15 dB attenuation at the limits of the stopband which is inadequate. On the other hand, the fourth-order solution with $k = 5.58$, $\lambda = 1.92$, $\xi = 3.80$ yields the maximum passband delay distortion $\Delta\tau_0 = 0.049$ ns, the maximum passband attenuation 0.14 dB and the minimum stopband attenuation 20.5 dB. Hence, all specifications are fulfilled and $n = 4$ can be chosen. The attenuation and delay characteristic of the filter are shown in Fig. 8. It is interesting to note that these specifications cannot be met either by the maximally-flat type of delay approximation or by the method described in Reference 10 even if $n = 12$ is chosen.

In the case of low-pass filter functions the frequency corresponding to the end of the useful range should be

substituted for f_0 . Suppose a low-pass filter which consists of a cascade of n lines is required with the following specifications²:

passband 0–900 MHz; maximum passband attenuation 3 dB; maximum passband delay distortion 0.015 ns; minimum attenuation of 30 dB at $f_s = 2400$ MHz; the end of the useful range 3000 MHz.

Now, we choose $f_0 = 2000$ MHz, the fractional bandwidth $w = (2 \times 900)/3000 = 0.6$ and $r = n$, since for lowpass filters with cascaded lines $S_{12}(p)$ must be of the form

$$S_{12}(p) = \frac{(1-p^2)^{n/2}}{h_n(p)} \quad \dots\dots(25)$$

It has been found that the above specifications are fulfilled with $n = 5$, $k = 7.79$, $\lambda = 1$, $\xi = 6.15$ while, if the Bernstein approximation discussed by Carlin and Zysman² is used, the number of cascaded lines must be $n > 10$ to meet the same specifications.

A more economical solution with regard to the total number of transmission line elements can be obtained if the filter is realized as a cascade of transmission lines with shunt open-circuited stubs which introduce a zero of transmission at the quarter-wavelength frequency. In this case the generic function of $S_{12}(p)$ must be of the form

$$S_{12}(p) = \frac{(1-p^2)^{(n-q)/2}}{h_n(p)} \quad \dots\dots(26)$$

where q is the order of transmission zero at the $\pi/2$ point. Now, using the same passband specifications as in the above example we find that a fourth-order network consisting of three cascaded lines and one stub ($q = 1$) with $k = 5.93$, $\lambda = 1$ and $\xi = 5.92$ yields minimum stopband attenuation $A_s = 31$ dB. This represents a marked improvement over 24 dB minimum stopband attenuation obtained with same filter structure when the method of determining the transmission factor described by Carlin and Zysman² is used.

Another important but clearly distinct aspect of filter synthesis is the realization of the transmission factor once the pole positions have been calculated.

Table 2. Comparison of bandpass filters with fractional bandwidth $w = 0.5$

	$n = 5$ Maximally flat	$n = 5$ Reference 10	$n = 5$ Described system	$n = 5$ Described system
Maximum passband attenuation (dB)	3.00	3.01	2.93	1.09
Minimum stopband attenuation (dB)	12.56	18.66	28.02	20.2
Maximum passband delay distortion (ns)	0.007	0.094	0.095	0.092

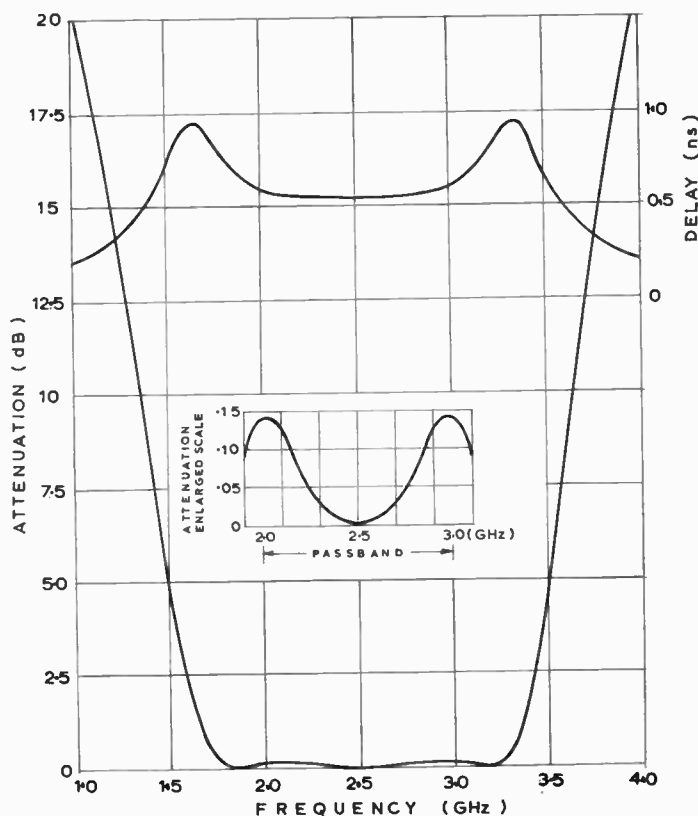


Fig. 8.
Attenuation and group delay characteristics of the fourth-order wide-band filter ($w = 0.4$).

The realization of these transmission factors in the form of interdigital networks and various other transmission line structures including cascaded transmission lines with or without shunt and/or series stubs, $\lambda/4$ transformers, etc. has been covered adequately in the literature together with the examples¹⁷⁻²³ and will not be considered in this paper.

5. Conclusion

A technique for determining transmission factors of a general class of microwave filters satisfying the prescribed attenuation and group delay characteristics has been presented. The method is based on the conformal mapping of a class of network functions, previously introduced by one of the present authors, by means of a function known from the theory of lumped element networks. A comparison of the results with those obtained by any previous method of which the authors are aware has revealed that the procedure reported here leads to a considerable improvement of the magnitude response of the filter while an excellent delay characteristic in the passband is still retained. Hence, using these new transmission factors, practical filters can be constructed in more compact form without need for separate delay

equalization which otherwise may require a rather complicated equalizer unit particularly if stringent requirements on phase linearity are prescribed.

The approximation problem, as an important part of network synthesis, is often sophisticated and based on rather complicated mathematical procedures with which the practising filter designers might not be well acquainted. Another advantage of the method described stems from the fact that it requires no formal training in approximation theory since its application has been shown in the development of a program for automatic computation of transmission factors of minimum complexity that will match the prescribed filter specifications.

6. Acknowledgments

The authors are much indebted to the Research Fund of S. R. Srbija for the financial support of work of which that described forms a part.

7. References

1. Richards, P. I., 'Resistor-transmission-line circuits', *Proc. Inst. Radio Engrs*, 36, pp. 217-20, February 1948.
2. Carlin, H. J. and Zysman, G. L., 'Linear phase transmission line networks', *Proc. Polytechnic Institute of Brooklyn Symp. on Generalized Networks* April 1966 pp. 193-226.

3. Abele, T. A., 'Transmission line filters approximating a constant delay in a maximally flat sense', *Trans. Inst. Elect. Electronics Engrs on Circuit Theory*, CT-14, pp. 297-306, September 1967.
4. Scanlan, S. O. and Rhodes, J. D., 'Microwave networks with constant delay', *Trans. I.E.E.E. on Circuit Theory*, CT-14, pp. 290-7, September 1967.
5. Rakovich, B. D., 'Microwave filters with maximally flat type of approximation to the constant group delay characteristic', Confer. of the Yugoslav Committee for ETAN, Rijeka (June 1968), (in Serbian).
6. Rhodes, J. D., 'The design and synthesis of a class of microwave bandpass linear phase filters', *Trans. I.E.E.E. on Microwave Theory and Techniques*, MTT-17, pp. 189-204, April 1969.
7. Rakovich, B. D. and Rabrenovich, D. M., 'Method of synthesis of phase correcting networks', *Proc. Instn Elect. Engrs*, 115, pp. 57-67, January 1968.
8. Allemandou, P., 'Low-pass filters—approximating—in modulus and phase—the exponential function', *Trans. I.E.E.E. on Circuit Theory*, CT-13, pp. 298-301, September 1966.
9. Rhodes, J. D., 'The theory of generalized interdigital networks', *Trans. I.E.E.E. on Circuit Theory*, CT-16, pp. 280-8, August 1969.
10. Rakovich, B. D., 'Linear phase transmission line filters with increased selectivity', *Trans. I.E.E.E. on Circuit Theory*, CT-17, pp. 41-5, February 1970.
11. Scentermai, G., 'The design of arithmetically symmetrical bandpass filters', *Trans. I.E.E.E. on Circuit Theory*, CT-10, pp. 367-75, September 1963.
12. Darlington, S., 'Network synthesis using Tchebycheff polynomial series', *Bell Syst. Tech. J.*, 31, pp. 613-66, July 1952.
13. Rakovich, B. D., 'Optimisation technique in delay-equaliser design by digital computer', *Electronic Letters*, 4, pp. 123-6, 5th April 1968.
14. Rakovich, B. D., 'Transfer functions approximating to a constant group delay': Part I, *Electronic Engng*, 40, pp. 242-6, May 1968; Part II, *Electronic Engng*, 40, pp. 326-8, June 1968.
15. Cartianu, G. and Constantin, I., 'Transfer functions obtained by a transform derived from the Darlington transform', *Electronic Letters*, 4, pp. 327-8, 9th August 1968.
16. Cartianu, G. and Constantin, I., 'Transfer functions with quasi Chebyshev-type characteristics obtained by a Darlington-derived transformation', *Electronic Letters*, 4, pp. 328-31, 9th August 1968.
17. Wenzel, R. J., 'Exact theory of interdigital bandpass filters and related coupled structures', *Trans. I.E.E.E. on Microwave Theory and Techniques*, MTT-13, pp. 559-75, September 1965.
18. Grayzel, A. I., 'A synthesis procedure for transmission line networks', *I.R.E. Trans. on Circuit Theory*, CT-5, pp. 172-81, September 1958.
19. Matthaei, G. L., 'Interdigital bandpass filters', *I.R.E. Trans. on Microwave Theory and Techniques*, MTT-13, pp. 479-491, November 1962.
20. Matthaei, G. L., Young, L. and Jones, E. M., 'Microwave filters', in 'Impedance Matching, Networks and Coupling Structures' (McGraw-Hill, New York, 1964).
21. Carlin, H. J. and Kohler, W., 'Direct synthesis of bandpass transmission line structures', *Trans. I.E.E.E. on Microwave Theory and Techniques*, MTT-13, pp. 283-97, May 1965.
22. Horton, M. C. and Wenzel, R., 'General theory and design of optimum quarter-wave TEM filters', *Trans. I.E.E.E. on Microwave Theory and Techniques*, MTT-13, pp. 316-27, May 1965.
23. Carlin, H. J., 'Synthesis of transmission line networks', Summer School on Circuit Theory (Prague 1968).
24. Steenaart, W., 'A contribution to the synthesis of distributed all-pass networks' (Appendix II), Proc. Polytechnic Institute of Brooklin Symp. on Generalized Networks (April 1966), pp. 173-191.

Manuscript first received by the Institution on 24th November 1969 and in final form on 29th May 1970. (Paper No. 1340/CC86).

© The Institution of Electronic and Radio Engineers, 1970

Standard L.F. Noise Sources using Digital Techniques and their Application to the Measurement of Noise Spectra

By

Professor H. SUTCLIFFE,

M.A., Ph.D., C.Eng., F.I.E.E.†

and

K. F. KNOTT, B.Eng., Ph.D.†

Reprinted from the Proceedings of the I.E.R.E. Conference on 'Digital Methods of Measurement' held at the University of Kent at Canterbury on 23rd to 25th July 1969.

Circuits of a random noise source and a pseudo-random noise source, both using digital techniques, are described and an account is given of their auto-correlation functions and power spectra. Their value as standard signal sources in noise power measurements is discussed and conclusions are reached about their relative merits.

1. Introduction

Within the general field of electronic instrumentation there is a trend towards the replacement of analogue circuits by their more precise digital versions, and the provision of precise noise sources is no exception. Especially at low frequencies the construction of a random noise source of known spectral intensity can be achieved more effectively by employing digital techniques than by attempting to exploit physical effects such as thermionic shot noise. In Section 2 of this paper a standard random low-frequency noise source based on digital techniques is described and in Section 3 a rather similar circuit employing pseudo-random sequences is discussed briefly. A comparison of the performances of these two circuits, together with some comments on experimental techniques concerning the measurement of noise spectra, appears in Section 4.

2. A Random Low-frequency Standard Noise Source

Of all the qualities needed of a standard noise source the accuracy of its spectral intensity is the most important. The distribution of amplitude is of little significance for measurements of noise spectra, since filtering action during the measurements has the effect of producing a normal distribution. This statement is justified by considering the response of a frequency selective filter in the time domain when the input function is wide-band in the frequency domain. The filter output at any instant may be regarded as the sum of a large number of independent effects, a situation typical of the normal distribution.

It is appropriate, therefore, to use as a basis for the noise source a random waveform which can be generated by precise processes, and the *random telegraph signal* is a natural choice. This waveform

† Department of Electrical Engineering, University of Salford, Salford M5 4WT.

was described in the early literature on random fluctuations, possibly because the ease of deducing its auto-correlation function makes it a good example for demonstrating the Wiener-Khinchine theorem.¹ This theorem is given in equations (1) and (2) in its practical form, that is for real sinusoids of positive frequency:

$$v(t)v(t-\tau) = R(\tau) \\ = \int_0^{\infty} G(f) \cos 2\pi f \tau \, df \quad \text{volts}^2 \quad \dots(1)$$

$$G(f) = \int_0^{\infty} 4R(\tau) \cos 2\pi f \tau \, d\tau \quad \text{volts}^2/\text{Hz} \quad \dots(2)$$

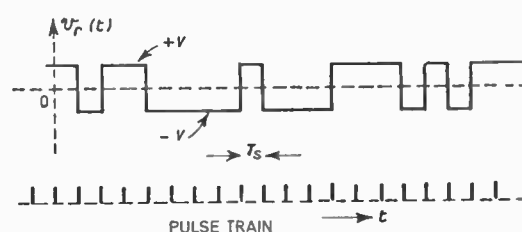


Fig. 1. Random telegraph signal.

The random telegraph wave $v_r(t)$ is shown in Fig. 1, together with a train of clock pulses with which it is associated. On the arrival of a clock pulse at intervals T_s , a random choice is imposed on $v_r(t)$, whether to acquire the value $+V$ or $-V$ for the duration of the next interval T_s . The autocorrelation function $R_r(\tau)$ of $v_r(t)$ is derived by simple reasoning and is shown in Fig. 2 together with the corresponding spectral intensity $G_r(f)$ derived from equation (2). For frequencies small compared with $f_s = 1/T_s$, the value of $G_r(f)$ is $2V^2T_s$. In this expression V is the amplitude of $v_r(t)$ and can be defined with great precision by simple circuits. T_s is also precise since it is simply the period of a regular pulse train. The

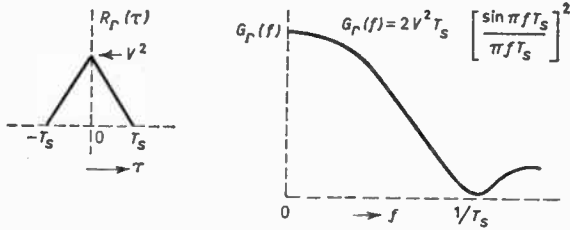


Fig. 2. Auto-correlation function of $v_r(t)$ and its spectral intensity.

remaining feature open to question in the system is the method of realizing the random choice $\pm V$ at the pulse instant. The method favoured by the authors is illustrated in Fig. 3.

In Fig. 3 the inclusion of bistable C ensures that v_3 has equal probability of '1' or '0'. The broadband h.f. waveform v_1 ensures that there is no correlation in waveform v_3 between adjacent pulse intervals. Thus if the minimum setting of T_s is 100 μ s, the bandwidth of v_1 is adequately large if it is in the region of 1 MHz. An exact analysis of this aspect of the design presents difficulties and would provide an interesting problem for theorists.

3. A Pseudo-random Noise Source

In recent years there has been an abundance of publications on the subject of pseudo-random binary signals generated as 'm-sequences' or maximum length sequences.² The conventional method of generating these is shown in Fig. 4.

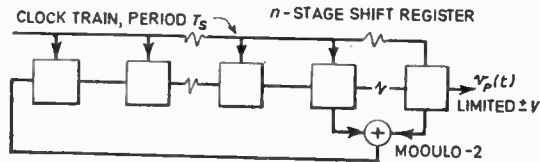


Fig. 4. Generation of p.r.b.s. as maximum length sequences.

Circuits of this type include the following features among their properties. The output waveform $v_p(t)$ often has a similar appearance, when viewed on an oscilloscope, to the random waveform $v_r(t)$ described in the previous Section. The difference is that waveform $v_p(t)$ is periodic and is repeated for every $(2^n - 1)$ intervals between clock pulses. This number is called the sequence length L , thus if the clock pulse interval is T_s , the fundamental period of $v_p(t)$ is $T_s \times L$ or $T_s(2^n - 1)$. The auto-correlation function of $v_p(t)$ is shown in Fig. 5, together with the power spectrum $G_p(f)$. It is interesting to note that if L is large, then

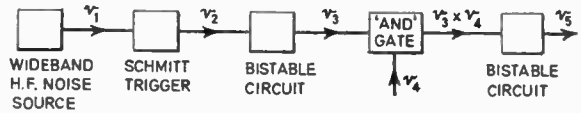


Fig. 3. Block diagram of circuit for generating random voltages.

- v_1 wideband h.f. noise
- v_2 h.f. random pulse train
- v_3 h.f. random square wave
- v_4 pulse train of interval T_s
- v_5 output (becomes $v_r(t)$ after clipping)

$G_p(f)$ approximates closely to a continuous power spectrum $S_p(f)$ such that:

$$S_p(f) = 2V^2T_s \left[\frac{\sin \pi f T_s}{\pi f T_s} \right]^2 \text{ volts}^2/\text{Hz}$$

This expression is similar to that of $G_r(f)$ in the previous section.

Instruments embodying these two types of noise source have been constructed and used extensively in noise spectrum investigations by the authors. Comments on the performance of the two types will be made in the next Section.

4. Application to Noise Spectrum Measurement

The first part of this paper described the different properties of true random and pseudo-random pulse trains. Let us now consider the different behaviour of the two types of waveform in l.f. noise spectra measurement. For a complete understanding of their application in this field it is perhaps useful to describe briefly a typical system used for measuring l.f. noise spectra. For the greatest accuracy it is preferable to compare the unknown source of noise with a calibrated noise generator. The block diagram of Fig. 6 illustrates such a system. The noise to be measured is amplified, filtered and detected. The change in

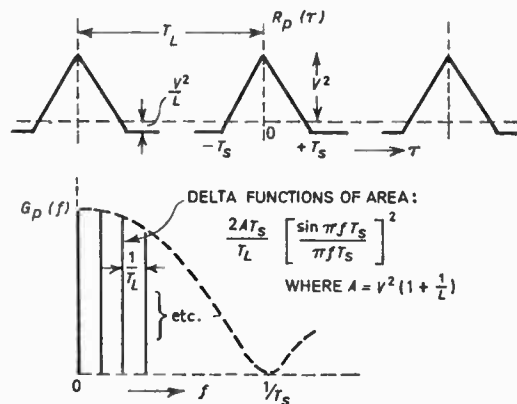


Fig. 5. Auto-correlation function of $v_p(t)$ and its power spectrum.

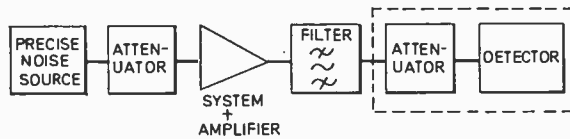


Fig. 6. Noise measurement system.

the detector reading when a known amount of additional noise is fed into the system then gives the value of the unknown noise. The detector may be a true r.m.s. meter but distinct advantages are gained if a transit-counting detector³ is used instead. The main advantages of this type of detector are that the period of integration may be set to an arbitrary length and that the reading is immune from drift in the datum level of the signal. However, this type of detector gives a reading which is proportional to the mean magnitude of the signal. For noise which has a Gaussian amplitude probability density function there is a constant relationship between mean magnitude and r.m.s. values. For almost all types of noise found in practice the amplitude probability density function approximates closely to Gaussian after band-pass filtering so that this is not a serious restriction. A further advantage of this type of detector is that the output is in digital form and hence the accuracy of the reading can be high. The factors governing the use of the two types of digital noise generators for narrowband noise spectra will now be considered.

If a true random generator is used for calibrating the system then the detector output must be integrated over a period of time which satisfies the usual standard deviation law

$$\sigma \ll 1$$

where

$$\sigma = \left(\frac{1}{BT} \right)^{\frac{1}{2}}$$

In this expression B is the bandwidth of signal, T is the time of observation and σ is the fractional standard deviation of the observed noise power. For example, if a measurement were being made to an accuracy of 10% over a bandwidth of 0.1 Hz centred at 1 Hz, the observation time would be 1000 seconds.

If a pseudo-random generator is used for calibrating the system there may be a saving in calibration time depending on the number of discrete lines required in the frequency spectrum of the system. The spacing of the lines depends on the sampling frequency and the length of the sequence of pulses. To avoid errors due to the shape of envelope of the power spectrum of the generator, the sampling frequency should be set to about 20 times the frequency of the noise measurement so that for a fixed

frequency of measurement the lower limit of sampling frequency is fixed. This, therefore, means that the time for one sequence is inversely proportional to the spacing of the lines in the power spectrum and since the calibration time need be only one sequence period it follows that it also is inversely proportional to the line spacing.

There will thus be a saving of time if the line spacing can be increased, but this raises the question of the effect on the apparent power spectral density of reducing the number of discrete lines. This effect will depend on the shape of the frequency response of band-pass filters. Consider the ideal rectangular band-pass response as shown in Fig. 7(a), then if the centre frequency is varied slightly a serious error could be introduced in calculating the power in the bandwidth. Whereas, if a filter with a less sharp cut-off characteristic as in Fig. 7(b) were used, there would be no abrupt changes in the power contained in the bandwidth. Figure 7(b) is perhaps more typical, since it is usual in l.f. noise spectra measurements to use only two single-tuned circuits in cascade for the filter.

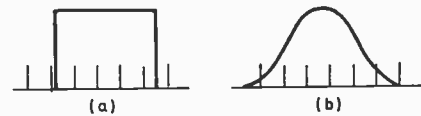


Fig. 7. Possible situations in bandpass systems.

A further question needs to be considered, namely the errors introduced by the departure of the amplitude probability density function from Gaussian when the number of lines in the spectrum is small. This will conceivably introduce errors if a transit counting or other mean magnitude method is used. It is also of interest to discover whether errors introduced in this manner are more serious than those that would be obtained from the effect described in the previous paragraph even if a true r.m.s. meter were used.

It has not yet been possible to find a satisfactory analytical solution to the problem. To obtain a measurement of the likely errors involved the following experiments were carried out using a true random noise generator and a pseudo-random generator of equal power spectral density as calculated from the sampling frequency and amplitude of the pulses. The two generators were compared by feeding them into a Bruel and Kjaer audio frequency analyser type 2107 and measuring the output with a transit counting detector. In the first experiment the setting of the analyser was kept constant at 500 Hz centre frequency and 21% bandwidth and the sampling frequency of

the noise generators was kept constant at 10 kHz. The sequence length of the pseudo-random generator was varied to give from 2.5 lines to 80 lines in this bandwidth. The reading of the detector was taken over one sequence period for the pseudo-random generator and for the random generator the reading was averaged over a length of time, such that the fractional standard deviation was less than 0.01. The readings agreed to within 2% for number of lines as low as 5.

In the second experiment the frequency of the analyser was varied gradually to $\pm 10\%$ of the initial value and also the bandwidth was varied from 21% to 12%. The response of the detector was plotted for the pseudo-random generator set to give 2.5 to 90 lines in the bandwidth. The graph of Fig. 8 illustrates the results obtained for 5 lines in 3 dB bandwidths of 21% and 12%. Also on this graph is the response of the detector for true random noise.

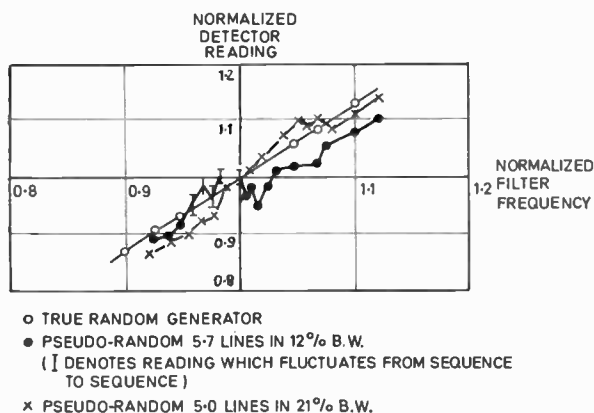


Fig. 8. Variation in detector reading for $\pm 10\%$ variation in filter frequency.

Considering the results for the pseudo-random generator it is seen that there are sharp fluctuations in the detector reading for quite small changes in frequency. Another thing to note is that when using narrow bandwidth and few lines one can encounter certain conditions under which the reading varies between sequences. This last result is rather unexpected since all input sequences are identical. A possible explanation of the result is that when there are only a few lines in the signal spectrum their phase relation is critical. If there were any slight changes in the phase response of the filter with time the reading might be affected.

Curves such as those in Fig. 8 were plotted for each value of bandwidth and for numbers of lines between 2.5 and 90. To obtain a comparison between the results the response of the detector to true random

noise was taken as a reference. The greatest deviation in the detector response encountered in the $\pm 10\%$ frequency range was then measured for each value of the number of lines. In Fig. 9 this deviation, expressed in noise power, is plotted as a function of the number of lines for the two values of bandwidth used. Also shown in this figure are experimental points obtained using a true r.m.s. meter.

It is seen from Fig. 9 that firstly there is little difference between the r.m.s. meter and the transit-counting detector for numbers of lines greater than 10. Secondly that for numbers of lines above 20 the deviation between the pseudo- and true-random generators is small. For a maximum deviation of 1% the number of lines required appears to be about 100.

For a 10% measurement 10 lines would be adequate. In this case the time taken for one sequence, assuming the bandwidth is 10% and the sampling frequency is 20 times the frequency of the measurement, is

$$T = 100/f_0$$

where

$$f_0 = \text{centre frequency}$$

For example, in a measurement over a bandwidth of 0.1 Hz the calibration time would be reduced from the previous 1000 s to 100 s.

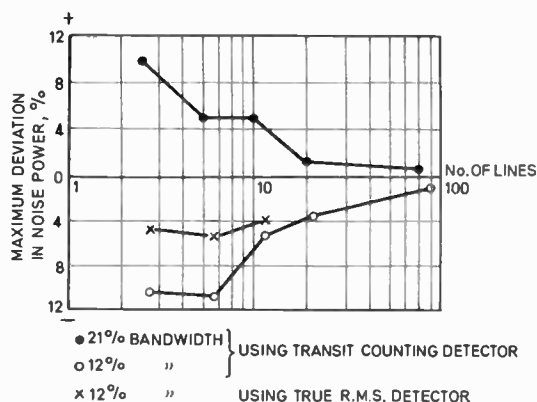


Fig. 9. Deviation between pseudo- and true-random generators for $\pm 10\%$ change in filter frequency.

Although the reduction in calibration time appears to be a big advantage for l.f. measurements there is one point to bear in mind if a pseudo-random generator is used. To obtain this reduction in calibration time the level of the pseudo-random signal fed into the system must be much larger than the inherent noise in it. This means that one must be careful to ensure that the signal at all points in the system is within the linear range of the constituent

parts of the system. This restriction does not apply when a true random signal is used since this can be fed into the system at the same level as the inherent noise.

In conclusion, the discussion has shown that the pseudo-random noise source has theoretical advantages, in that the time of calibration of systems may be reduced. This advantage is partly offset by the added complexity of the experimental procedures.

5. Acknowledgments

The authors would like to express their thanks to their colleague Mr. H. Dunderdale for his collaboration and advice on pseudo-random sources.

6. References

1. Rice, S. O., 'Mathematical analysis of random noise', *Bell Syst. Tech. J.*, 23, p. 282, 1944 and 24, p. 46, 1945.
2. Kramer, C., 'A low frequency pseudo-random noise generator', *Electronic Engineering*, 37, p. 465, July 1965.
3. Sutcliffe, H., 'Mean detector for slow fluctuations', *Electronics Letters*, 4, No. 6, p. 97, March 1968.

Manuscript first received by the Institution on 19th May 1969 and in final form on 14th July 1970. (Paper No. 1341/CC87).

© The Institution of Electronic and Radio Engineers, 1970

STANDARD FREQUENCY TRANSMISSIONS—August 1970
(Communication from the National Physical Laboratory)

August 1970	Deviation from nominal frequency in parts in 10 ¹¹ (24-hour mean centred on 0300 UT)			Relative phase readings in microseconds N.P.L.—Station (Readings at 1500 UT)		August 1970	Deviation from nominal frequency in parts in 10 ¹¹ (24-hour mean centred on 0300 UT)			Relative phase readings in microseconds N.P.L.—Station (Readings at 1500 UT)	
	GBR 16 kHz	MSF 60 kHz	Droitwich 200 kHz	*GBR 16 kHz	†MSF 60 kHz		GBR 16 kHz	MSF 60 kHz	Droitwich 200 kHz	*GBR 16 kHz	†MSF 60 kHz
1	-300.0	+0.1	+0.1	629	615.5	17	-300.1	-0.2	0	641	627.2
2	-300.0	0	+0.1	629	615.4	18	-299.9	+0.1	0	640	625.3
3	-300.0	0	+0.1	629	615.0	19	-300.0	+0.1	0	640	624.8
4	-300.0	0	+0.1	629	614.8	20	-300.0	-0.1	0	640	625.6
5	-300.4	0	+0.1	630	614.9	21	-300.0	0	0	639	625.8
6	-300.0	0	+0.1	630	615.0	22	-299.9	0	0	638	625.5
7	-300.0	0	+0.1	630	615.1	23	-299.9	0	+0.1	637	625.1
8	-300.0	-0.1	+0.1	630	616.1	24	-299.9	+0.1	+0.1	636	624.6
9	-300.1	-0.1	+0.1	631	617.1	25	-299.9	+0.1	+0.1	635	623.8
10	-300.1	-0.1	+0.1	632	617.8	26	-299.9	0	+0.1	634	623.6
11	-300.1	-0.1	+0.1	633	619.0	27	-300.1	+0.1	+0.1	635	623.0
12	-300.1	-0.1	+0.1	634	620.4	28	-300.1	0	+0.1	636	623.0
13	-300.1	-0.1	+0.1	635	621.1	29	-300.0	0	0	636	623.4
14	-300.2	-0.1	+0.1	637	622.0	30	-300.0	-0.1	+0.1	636	624.4
15	-300.1	-0.1	+0.1	638	623.4	31	-300.1	0	+0.1	637	624.5
16	-300.2	-0.2	0	640	625.4						

All measurements in terms of H.P. Caesium Standard No. 334, which agrees with the N.P.L. Caesium Standard to 1 part in 10¹¹.

* Relative to UTC Scale; (UTC_{NPL} - Station) = + 500 at 1500 UT 31st December 1968.

† Relative to AT Scale; (AT_{NPL} - Station) = + 468.6 at 1500 UT 31st December 1968.

CHANGES TO THE MSF MODULATION SCHEDULE

The following changes to the MSF modulation schedule will be made during September 1970.

1. The A2 modulation at present carried by the 60 kHz transmission between 1430 and 1530 GMT will be abandoned, and the A1 pulse modulation extended to 24 hours a day.
2. The MSF station identification will be emitted twice only at ten-minute intervals on the h.f. transmissions and during the five seconds preceding every hour on the 60 kHz transmission.

Pulse Counting and Encoding Systems used on a Rocket-borne Spectrophotometer

By

D. H. BEATTIE, Graduate†

and

C. H. PATERSON, B.Sc.‡

Problems encountered making photoelectric observations of stars from spinning rockets have been solved by using photon pulse counting. The counter states were suitably encoded so as to permit the transmission of all useful data within the limits of two analogue telemetry channels of restricted bandwidth. The logic system is described and resulting data shown.

1. Introduction

The middle and far ultra-violet radiation from stars, which is absorbed by the Earth's high atmosphere, is of considerable interest to astronomers. The Space Research Division of the Royal Observatory, Edinburgh, has therefore been engaged for some years in the construction and flight of rocket-borne instruments to make observations in this wavelength region. The instruments have been carried on *Skylark* rockets which, for technical reasons, are uncontrolled in their motion while above the atmosphere, except that a predetermined spin rate can be obtained by means of a gas jet unit. During the flight, telescopes looking sideways from the rocket therefore make a series of scans across the sky as the rocket spins and precesses, and stars can be observed photoelectrically as they drift across the fields of view.

The system described in this paper was devised as part of an instrument¹ in which the starlight collected by the primary telescope mirror during a scan is reflected off a plane diffraction grating, so that in the focal plane each star is imaged as a spectrum as well as a 'white light' point image (zero spectral order). When a photomultiplier is mounted behind a slit which is in this focal plane, the spectrum can be measured photoelectrically because the spin of the rocket causes these images to be swept across the slit. The entire section of spectrum to be observed (150 nm to 300 nm) takes about 85 ms to cross the slit and the faintness of even bright stars, combined with the requirement to observe a reasonable number in each flight, means that the spectral resolution provided by the slit must be quite low, some 20 nm in this case. Nevertheless, the sampling frequency needed to transmit, and to enable full reconstruction of the spectrum in wavelength and intensity at this resolution, proves to be quite high, about 360 s⁻¹.

Most detection systems flown by the Division have used electrometer amplifiers with large value feedback

resistors to measure the small currents (typically 10⁻⁹ A) available at the anodes of the photomultipliers. However, the dynamic range required to observe a substantial number of stars with tolerable accuracy is in the order of 1000 : 1. This can in principle be achieved by logarithmic and quasi-logarithmic feedback elements, but difficulties in resolution are then introduced owing to the compression of the data into the 2 V full scale of the telemetry channel, the absolute resolution of the received signal being some 2% of full scale. The problem of obtaining a wide frequency response at high sensitivity is also substantial, while the necessity of running the photocathode at high potential can cause spurious noise.

The system described in this paper overcomes these problems by detecting the pulses produced at the anode of the photomultiplier corresponding to the arrival of photons.² Due to restrictions in the bandwidth of the telemetry channels available§ it is not possible to transmit the count in raw digital form. An encoding system was therefore developed which enabled these two channels to transmit all of the necessary experiment output in an analogue coded form in which a logarithmic type of characteristic arises naturally. (Such a pulse counting mode also gives a very valuable increase³ in the information content compared with the direct current mode of measurement.)

2. System Description

Figure 1 shows a block diagram of the main components of the system. The photon-generated pulses from the photomultiplier (an E.M.I. type D.104, having a Cs-Te photocathode with quartz window), pass through the pre-amplifier to the amplifier and threshold discriminator. Pulses satisfying the threshold conditions are allowed into the counter during the counting interval and gated out during the sample-

† Space Research Division, Royal Observatory, Edinburgh.
‡ Formerly at the Royal Observatory, Edinburgh; now at Electrical Engineering Department, University of Edinburgh.

§ The telemetry used was I.R.I.G. f.m./f.m. Two sub-carrier channels were available for the experiment, channel 16 which has a bandwidth of 600 Hz and channel 15 which has a bandwidth of 450 Hz.

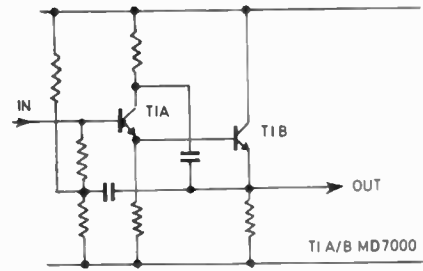
and-reset interval. The condition of the counter is represented by a 12-bit binary number which is presented to the encoder and converted to the form described later. The encoder outputs are fed to the sample-and-hold amplifiers whose outputs then pass to telemetry channels 15 and 16.

2.1 Pre-amplifier, Amplifier and Threshold Discriminator

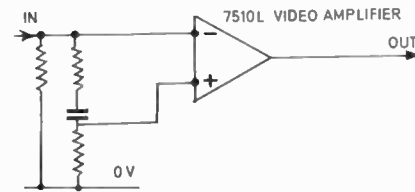
The pre-amplifier, Fig. 2(a), is a White follower and uses a dual n-p-n transistor. It is mounted next to the photomultiplier and drives the length of cable necessary to connect to the amplifier, Fig. 2(b), which is a commercial integrated circuit video amplifier used without external feedback in the inverting mode. The threshold discriminator, Fig. 2(c), is an i.c. voltage comparator connected as a Schmitt trigger. The threshold level is variable between 0-5 V and the strobe terminal of the i.c. is used to gate the input to the counter.

2.2 The Counter

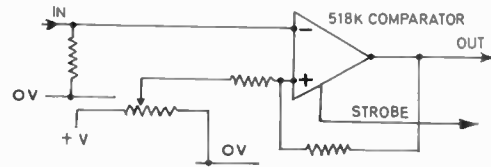
The 4-bit i.c. binary counters, Fig. 3, are cascaded to form a 12-bit binary counter. The relationship between the capacity of the counter, the sampling frequency (see introduction), and the resolution of the counter was chosen so that stars of apparent visual magnitudes from zero to 6, (a range of 250 : 1) could be observed without substantial loss of accuracy when compared with the stochastic variation (see Sect. 3). The expected photon counts were taken from Houziaux.⁴ When the counter is full, a carry pulse from the last stage latches flip-flop 1 which causes channel 15 to show an overload level.



(a) Pre-amplifier.



(b) Amplifier.



(c) Threshold discriminator.

Fig. 2.

2.3 The Encoder

Only the ten most significant bits of the counter are used by the encoder and are shown as P1 to P10 (Fig. 3). The functions of the encoder are:

- (i) to select the five most significant bits from the 10-bit number represented by P1 to P10;
- (ii) to indicate the position of the five bits within the 10-bit number;
- (iii) to form analogue voltages corresponding to (i) and (ii) in the form shown in Fig. 5. Functions (i) and (ii) are illustrated in Table 1 and Table 2.

In Table 1 the six 5-bit words corresponding to the possible conditions of the upper 10-bits of the counter are shown. It is convenient to generate six quantities Q1 to Q6, which characterize each of the six possible words. Table 1 is a truth table forming Q1 to Q6 as functions of P1 to P10. Table 2 giving Q1 to Q6 in minterm form is written down and the combinational logic corresponding to these terms is carried out by gates 1 to 6 of Fig. 3. Gates 13 to 42 form the Boolean products of the components of the required 5-bit

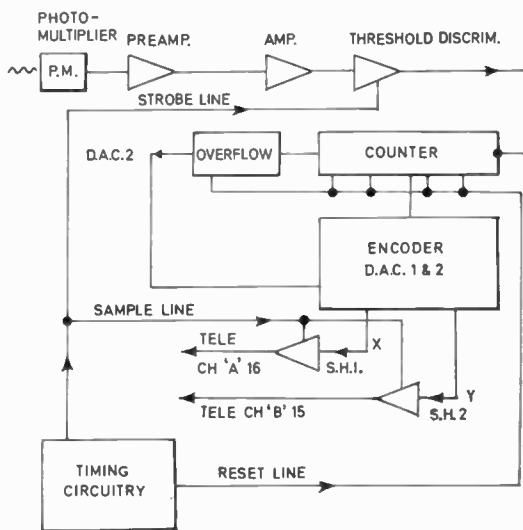


Fig. 1. Block diagram of system.

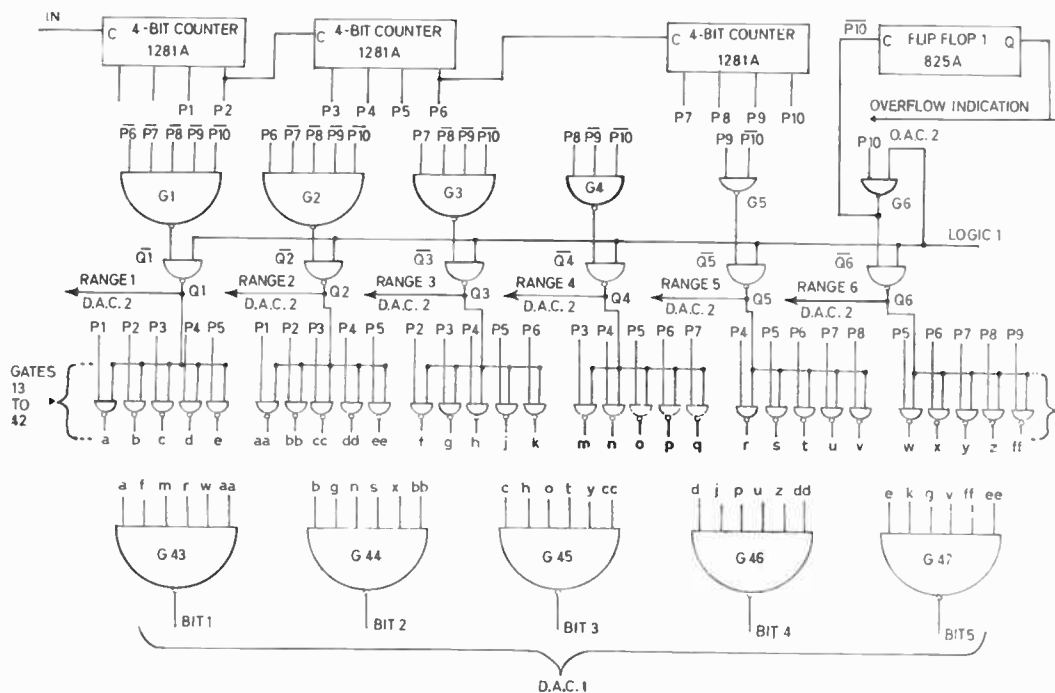


Fig. 3. Logic diagram. The inverting stages between counters and gates are not shown. Signetics 8000 series logic was used throughout.

Table 1

P1	P2	P3	P4	P5	P6	P7	P8	P9	P10	Q6	Q5	Q4	Q3	Q2	Q1
				X	X	X	X	X	1	1	0	0	0	0	0
			X	X	X	X	X	1	0	0	1	0	0	0	0
		X	X	X	X	X	1	0	0	0	0	1	0	0	0
X	X	X	X	X	1	0	0	0	0	0	0	0	0	1	0
X	X	X	X	X	0	0	0	0	0	0	0	0	0	0	1

words with the quantities Q1 to Q6, so that the points, a, b, c, d and e, for example, indicate the first word only when Q1 is high, and are at logical 1 otherwise. Gates 43 to 47 combine the five words serially.

Figure 4(a) shows the digital-to-analogue converter (D.A.C.1). The basis of the converter is a ladder network⁵ with two resistor values, 10k and 20k which converts the 5-bit binary number formed by the outputs of gates 43 to 47 to an analogue form available at the point X in Fig. 1. A second digital-to-analogue converter (D.A.C.2), Fig. 4(b), generates voltages which identify each of the quantities Q1 to Q6 and these voltages are available at point Y in Fig. 1. The points X and Y form the outputs of the encoder and the relationship between the voltages at these points and the count accrued during a sample period is shown in Fig. 5. The output level of the overflow flip-flop in Fig. 1 overrides inputs Q1 to Q6 to produce

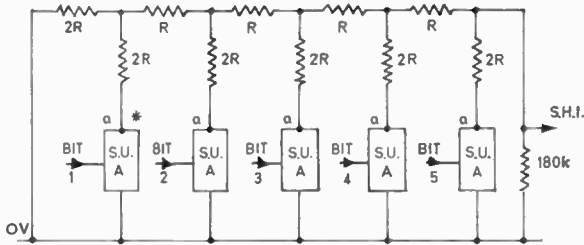
Table 2

Q1 =	$\overline{P6}$.	$\overline{P7}$.	$\overline{P8}$.	$\overline{P9}$.	$\overline{P10}$.
Q2 =	$\overline{P6}$.	$\overline{P7}$.	$\overline{P8}$.	$\overline{P9}$.	$\overline{P10}$.
Q3 =		$\overline{P7}$.	$\overline{P8}$.	$\overline{P9}$.	$\overline{P10}$.
Q4 =			$\overline{P8}$.	$\overline{P9}$.	$\overline{P10}$.
Q5 =				$\overline{P9}$.	$\overline{P10}$.
Q6 =					$\overline{P10}$.

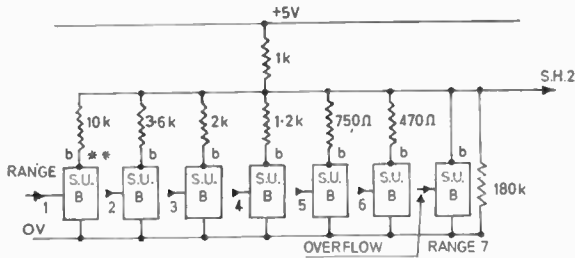
the overload indication shown as range 7 in output Y of Fig. 5.

2.4 Sample and Hold Systems

S.H.1 and S.H.2 are identical; Fig. 6 shows the circuit details. In the flight instrument a sampling frequency of 360 per second was used. During the hold phase the f.e.t. switch is held off by the -12.5 V gate potential. The sample voltage is maintained across the capacitor by the operational amplifier until



(a) Digital-to-analogue converter 1.



(b) Digital-to-analogue converter 2.

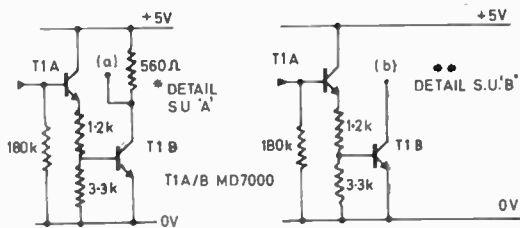


Fig. 4.

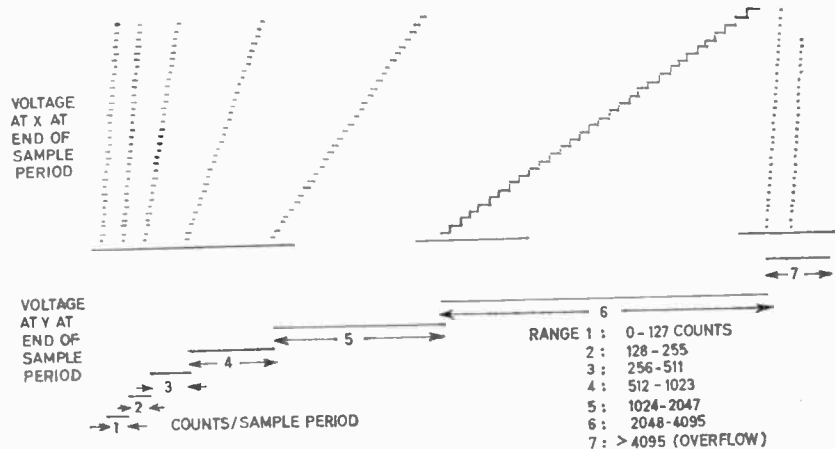


Fig. 5. Voltages at points X and Y at end of sample period.

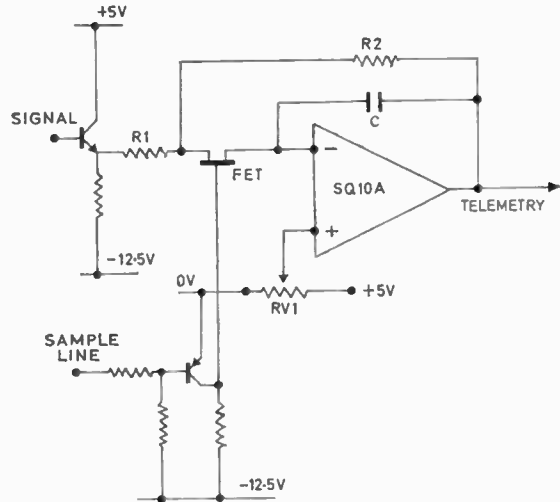


Fig. 6. Sample and hold system.

the end of the hold phase at which time the sampling line rises to near zero and switches on the f.e.t. In the sample phase the system operates as an inverting amplifier with a gain set primarily by R1 and R2. Positive or negative excursions are made possible by adjusting the centre tap of RV1. To ensure compatibility with the telemetry system, the output was adjusted to occupy the range 0 to +2 V. The sampling interval occupies about 6% of the total cycle.

2.5 Timing Circuitry

Three functions are required by the system:

- (i) to obtain the sample rate;
- (ii) to isolate the counter input during the sample and reset period;
- (iii) to reset the counter.

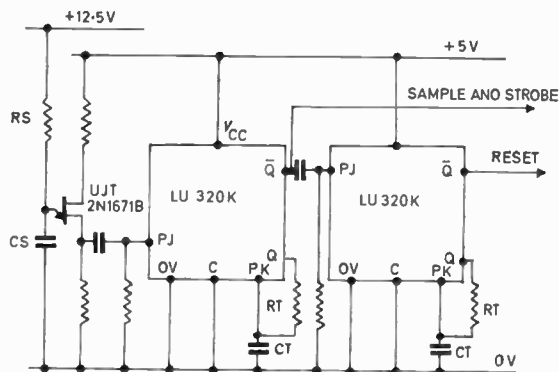


Fig. 7. Timing circuitry using two Signetics t.t.l. flip-flops (L.U. 320K).

The circuit details are shown in Fig. 7. The basic timing pulse is obtained from a unijunction transistor which drives two cascaded t.t.l. flip-flops connected as monostables. The output pulse widths of the monostables are fixed by C_T and R_T . The width of the 'sample' pulse is 200 μ s, and that of the reset pulse is 10 μ s.

3. Accuracy and Performance

Any discussion of accuracy in a photon counting system should consider the effect of photon statistics which set an intrinsic upper limit to what is possible.

If N photons with a Poissonian distribution are counted, the probable standard deviation of the mean count is $N^{1/2}$. Therefore the proportional statistical error for a count of N photons is,

$$100/N^{1/2}\% \quad \dots\dots(1)$$

A total of 90 spectra, corresponding to some 75 stars, has been obtained from the records. Of these nearly half have sufficient signal strength and are otherwise suitable to be of scientific value.¹ Of the useful spectra, 22 (those of the brightest and hottest stars) are rated as good or very good, having r.m.s. deviation due to the photon statistics at each recorded spectral wavelength, of less than 10%, i.e. more than 100 stellar photons were counted during the integration period. In order that such recorded data are not appreciably degraded, the resolution of the counting system must be appreciably (two or three times) better than this. Such resolution is then better than is necessary for the fainter and cooler stars.

3.1 Encoder Accuracy

Single step voltage changes at the input of S.H.1 (Fig. 1) only occur for every n counts when $n = 4, 4, 8, 16, \dots$ according to the range. Therefore there is an uncertainty in the exact number of counts corresponding to a given level L . Following the procedure of Cliff⁶ we may minimize the error arising from this by assigning a nominal count value \bar{C} to this level such that the possible uncertainty in the count is symmetrical about \bar{C} . We have then two expressions for \bar{C} : on range 1, equation (2) applies, while on ranges 2 to 6 we use equation (3).

$$\bar{C} = 2 \times (2L - 1) - \frac{1}{2} \quad \dots\dots(2)$$

$$\bar{C} = 2^R \times (L + 32 - \frac{1}{2}) - \frac{1}{2} \quad \dots\dots(3)$$

where \bar{C} = nominal count value

R = range

L = level number (from 1 to 32).

The actual count C in range 1 and ranges 2 to 6 can differ from the nominal count within the limits shown

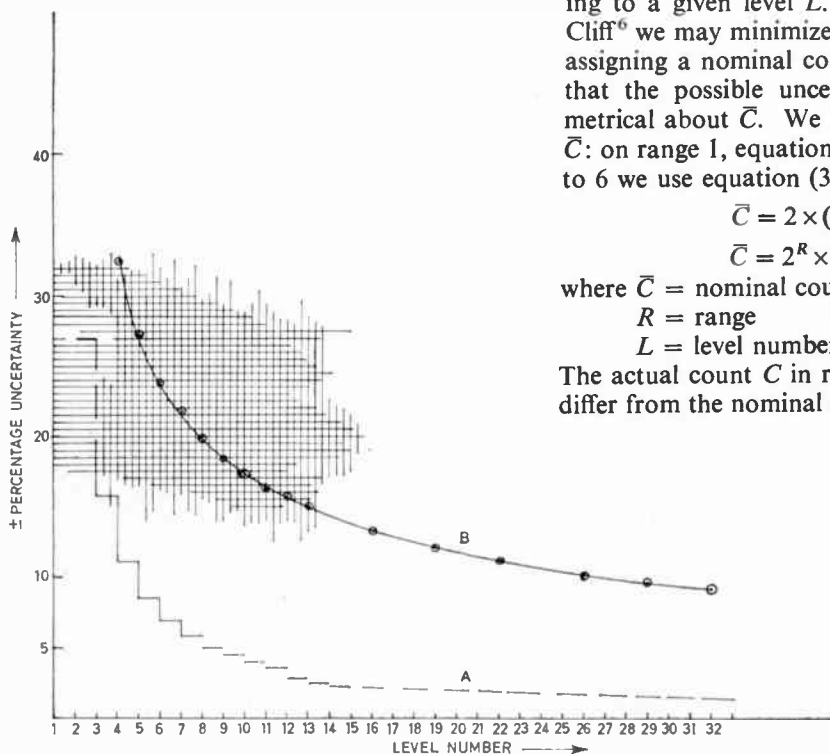


Fig. 8. Percentage uncertainties of counting system (A) and photon statistical accuracy (B) on Range 1.

in (4) and (5) respectively.

$$2^R \times (2L-2) \leq C \leq 2^R \times (2L)-1 \quad \dots\dots(4)$$

$$2^R \times (L+31) \leq C \leq 2^R \times (L+32)-1 \quad \dots\dots(5)$$

Accuracy: Range 1:

From (4) it can be seen that the count on any level L in range 1 can vary by $\pm 1\frac{1}{2}$ counts and that the possible uncertainty will be

$$\pm \frac{1\frac{1}{2} \times 100}{2 \times (2L-1) - \frac{1}{2}} \% \quad \dots\dots(6)$$

In Fig. 8 these uncertainties are plotted as a function of level number at A. The percentage photon statistical accuracy as obtained by substituting in (1) the value $N = 2 \times (2L-1) - \frac{1}{2}$ is plotted on the same graph at B.

Accuracy: Ranges 2 to 6:

In ranges 2 to 6 using (3) the nominal count is

$$\bar{C} = 2^R \times (L+32 - \frac{1}{2}) - \frac{1}{2}$$

From (5) the actual count can differ from the nominal by at most $\pm \frac{1}{2}(2^R - 1)$, so that the percentage uncertainty on any step is at most:

$$\pm \frac{\frac{1}{2} \times (2^R - 1) \times 100}{2^R \times (L+32 - \frac{1}{2}) - \frac{1}{2}}$$

which simplifies to

$$\frac{(2^R - 1) \times 100}{2^R \times (2L+63) - 1} \% \quad \dots\dots(7)$$

In Fig. 9 the percentage uncertainties are plotted as functions of level and range number. The photon statistical accuracy is plotted on the same graph.

3.2 Channel Utilization

The coding system of Fig. 5 defines an alphabet of 32×6 characters, so that each character represents $\log_2(192) = 7.59$ bits of information. For a sample rate of 360 per second, the information output is a constant 2732 bit/s. If the two channels are considered separately, the information is divided:

Channel 16 : 1800 bit/s

Channel 15 : 932 bit/s

The cost of rocket astronomy is considerable, and the experimental environment is severe. It is common therefore for generous redundancy to be allowed when matching telemetry channels to experiments. Final pre-flight checks may be made more confidently and the effects of certain minor system failures during the flight minimized, if the information output is deliberately limited to a value appreciably below the theoretical channel maximum. Thus the spacing of steps in channel 16 was set at 50 mV at the input to the telemetry sender, although the specified telemetry resolution implies the possibility of steps spaced by only 25 mV. Ten per cent of the 2 V telemetry scale was reserved on each channel for a 'channel live' indication, so that the maximum number of steps possible was $1.8/0.0250 = 72$ levels. Each level would have a value of 6.18 bits compared with the 5 bits per level value for channel 16 in the flight instrument.

Similar considerations apply in the selection of a sample rate. The maximum possible rate is set by the bandwidth of channel 15. In practice a considerably lower rate was chosen, which allowed the shape of the output histograms to be resolved easily.

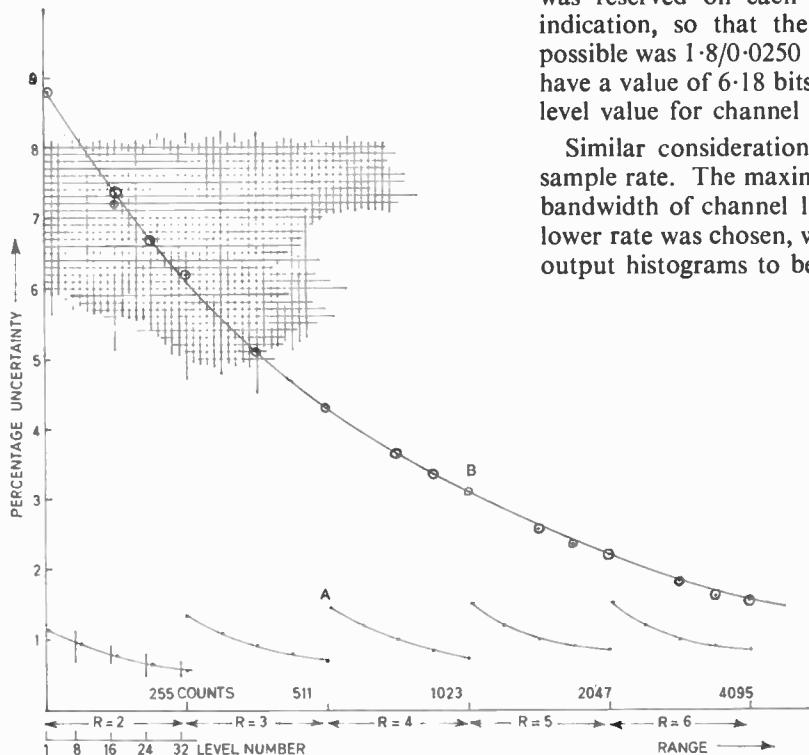


Fig. 9. Percentage uncertainties of counting system (A) and photon statistical accuracy (B) on Ranges 2 to 6.

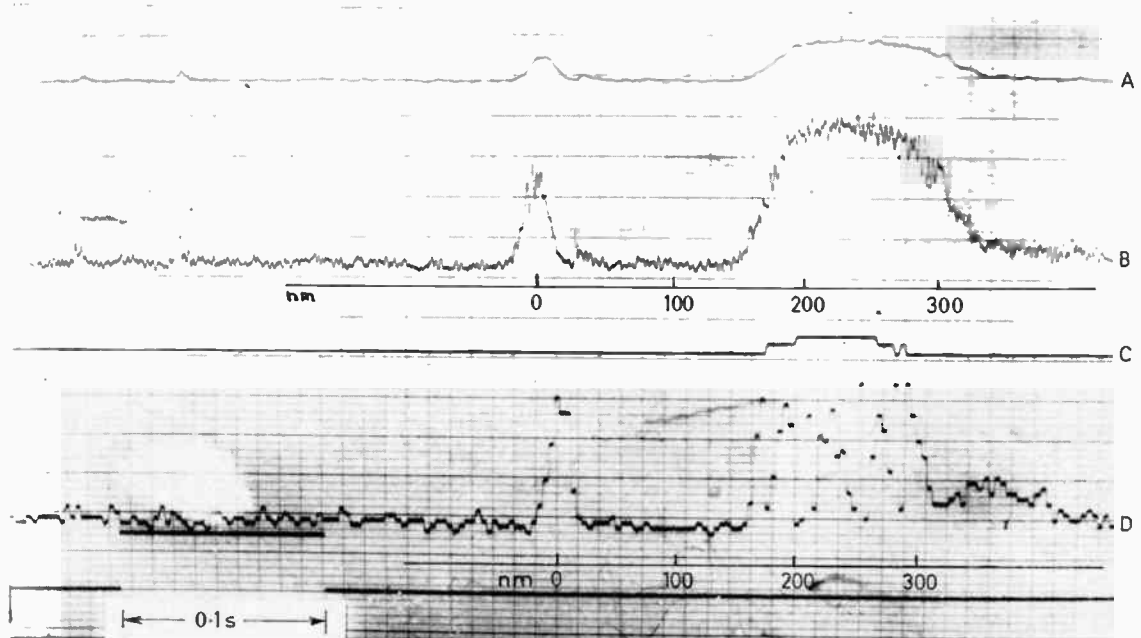


Fig. 10. Photograph of part of the telemetry record showing typical stellar spectrum.

3.3 Flight Performance

A *Skylark* rocket containing the system described in this paper was successfully flown from the Salto di Quirra range in Sardinia on 3rd December, 1968. The initial spin rate generated by the roll control unit was just over 15° per second, decaying to nearly 14° per second by re-entry and giving a dispersion of 1.75 nm per ms. Figure 10 is a photograph of part of the telemetry record showing a typical stellar spectrum (the star is Beta Cephei). Trace C shows the range indication channel whilst trace D is the count rate channel. Traces A and B show the same spectrum as obtained from d.c. amplifiers which were used on a separate spectrophotometer mounted on the same payload.

Reading from left to right on Fig. 10, all traces A, B, and D show initially the 'zero order' followed by the star's spectrum. Comparison of trace C with D shows that the zero order did not exceed the capacity of the first range (i.e. 127 counts maximum) whilst the spectrum gave counts up to the third range (i.e. 511 counts).

In addition to providing good spectral data, the low dark count of the Cs-Te photomultiplier enabled absolutely quantified and reliable data to be obtained between the stellar signals. These data have been valuable in estimating zodiacal and galactic fluxes and brightness distributions near 240 nm.⁷

4. Acknowledgments

The authors wish to thank G. C. Sudbury and Dr. H. E. Butler for their encouragement, and N. P. Bennett for much helpful discussion during the design of the system. They also wish to thank the Astronomer Royal for Scotland, Professor H. A. Brück, C.B.E., and the Science Research Council for permission to publish this paper.

5. References

1. Sudbury, G. C., 'A rocket-borne photoelectric spectrophotometer using convergent beam dispersion to observe far ultraviolet stellar spectra', *Applied Optics*, 8, No. 10, p. 2013, October 1969.
2. Baum, W. A., 'Counting photons', *Sky and Telescope*, p. 265, May 1965.
3. Alfano, R. R. and Ockman, N., 'Methods for detecting weak light signals', *J. Opt. Soc. Amer.*, 58, No. 1, p. 90, January 1968.
4. Houziaux, L., E.S.R.O. Scientific Memorandum, E.S.R.O. S.P.-3, June 1965.
5. Freeman, J., *Electronic Design*, p. 68, 5th July 1967.
6. Cliff, R. A., 'Logarithmic compression of digitally telemetered data', *I.E.E.E., Trans. on Aerospace and Electronic Systems*, AES-3, No. 4, p. 712, July 1967.
7. Sudbury, G. C. and Ingham, M. F., 'The background brightness of the Milky Way near 2500 Å between $l^{\text{II}} = 72^\circ$ and 126° ', *Nature*, 226, No. 5245, p. 526, 9th May 1970.

Manuscript first received by the Institution on 23rd December 1969 and in final form on 18th May 1970. (Paper No. 1342/AMMS31).

© The Institution of Electronic and Radio Engineers, 1970

Contributors to the Journal



Mr. C. H. Paterson graduated as B.Sc. in electrical engineering at the University of Edinburgh in 1964. He then joined the staff of Nuclear Enterprises (G.B.) Ltd., where he worked until 1966 on nuclear radiation detection systems and data logging equipment. As a member of the Space Research Group of the Royal Observatory, Edinburgh from 1966-68, he was involved in

the electronic design of rocket-borne astrophysical experiments. He is now doing Ph.D. research at the University of Edinburgh on surface effects in the conductivity of ultra-thin metal films.



Mr. D. H. Beattie (G. 1967) received his technical education at the Technical Colleges of Carlisle and Workington, and the Napier College of Science and Technology, Edinburgh. He worked initially with the North Western Electricity Board, moving in 1959 to the Spadeadam Rocket Establishment, Cumberland, where he was with the Instrumentation Division of Rolls-Royce. In

1964 he took up his present appointment of experimental officer in the Space Research Division of the Royal Observatory, Edinburgh, where he works on the electrical and electronic design requirements of rocket and satellite-borne experiments and their associated ground-based test equipments.



Dr. A. P. Clark (M. 1964) graduated in natural sciences at the University of Cambridge in 1954 and joined the Cambridge Instrument Company as an assistant physicist. In October 1955 he was appointed as an engineer with British Telecommunications Research Limited, Taplow, now Plessey Telecommunications Research. Since then he has worked on a variety of radio engineering

projects, principally concerned with data transmission techniques and in July 1970 he was promoted to senior principal engineer. From 1965 to 1968 Dr. Clark held an Industrial Fellowship in the Department of Electrical Engineering at Imperial College, working on the trans-

mission of digitally-coded signals by means of random access discrete address systems. Subsequently he applied these techniques to adaptive detection processes and this work provided the basis of his doctoral thesis. Dr. Clark is author of a number of papers and articles and holder of several patents. His present paper is the third to be published in the Institution's *Journal*. On October 1st Dr. Clark takes up an appointment as lecturer in the Department of Electronic and Electrical Engineering at Loughborough University of Technology.



Professor Branko D. Rakovich received his bachelor's and doctor's degrees from the University of Belgrade in 1948 and 1955 respectively. From 1948 to 1951 he served as a research assistant at the Institute for Telecommunications of the Serbian Academy of Science, Belgrade. He then joined the Faculty of Electrical Engineering at the University of Belgrade as a teaching

assistant until 1954 when he became an assistant professor. From 1954 to 1955 he held a National Education and Research Council of Yugoslavia Fellowship at Marconi College, Chelmsford, England, and at Imperial College, London. In 1959 he became an associate professor of electronics and later a full professor of electrical engineering. Currently he is also the head of the Department of Electronics of the Faculty of Electrical Engineering, University of Belgrade. Professor Rakovich has been a consultant to the Federal Research Council of Yugoslavia since 1960. His research interests lie in the field of active and passive network synthesis, particularly in wide-band amplifier theory and design.



Mr. A. D. Jovanovich graduated with a B.E.E. degree from Belgrade University in 1960. He is at present studying for his Ph.D. in the Department of Electronics of the Faculty of Electrical Engineering, University of Belgrade, where since 1965 he has held an appointment as a teaching assistant in electronics. He worked as a research engineer at the Institute for Automation

and Telecommunications 'Mihailo Pupin' in Belgrade from 1962 to 1965. Mr. Jovanovich has published papers on computer-aided design of electrical circuits and other computer applications.

High Accuracy Digital Linearization of Frequency Signals of Transducers

By

G. MAYER,

Dipl.Ing.,†

L. SIMONFAI,

Dipl.Ing.†

and

P. PÓTY,

Dipl.Ing.†

Reprinted from the Proceedings of the I.E.R.E. Conference on "Digital Methods of Measurement" held at the University of Kent at Canterbury on 23rd to 25th July 1969.

The signals from transducers having non-linear frequency output have to be linearized by a signal converter. The paper describes first some accepted linearizing method and afterwards the digital linearizer developed in the Central Measurement Research Laboratory, Hungary, will be discussed. The DENSITON electronic units connected to the frequency output density transducer correct digitally the error curve by pulse-rate modulation of clock-pulses. The system design of the instrument built up of integrated circuits is outlined, and various problems regarding the settings of the linearizer discussed.

1. Introduction

The output signal of analogue industrial transducers generally has a linear relation to the primary quantity to be measured. Linearity has been of prime importance with the development of transducers to apply simple signal processing unit. The growing use of digital data loggers and process control computers has drawn attention to non-coded digital output transducers.¹ The repeatability and long term stability of these frequency output transducers (e.g. vibrating wires, beams, cylinders and quartz crystal oscillators) is very good, but the output signal is not linear. The linearity, however, need not be a prime factor in the future when designing transducers, and that fact makes possible the application of physical principles where higher repeatability of frequency output or other benefits may be ensured.

Figure 1 represents some possibilities of the system engineering to accomplish the linearization. Analogue linearizers of various solutions are well known, most of them perform the required straight-line approximation by using diodes. Their common disadvantage is that the elimination of temperature dependence is rather difficult and the practical accuracy limited.

The economics and reliability of new electronic components—first of all integrated circuits—make possible the building of a high accuracy digital linearizer of non-linear functions at moderate price. Digital linearizer may be used advantageously with the solutions shown on the last two rows of Fig. 1.

A digital linearizer called the DENSITON has been developed, which reduces the non-linearity of quantities represented either by frequency or by pulse number, to about 1 : 100. The task of the research team was to develop an electronic unit having analogue

† Central Measurement Research Laboratory, Budapest, Hungary.

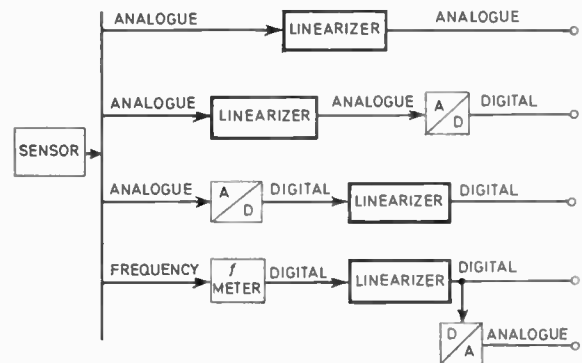


Fig. 1. Some possible methods of linearization.

and/or digital output for a vibrating tube gas or liquid density transducer. The last row of Fig. 1 represents that solution.

2. Survey of Various Solutions

The relation between the output frequency of the transducer and the density measured may be expressed with the following equation:

$$f = f_0 \sqrt{\frac{1}{1 + d/d_0}} \quad \dots (1)$$

where f = output frequency of the transducer,

f_0 = output frequency of the transducer at density $d = 0$,

d_0 = dimensional constant,

d = density to be measured.

The f and f_0 frequencies are of kilohertz order, the measurable maximum density is 120 mg/cm^3 with gases, and 1.5 g/cm^3 with liquids. The tolerance of d_0 and f_0 constants may be as high as $\pm 10\%$. The electronic unit connected to the frequency output

transducer must be set between $d_2 - d_1$ span in the light of tolerance of d_0 and f_0 .

With the digital linearization, the function to be linearized is usually approximated by straight lines. Conforming to the slope of the individual straight segments, pulses are inhibited from or added to the input quantity represented by the pulse number.

Three different methods may be used for digital linearization:

- (1) Application of a special function generator, which inverts the function to be linearized by digital method. In that case there is a necessity for a square generator. That solution is applied in the mass flow computer of Electronic Flow-Meters Ltd.² The method is suitable only for linearizing of special functions, and it is relatively expensive to achieve a given accuracy.
- (2) Linearization of the characteristics. The curve to be linearized is approximated by straight segments and the corrections necessary to obtain the different slopes of the segments are realized by binary multipliers. Programming is easy but the method does not result in a very exact linearization.³
- (3) Linearization of the error curve. The error curve, i.e. the difference between the non-linear characteristics and the desired linear correlation is approximated by straight segments, and the correction necessary for the segments having different slopes is computed.⁴ The method results in a very high accuracy, but it is relatively hard to program. The necessary hardware may be decreased and the programming may be simplified if the error curve has an axis of symmetry. It is often advisable to write the general error function into a Taylor sequence up to the term of the second order. The resulting approximating error function will be symmetrical, and the above advantages may be realized.

3. Description of the DENSITON Instrument

3.1. System Engineering and Operation

The block diagram of DENSITON electronics is shown on Fig. 2. The amplified and squared signal from the transducer is fed to the PR1 binary counter. The clock generator consists of a high-stability quartz oscillator. During the first cycle of the measurement the PR3 counter operates as a preset counter, during the second cycle it serves for counting the linearized result-pulses and also controls the linearizer.

All counters are reset at the beginning of the measurement and the signal of the clock generator is inhibited. On receiving the first pulse from the transducer the

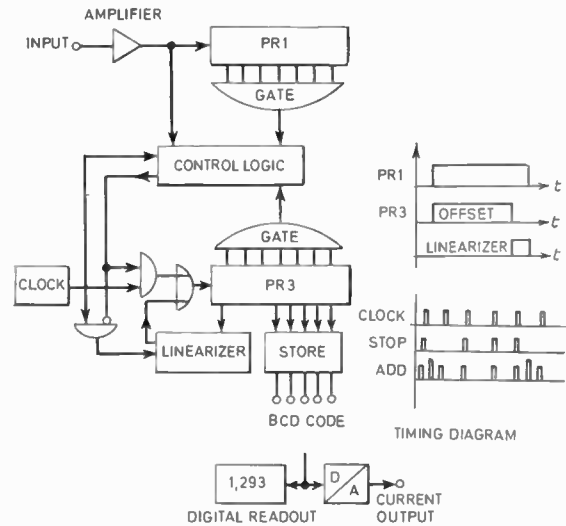


Fig. 2. Block diagram of DENSITON instrument.

control unit feeds clock pulses to the PR3 preset counter. This counter generates the zero offset conforming to the T_1 periodic time appropriate to the lower limit of the $d_2 - d_1$ span. During that measuring cycle the linearizer is inhibited. When the PR3 counter arrives at the adjusted preset value, a change takes place in the output of the decoding gate. The control unit then resets the PR3 counter, opens the linearizer and gates the clock pulses. During the second cycle the clock pulses are fed through the linearizer to the PR3 counter. At certain intervals a pulse is inhibited from or added to the clock pulses depending on the slope of the straight line polygon approximating the error curve. The measurement comes to an end when the PR1 preset counter arrives at the adjusted value. Now the control unit writes the result from the PR3 counter into the stores and resets all the counters. The contents of the stores are fed into a decoder and/or digital-to-analogue converter. The measurement is repeated automatically.

The requirements of the linearizer are the following:

- (1) The accuracy of the linearization should be of a given value, in this case 0.1%.
- (2) Error of control settings should be minimum.
- (3) The electronic unit should make allowance for the tolerance of parameters.
- (4) The programming should be as simple as possible. With respect to easy-to-set hardware, the question of built-in redundancy must arise in the interest of the instrument.
- (5) The value measured must appear in engineering units (g/cm^3 , mg/cm^3), in code form.

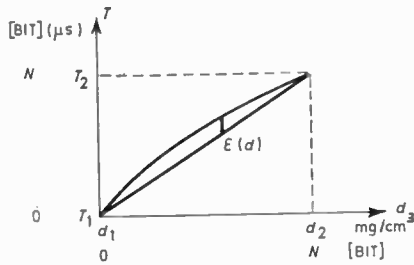


Fig. 3. Characteristic of sensor for determining the error curve.

It is the task of the system engineer to achieve a compromise between the above requirements. Before outlining the construction and operation of the linearizer, we will deal briefly with the determination of the error curve, and the mathematical description of the polygon approximating the error curve.

3.2. Determination of Error Curve and the Amount of Maximum Error

(1) Expressed as the periodic time, becomes equation

$$T = T_0 \sqrt{1 + \frac{d}{d_0}} \quad \dots\dots(2)$$

The lower limit of the density span to be measured is d_1 and the related periodic time is T_1 ; the upper limit of the density is d_2 and the related periodic time is T_2 . Plot the equation (2) to be the d_1 density and T_1 periodic time in the origin of the coordinates (Fig. 3). Without any linearization the measurement of periodic time would take place along the straight line connecting the two end-points of the curve. The error curve is equal to the difference between the non-linear function and the straight line.

On the output of the digital linearizer a code proportional to the density measured appears for which there is a pulse number proportional to the density. It is necessary to scale these to give engineering units. Suppose 0 pulse represents d_1 density and N pulses the d_2 density. Multiply the non-linear function of equation (2) by a constant to give a physical dimension which is also a density and/or pulse number. Regarding the extreme values of the function, the constant value will be

$$c = \frac{d_2 - d_1}{T_2 - T_1} \quad \dots\dots(3)$$

The equations of the non-linear curve and the straight line in the plotted coordinate system are:

$$T_s = T_1 + \frac{T_2 - T_1}{d_2 - d_1} (d - d_1) \quad \dots\dots(4)$$

$$T_c = T_0 \sqrt{1 + \frac{d}{d_0}} \quad \dots\dots(5)$$

The following relation may be established for determining the error curve:

$$\epsilon(d) = T_c - T_s = T_0 \sqrt{1 + \frac{d}{d_0}} - T_1 - \frac{T_2 - T_1}{d_2 - d_1} (d - d_1) \quad \dots\dots(6)$$

which may be converted to the following equation:

$$\epsilon(d) = \frac{N}{T_2 - T_1} \left(T_0 \sqrt{1 + \frac{d}{d_0}} - T_1 \right) - n \quad \dots\dots(7)$$

where n is the number of pulses conforming to a density of d . When determining the error curve the error was referred to the result impulse. It is possible to refer the error to the clock pulses, and in that case the inverse function must be derived from equation (2), and the equation of the $\epsilon(T)$ error curve must be determined by the method described above.

$$\epsilon(T) = \frac{d_0}{T_0^2} [T^2 - (T_1 + T_2) \cdot T + T_1 T_2] \quad \dots\dots(8)$$

Note that the error curve resulting from the first method comprises a square root member, and therefore it is not geometrically symmetrical. The mathematical formula of the error curve when referred to the clock pulses is a curve of the second order, which is always symmetrical. The maximum value of the relative error is

$$h_{max} = 25 \frac{T_2 - T_1}{T_2 + T_1} \% \quad \dots\dots(9)$$

The tolerance of T_0 and d_0 transducer parameters is $\pm 10\%$. The effect of tolerance shows in the amount of maximum error. It may be proved that the maximum value of error is independent of the tolerance of T_0 , and the tolerance of d_0 parameters alters the relative value of maximum error:

$$\frac{\Delta h_{max}}{h_{max}} = - \frac{1}{\sqrt{1 + d_1/d_0} \cdot \sqrt{1 + d_2/d_0}} \cdot \frac{\Delta d_0}{d_0} \quad \dots\dots(10)$$

After determining the equation of the error curve and the place and amount of maximum error, the symmetrical error curve according to equation (8) is plotted in a new coordinate system. Both coordinates must be related to the pulse dimension. The equipment approximates the error curve with a straight line polygon. The desired straight line characteristic to ensure the desired difference must be traced in the symmetrical error curve (Fig. 4).

The straight lines are characterized by their slopes. Pulses must be inhibited from the incoming clock pulses in the segment, where the slope is negative. At the segments where the slope is zero, the clock frequency has not to be corrected. At the segments which have a positive slope, pulses must be added to the clock pulses. The slopes of the straight lines are characterized by an expression represented by a

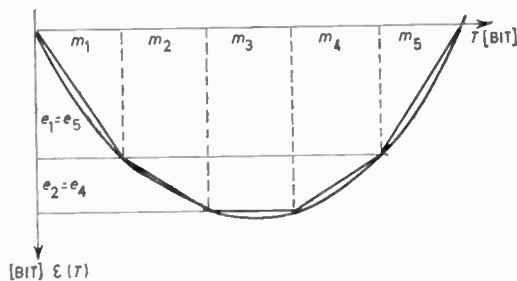


Fig. 4. Straight line approximation. Density sensor curve is plotted then segmented.

fractional number, in which the numerators of the numbers representing the segments having negative slope determine the number of pulses to be inhibited in the given segment. The denominator represents the number of clock pulses, i.e. the sum of the result-pulses and inhibited pulses.

Thus

$$\frac{L_1}{K_1} = \frac{e_1}{m_1 + e_1} \dots\dots(11)$$

where e_1 = error pulses in m_1 segment (inhibited pulses),

and m_1 = length of the first segment (the number of the corrected result-pulses).

The numerator of the fractional number representing segments having positive slope determines the number of pulses to be added to the clock pulses and the denominator represents the number of clock pulses, i.e. the difference between the result-pulses and added pulses.

$$\frac{L_4}{K_4} = \frac{e_4}{m_4 - e_4} \dots\dots(12)$$

where e_4 = error pulses in m_4 segment (added pulses), and m_4 = length of the fourth segment in pulses (the number of the corrected result-pulses).

The number of segments has to be determined when plotting the straight lines. If the number of segments is reduced, the deviation from the straight lines becomes greater, but the control is easier. It is easy to take the individual characteristics which are due to the parameter tolerance into consideration.

By increasing the number of segments a linearization of any desired accuracy may be achieved, but in that case the control is complicated and the cost of the instrument increases. From 5 to 7 segments are regarded as practical value, and enable the non-linearity of about 10% to be reduced to 0.1%. It is expedient to introduce a segment of zero slope in the middle part of the error curve, since that is the cheapest segment and does not complicate the control.

3.3. Linearizer (Fig. 5)

During the second cycle of the measurement the PR3 counter is reset again after adding the offset, and the clock pulses are fed to the linearizer.

The addition and subtraction depending on the slope of the given straight line is controlled by the content of the result counter through the segment-limit decoding and storing circuits.

Suppose the binary counter to be k bits. The L_n/K_n fractional number determined from the error curve is reduced to the form of

$$\frac{L_n}{K_n} = \frac{P_n}{Q_n} \text{ where } Q_n \leq 2^k - 1.$$

The length of the binary counter determines the maximum of the denominator in the fractional number representing the slope of the segment. Theoretically defined fractional numbers must be realized for the error curve in a way that the fraction is selected which is nearest to the calculated value. If the numerator is short, the denominator is a small number, therefore individual fractions are hard to approximate. The authors' experiences proved that a numerator of 6-7 bits gives a satisfactory result.

The combinational network connected to the outputs of the binary counter realizes the P_n/Q_n fractional numbers. In the n th segment the cycle length of the binary counter is Q_n , and during that time coincidence appears P_n times on the outputs of the combinational network. P_n must be selected from Q_n state of the counter, so that the location of the decoded positions will be equalized on the Karnaugh map. Often $P_n = 1$, and in those cases the binary counter is reduced for Q_n period to be preset counter, which modifies the number of clock pulses producing one pulse for every Q_n .

Generally the combinational network may be simply constructed by using the Karnaugh diagram.

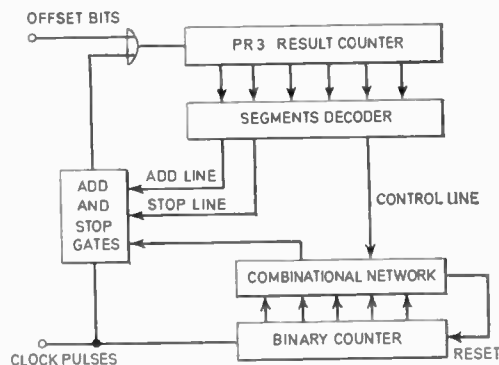


Fig. 5. Block diagram of linearizer.

Table 1
The main characteristics of the linearizer for various density ranges

	$d_2 - d_1$	h_{\max} %	Number of segments	Number of fractional numbers	Calculation of fractional numbers
Gas sensor	10 mg/cm ³	0.85	3	1	Closed equation
$T_0 = 260 \mu\text{s} \pm 10\%$	20 mg/cm ³	1.52	3	1	Closed equation
	40 mg/cm ³	2.75	5	2	Closed equation
$d_0 = 60 \text{ mg/cm}^3 \pm 10\%$	80 mg/cm ³	4.5	7	3	or plotted error curve
Liquid sensor	0.1 g/cm ³	0.91	3	1	Closed equation
$T_0 = 180 \mu\text{s} \pm 10\%$	0.2 g/cm ³	1.56	3	1	Closed equation
	0.4 g/cm ³	2.90	5	2	Closed equation
$d_0 = 210 \text{ mg/cm}^3 \pm 10\%$	0.8 g/cm ³	4.5	7	3	or plotted error curve

Suppose $P_n/Q_n = 4/17$, and the four states prescribed by the counter to be 3, 7, 11 and 15, which may be decoded by a single, two-input NAND gate.

The 17th state may be decoded by another two-input NAND gate, and the counter may be reset from the output of that gate. When approximating segments, there is a possibility for using the same decoding gates for generating straight lines with different slopes.

On starting the linearization the segment-limit decoder and storage units enable the stop line and the gates realizing the $P_1/Q_1 = L_1/K_1$ fractional number.

During a segment of zero slope an inhibiting signal is fed to the stop and add lines, and the pulse-rate is not changing. When generating straight lines of positive slope the segment-limit decoder feeds an open signal to the add line.

Application of equations (1) and (8) of the transducer provide an opportunity to receive symmetrical error curve. In a polygon giving a symmetrical error curve, every straight line having positive slope has its counterpart of negative slope, therefore the number of fractional numbers to be realized decreases by half. The designer must decide how many pulses correspond to the full scale deviation. That decision needs a careful calculation.

The 0.1% accuracy required gives the minimum number of pulses (N_{\min}) and conforms to the full scale deviation of 100 pulses. The maximum permissible frequency of the clock generator, and the maximum average time of the input signal determines the value of N_{\max} . The requirement for engineering units determines the possible N values between N_{\max} and N_{\min} .

One may select the integral multiple of these values. The predetermined count of preset counter 1, which forms the average periodic time, and the length of the counter is influenced by N . With the increasing of N , the setting of the value of the preset counter will be more accurate. According to the authors' experiences it is not advisable to select a pulse number less than $N = 4000$ for an accuracy of 0.1%.

The maximum clock frequency limits the decrease of measuring period at a given error of linearity. The maximum frequency of the clock is determined by the dynamic performance of the applied integrated circuits. When using a high clock frequency the noise problems of the system are pronounced.

The authors' examinations regarding the stability of clock frequency proved that 100 parts in 10^6 stability causes only a negligible error.

3.4. Setting

The calculations of the preset values of the PR1 and PR3 counters are performed according to the methods applied with average periodic time meters. The programming of the linearizer is performed so that after determining the segment limits the gates consisting of the individual fractional numbers are connected.

Table 1 represents the maximum error as the function of span, the number of segments and the number of fractional numbers to be realized as well as the method of calculation of the fractional numbers with regard to the parameters of the density to be measured. The tolerance of d_0 parameters of the transducer involves the modification of the preset values of PR1 and PR3 counters, and the change of h_{\max} . The tolerance of the T_0 parameter modifies only the preset of PR1 and PR3 counters.

When d_0 parameter alters, the value of P_n/Q_n alters too according to h_{max} . Our calculations and measurements have shown that the influence of d_0 tolerance exerted on the fractional numbers is quite negligible, there is no need for modifying the fractional numbers.

3.5. Main Characteristics of the DENSITON Instrument

Input signal: 1 Vp.p.–10 Vp.p., sinusoid or square signal.

Density range (depending on types): 0–120 mg/cm³ and 0.5–1.5 g/cm³

LIQUIDS	GASES	TYPE
0.1 g/cm ³	10 mg/cm ³	AR-102-1 or DR 104-1
0.2 g/cm ³	20 mg/cm ³	AR-102-1 or DR 104-1
0.4 g/cm ³	40 mg/cm ³	AR-102-2 or DR 104-2

Output signal: AR-102 analogue current 4–20 mA
0–10 mA
0–5 mA
10–50 mA

DR-104 digital display, 4 digits b.c.d. code
b.c.d. code output
impulse number output

Accuracy: 0.2% (±1 digit with Type DR-104) for span.

Construction: SN74N integrated circuits
silicon transistors.

Weight: 4.5 kg.

Consumption: 20 VA.

4. Conclusions

The following conclusions can be stated:

- (i) Linearization is accomplished according to the error curve of the output signal of the sensor.

- (ii) The error curve is approximated with a straight line polygon.
- (iii) The straight lines are accomplished by inhibiting pulses from or adding pulses to the clock signal.
- (iv) A symmetrical error curve has to be used to decrease the cost of components of the linearizer.
- (v) The 10% non-linearity can be decreased to 0.1% using the linearizer based on the principle described.

5. Acknowledgments

Grateful acknowledgments are due to Mr. T. Boromisza and Mrs. F. Urbán for contributing to the completion of practical design.

6. References

1. Collins, G. B., 'A survey of digital instrumentation and computer interface methods and development', Conference on Industrial Measurement Techniques for On-line Computers. I.E.E. Conference Publication No. 43, 11–13 June 1968.
2. 'Mass Flow Computer Type M.F.C.-1', Electronic Flow-Meters Ltd. Manual 1969.
3. Linearizer Unit Type LU 1965. Data Sheet No. IC69. The Solartron Electronic Group Ltd.
4. Kollataj, J. H. and Harkonen, T., 'Linearizing sensor signals digitally', *Electronics*, 41, 4th March 1968, p. 112.

Manuscript first received by the Institution on 16th June 1969 and in final form on 2nd June 1970. (Paper No. 1343/IC30).

© The Institution of Electronic and Radio Engineers, 1970

A Laguerre Series Approximation to the Ideal Gaussian Filter

By

N. B. JONES,
B.Sc., M.Eng., D.Phil.†

It is shown that the ideal Gaussian frequency response can be more accurately approximated by using a Laguerre series approach instead of a Taylor series approach involving the same number of components. The resulting filters are physically realizable and their responses compare favourably with those of the conventional approximations to Gaussian filters.

1. Introduction

If the modulus response of a filter can be constrained to have a true Gaussian shape as shown in Fig. 1, then the impulse response also has a Gaussian shape as shown in Fig. 2.

The Gaussian filter is unique in that the shape is preserved on transformation from the frequency domain to the time domain. Such filters are of considerable theoretical interest, for example, in the calculation of correlation functions¹ and physically realizable approximations to them have practical application. It is, for example, desirable that in pulse circuits the step response be rapid but have small or non-existent overshoots and Gaussian approximants have these properties.

Attempts to synthesize Gaussian filters started when it was shown² that a set of equal lags in series (a binomial filter) tends to have a Gaussian amplitude characteristic as the number of stages becomes large. Unfortunately, the convergence to the ideal is slow and the delay becomes excessive when the filter is of high order.

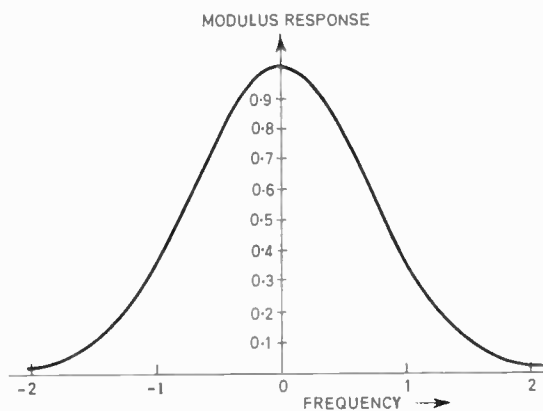


Fig. 1. Ideal modulus response as a function of angular frequency.

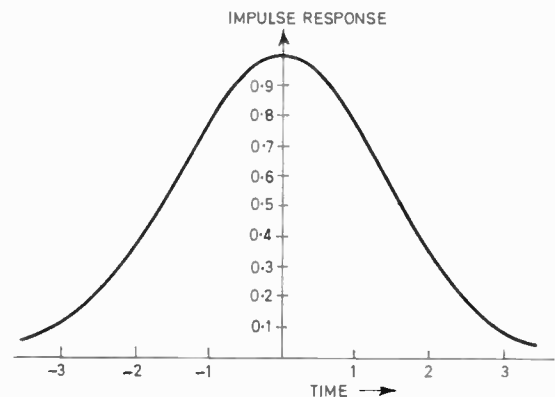


Fig. 2. Ideal impulse response as a function of time.

In 1959 Dishal³ approximated to the inverse square of the modulus function of the ideal by means of a truncated Taylor series. In this paper, however, the error weighting is not concentrated at the origin as with a Taylor series, but is more evenly distributed in the frequency range of interest. It is shown that both the frequency response and step response are in fact nearer to the ideal for the same number of components when a truncated Laguerre series is used as the alternative to the Taylor series. In this paper the Taylor and Laguerre types of filters are compared.

2. The Ideal Gaussian Characteristics

The ideal frequency response is taken to be $\exp(-\omega^2)$ and is shown plotted in Fig. 1 where ω is the normalized angular frequency. This function is unrealizable because there is no phase shift to any real frequency.

By transforming it can be shown (Appendix 1) that the normalized impulse response is

$$f(\tau) = \exp\{-0.721(\omega_{3dB}\tau)^2\}$$

which is shown plotted in Fig. 2. Here ω_{3dB} is the bandwidth and τ is the time.

Integration of the above impulse response gives a normalized step response of $\{\frac{1}{2} + \frac{1}{2} \operatorname{erf}(0.8491\omega_{3dB}\tau)\}$.

† Applied Sciences Laboratory, University of Sussex, Falmer, Brighton BN1 9QT.

3. Approximation to the Ideal

The approximation to the squared modulus curve $[\exp(-\omega^2)]^2$ for the low-pass case is to be of the form $1/P_n(\omega^2)$ where P_n is a polynomial of degree n . Once the coefficients of $P_n(\omega^2)$ have been found, the transfer function, $1/G_n(s)$, is available by taking factors of the form

$$\frac{1}{P_n(\omega^2)} = \frac{1}{G_n(j\omega)} \cdot \frac{1}{G_n(-j\omega)}$$

Here $G_n(s)$ is again a polynomial of degree n , but with the restriction that all the roots lie in the left half of the s -plane.

3.1. Taylor Series Approximants

In this case $P_n(\omega^2)$ is taken to be a truncation of the Taylor series expansion of $[\exp(\omega^2)]^2$ after $n+1$ terms, so that, for example, if $n = 4$ and ω^2 is replaced by x ,

$$P_4(x) = 1 + 2x + \frac{2^2}{2!} x^2 + \frac{2^3}{3!} x^3 + \frac{2^4}{4!} x^4.$$

By putting $-s^2 = x$ and factorizing, the roots of $G_n(s)$, which are also the poles of the transfer function, have been computed to give the following list:

Order (n)	Pole position
2	$-0.777 \pm j0.322$
3	-0.893 $-0.813 \pm j0.556$
4	$-0.958 \pm j0.232$ $-0.835 \pm j0.750$
5	-1.045 $-1.002 \pm j0.423$ $-0.851 \pm j0.919$

These filters are obviously all physically realizable.

3.2. Laguerre Series Approximants

In this case, in order to produce a more evenly distributed error, $[\exp(\omega^2)]^2$ is approximated by a finite series of Laguerre polynomials, namely

$$\exp(2x) \approx \sum_{m=0}^n a_m L_m(\alpha x) \quad \dots\dots(1)$$

where α is a scaling factor which is to be chosen and L_m is the m th order Laguerre polynomial.⁴

For uniform absolute convergence⁵

$$\int_0^{\infty} \exp(-\alpha x) [\exp(2x)]^2 dx < \infty$$

i.e. $\alpha > 4$.

By taking $\alpha = 7$ it can be demonstrated that the terms of the series decrease by more than a factor of two at each step and the coefficient integral,⁵ given by

$$a_m = \int_0^{\infty} L_m(x) \exp(-5x/7) dx$$

are simple to evaluate. The derivation of this coefficient formula is given in Appendix 2.

The first six coefficients calculated thus are,

- $a_0 = 1.4$
- $a_1 = -0.56$
- $a_2 = 0.244$
- $a_3 = -0.0896$
- $a_4 = 0.03584$
- $a_5 = -0.014336$

Alternatively these coefficients could have been evaluated from the formula $1.4 \times (-0.4)^m$ which results from Head's equation.⁶

$$a_m = \frac{(\alpha - \frac{1}{2}\beta)^m}{(\alpha + \frac{1}{2}\beta)^{m+1}} \quad \text{with } \alpha = \frac{3}{14} \quad \text{and } \beta = 1.$$

Hence the inverse modulus squared polynomials can be constructed from the series (1) and are as follows:

- $Q_2(x) = 1.064 + 0.784x + 5.488x^2$
- $Q_3(x) = 0.9745 + 2.6635x - 1.09025x^2 + 5.1164x^3$
- $Q_4(x) = 1.01034 + 1.66x + 4.178x^2 - 3.076x^3 + 3.587x^4$
- $Q_5(x) = 0.996 + 2.16174x + 0.666x^2 + 5.116x^3 - 3.586x^4 + 2.0078x^5$

These give the following set of poles for the transfer functions when factorized as before;

Order (n)	Pole position
2	$-0.5058 \pm j0.4295$
3	-0.5356 $-0.5315 \pm j0.7297$
4	$-0.5687 \pm j0.3174$ $-0.5475 \pm j0.9756$
5	-0.5989 $-0.5959 \pm j0.5720$ $-0.5588 \pm j1.1885$

These filters are also realizable.

3.3. Comparison of the two Approximants

Figures 3 and 4 show respectively the modulus responses and the step responses of the two sets of filters. In each case the existing Taylor approximant is marked T and the new Laguerre equivalent is marked L.

For the modulus curves of Fig. 3, the ideal Gaussian response is marked G and it is obvious that in all cases the Laguerre approximant to it is considerably superior to the Taylor approximant at all frequencies. It is also evident that the ideal curve does not go to zero modulus asymptotically whereas all real responses must do.

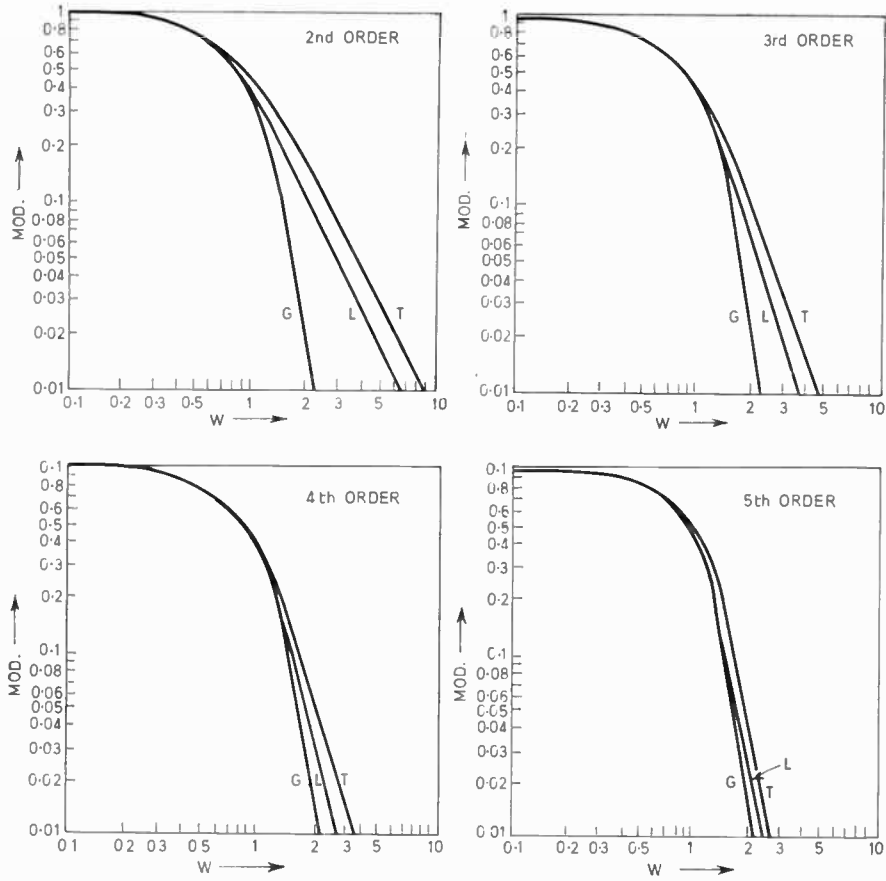


Fig. 3. The ideal modulus response compared with those of the approximants for various orders.

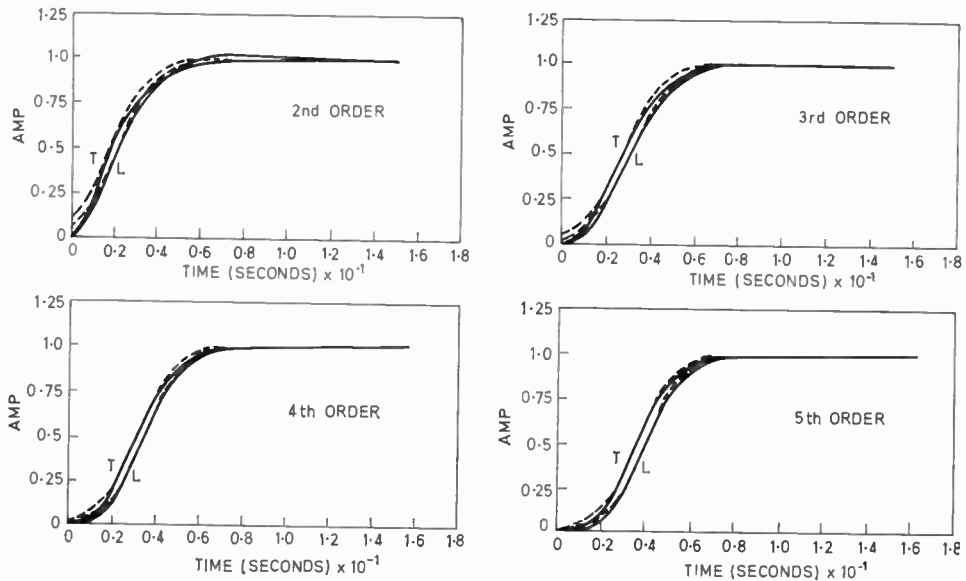


Fig. 4. The ideal step responses compared with those of the approximants for various orders.

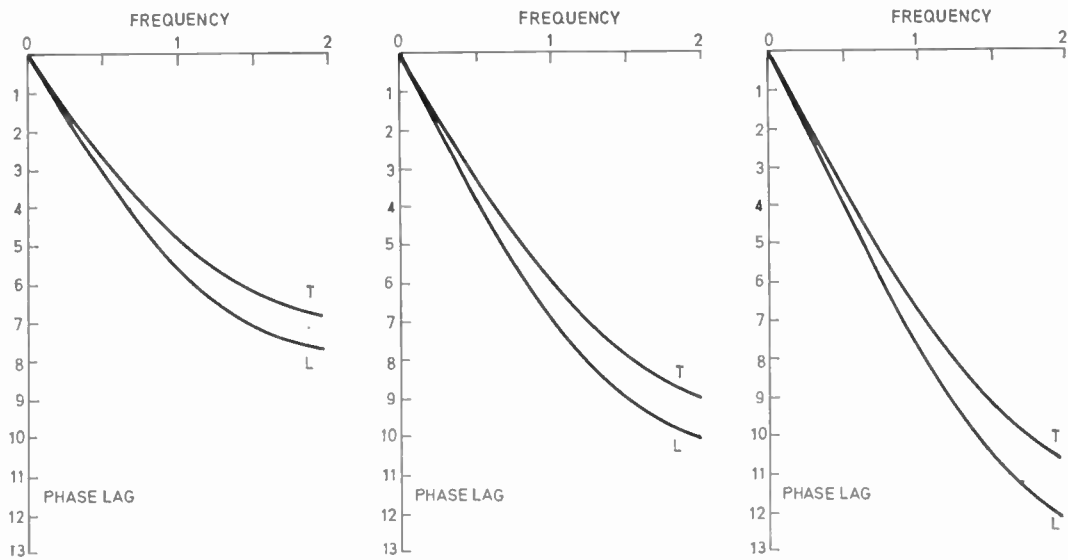


Fig. 5. The phase responses of the approximants.

The step responses of Fig. 4 are marked slightly differently in that the ideal response (shown dotted) has been matched at the asymptotes and at the half-way point on the responses of both the Taylor and Laguerre approximants. This is in order to facilitate the comparison of the shapes without consideration of the delay.

Here again the Laguerre-Gaussian filters are an improvement on the Taylor-Gaussian filters, although it is not as obvious as in the frequency domain.

It was stated by Dishal³ that the phase response of the original approximants to the ideal Gaussian filter were almost linear, that is the group delays

approximate to constants. Comparison of the phase responses are shown in Fig. 5 and it can be seen that the new Laguerre approximants are also accurately constant group delay filters in the pass bands.

This is not an altogether surprising result when the pole distribution is considered since the Laguerre-Gaussian filters have poles which are almost equally spaced and in a line parallel to the imaginary axis. Filters with equally-spaced poles in a line parallel to the imaginary axis are well known approximants to a pure group delay.⁷ The pole distributions are shown in Fig. 6.

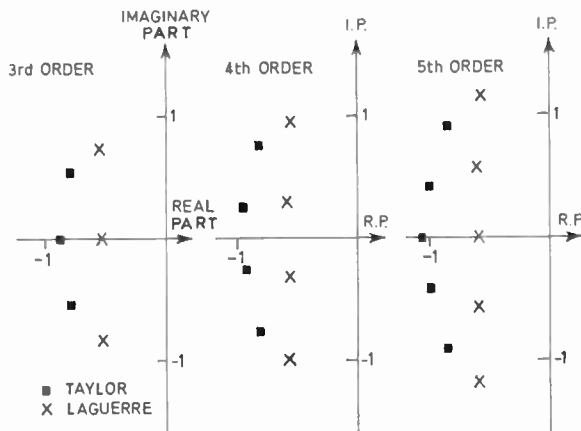


Fig. 6. The pole distribution of the approximants.

4. References

1. Goldman, H. D., 'On using Gaussian filter characteristics for approximate auto-correlation calculations', *I.E.E. Transactions on Communications Techniques*, COM-14, No. 6, pp. 865-6, 1966.
2. Valley, G. and Wallman, H., 'Vacuum Tube Amplifiers', pp. 721-7. (McGraw-Hill, New York, 1948).
3. Dishal, M., 'Gaussian-response filter design', *Electrical Communications*, 36, No. 1, pp. 3-26, 1959.
4. Abramowitz, M. and Stegun, I. A., 'Handbook of Mathematical Functions', p. 297. (Dover, New York.) 1964.
5. Sansone, G., 'Orthogonal Functions', p. 297. (Interscience, New York, 1959).
6. Head, J. W. and Proctor Wilson, W., 'Laguerre functions: tables and properties', I.E.E. Monograph No. 183R p. 7. June 1956. (*Proc. Instn Elect. Engrs*, 103, Part C, No. 4, pp. 428-40, 1956.)
7. Guillemin, E. A., 'Synthesis of Passive Networks', p. 639. (Wiley, New York, 1957).

5. Appendix 1: To find the impulse and step responses of the ideal Gaussian low-pass filter of bandwidth $\omega_{3\text{dB}}$

If the modulus function is taken as $\exp(-\omega^2)$, the bandwidth can be shown to be $(\frac{1}{2} \log 2)^{\frac{1}{2}}$ and so the filter of bandwidth $\omega_{3\text{dB}}$ has a modulus function of

$$\exp \left[- \left(\frac{\omega}{\omega_{3\text{dB}}} \right)^2 0.3466 \right]$$

Therefore the impulse response is given by

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \exp \left[- \left(\frac{\omega}{\omega_{3\text{dB}}} \right)^2 0.3466 \right] \exp(j\omega\tau) d\omega$$

$$= 1.201\omega_{3\text{dB}} \exp(-0.721(\omega_{3\text{dB}}\tau)^2)$$

The step response therefore is given by

$$1.201\omega_{3\text{dB}} \int_{-\infty}^{\tau} \exp[-0.721(\omega_{3\text{dB}}\tau)^2] d\tau$$

$$= 2.5076 \left\{ \frac{1}{2} + \frac{1}{2} \operatorname{erf}(0.8491\omega_{3\text{dB}}\tau) \right\}.$$

When normalized so that the maximum value in both cases is taken to be unity then:

$$\begin{aligned} \text{impulse response} &= f(\tau) \\ &= \exp[-0.721(\omega_{3\text{dB}}\tau)^2] \\ \text{step response} &= g(\tau) \\ &= \frac{1}{2} \{ 1 + \operatorname{erf}(0.8491\omega_{3\text{dB}}\tau) \}. \end{aligned}$$

6. Appendix 2: Evaluation of the coefficients of the Laguerre series

The ideal frequency response is taken to be $\exp(-\omega^2)$ and hence, for the low-pass case, it is necessary to find a polynomial approximation in ω^2 to $\exp(2\omega^2)$.

$$\begin{aligned} \text{Let } \exp(2\omega^2) &= \exp(2x) \\ &= \sum_{m=0}^n a_m L_m(\alpha x) + E(x) \end{aligned}$$

The term $E(x)$ is now minimized in a weighted mean square sense, taking as the weighting function $\exp(-\alpha x)$.

Thus,

$$\begin{aligned} [E(x)]^2 \exp(-\alpha x) &= \left[\exp(2x) - \sum_{m=0}^n a_m L_m(\alpha x) \right]^2 \exp(-\alpha x) \\ &= \exp[(4-\alpha)x] - 2 \exp[(2-\alpha)x] \sum_{m=0}^n a_m L_m(\alpha x) + \\ &\quad + \left(\sum_{m=0}^n a_m L_m(\alpha x) \right)^2 \exp(\alpha x). \end{aligned}$$

Let

$$\int_0^{\infty} [E(x)]^2 \exp(-\alpha x) d(\alpha x) = \bar{E}$$

Therefore

$$\begin{aligned} \frac{\partial \bar{E}}{\partial a_k} &= -2 \int_0^{\infty} \exp[(2-\alpha)x] L_k(\alpha x) d(\alpha x) + \\ &\quad + 2 \int_0^{\infty} \exp(-\alpha x) L_k(\alpha x) \sum_{m=0}^n L_m(\alpha x) d(\alpha x) \\ &= -2 \int_0^{\infty} \exp[(2-\alpha)x] L_k(\alpha x) d(\alpha x) + 2 a_k \end{aligned}$$

Since

$$\int_0^{\infty} \exp(-\alpha x) L_m(\alpha x) L_k(\alpha x) d(\alpha x) = \delta_{mk}$$

(δ_{mk} is the Kronecker delta)

Therefore

$$\frac{\partial^2 \bar{E}}{\partial a_k^2} = 2.$$

Hence a minimum value of \bar{E} is achieved when

$$a_k = \int_0^{\infty} \exp[(2-\alpha)x] L_k(\alpha x) d\alpha x.$$

Choosing $\alpha = 7$ gives

$$a_k = \int_0^{\infty} \exp\left(\frac{-5x}{7}\right) L_k(x) dx.$$

Manuscript first received by the Institution on 9th January 1970 and in final form on 23rd March 1970. (Paper No. 1344/CC88)

© The Institution of Electronic and Radio Engineers, 1970

Contributors to the Journal



Professor H. Sutcliffe graduated from the University of Cambridge in 1941 in mechanical sciences. His first post was in radar research and development at the then T.R.E., now R.R.E., Malvern. He entered university teaching in 1951 in the Dundee branch of St. Andrew's University, later joined Smith's Aircraft Instruments in Cheltenham, and re-entered academic life at

Bristol University in 1956. In 1963 he went to Salford University, first as reader and later as professor of electronic engineering. His research interests are mainly in the field of electronic instrumentation and he has published a number of papers principally in this field. He is also the author of a text book on electronics for mechanical engineers.

a lecturer in control engineering at the University of Sussex. He obtained his D.Phil. at Sussex in 1968 on aspects of the identification and synthesis problems.



Mr. Geza Mayer received his Diplom-Engineer in electronic engineering from the Politechnical University of Budapest in 1965, and then joined the Central Measurement Laboratory, Budapest. At the Laboratory his main concern was research and development of industrial electronic measuring equipments, notably the development of a digital linearizer equipment. He also

held an appointment at the Politechnical University of Budapest as an assistant lecturer. In March this year Mr. Mayer joined the Infelcor System Engineering Company in Budapest, as senior project engineer.



Dr. K. F. Knott received his degree of B.Eng. from the University of Liverpool in 1960 and his Ph.D. was awarded by the University of Salford in 1968. From 1960 to 1962 he was an instrument research engineer at A. V. Roe and Co. Ltd., Woodford, Cheshire. After a brief period with Ferranti Limited as a technical author, he received his present appointment as a

lecturer in the Electrical Engineering Department of the University of Salford. Dr. Knott has published many papers and articles in the field of low-frequency noise.



Mr. Laszlo Simonfai qualified in electronic engineering at the Politechnical University of Budapest in 1964, obtaining his Diplom-Engineer. He worked at the Central Measurement Research Laboratory, Budapest on the development of transducers and digital panel meters, and on the digital linearizer programme. Recently he joined the Infelcor System Engineering

Company as a senior project engineer.



Dr. N. B. Jones graduated in electrical engineering from Manchester University in 1962. He spent some time with B.I.C.C. as a graduate apprentice, and then took up a scholarship to study for an M.Sc. degree in electrical engineering at McMaster University, Hamilton, Ontario. After graduating he worked as a research and development engineer with Canadian

Westinghouse Limited. Dr. Jones returned to England and was employed firstly as a tutorial fellow and then as



Mr. Peter Potzy obtained his degree of Diplom-Engineer in electronic engineering in 1963, and completed a post-graduate course on control engineering in 1968. He joined the Central Measurement Research Laboratory, Budapest, Hungary in 1968, and participated in the development of the digital linearizer. Previously Mr. Potzy worked at Vilati on industrial automation, and he

is now with Infelcor System Engineering Company working on computer applications.