

THE BELL SYSTEM

Technical Journal

Volume 51

October 1972

Number 8

D2 CHANNEL BANK

System Aspects	H. H. Henning and J. W. Pan	1641
Per-Channel Equipment	C. L. Maddox and D. K. Thovson	1659
Multiplexing and Coding	C. L. Dammann, L. D. McDaniel, and C. L. Maddox	1675
Digital Functions	A. J. Cirillo and D. K. Thovson	1701
Power Conversion	S. D. Bloom and G. F. Swanson	1713
Physical Design and Introductory Program	D. J. VanSlooten and K. A. Gluckow	1725
Manufacturing and Testing	J. E. D. Batson, Jr., and J. W. Gorman	1743

GENERAL ARTICLES

Optical Power Flow in Multimode Fibers	D. Gloge	1767
Pulse Propagation in a Two-Mode Waveguide	D. Marcuse	1785
Fluctuations of the Power of Coupled Modes	D. Marcuse	1793
Higher-Order Scattering Losses in Dielectric Waveguides	D. Marcuse	1801
Higher-Order Loss Processes and the Loss Penalty of Multimode Operation	D. Marcuse	1819
A Linear Phase Modulator for a Short-Hop Microwave Radio System	S. R. Shah	1837
A Doped Surface Two-Phase CCD	R. H. Krambeck, R. H. Walden, and K. A. Pickar	1849
Almost-Coherent Detection of Phase-Shift-Keyed Signals Using an Injection-Locked Oscillator	M. Eisenberg	1867
Analysis of a Dual Mode Digital Synchronization System Employing Digital Rate-Locked Loops	R. W. Chang	1881
Contributors to This Issue		1913

(Continued on inside back cover)

THE BELL SYSTEM TECHNICAL JOURNAL

ADVISORY BOARD

D. E. PROCKNOW, *President, Western Electric Company*

J. B. FISK, *President, Bell Telephone Laboratories*

W. L. LINDHOLM, *Vice Chairman of the Board,
American Telephone and Telegraph Company*

EDITORIAL COMMITTEE

W. E. DANIELSON, *Chairman*

F. T. ANDREWS, JR.

A. E. JOEL, JR.

S. J. BUCHSBAUM

H. H. LOAR

R. P. CLAGETT

B. E. STRASSER

I. DORROS

D. G. THOMAS

D. GILLETTE

C. R. WILLIAMSON

EDITORIAL STAFF

L. A. HOWARD, JR., *Editor*

R. E. GILLIS, *Associate Editor*

H. M. PURVIANCE, *Production and Illustrations*

F. J. SCHWETJE, *Circulation*

J. W. PAN, *Coordinating Editor of D2 Channel Bank Articles*

THE BELL SYSTEM TECHNICAL JOURNAL is published ten times a year by the American Telephone and Telegraph Company, J. D. deButts, Chairman and Chief Executive Officer, R. D. Lilley, President, J. J. Scanlon, Vice President and Treasurer, R. W. Ehrlich, Secretary. Checks for subscriptions should be made payable to American Telephone and Telegraph Company and should be addressed to the Treasury Department, Room 2312C, 195 Broadway, New York, N. Y. 10007. Subscriptions \$10.00 per year; single copies \$1.25 each. Foreign postage \$1.00 per year; 15 cents per copy. Printed in U.S.A.

THE BELL SYSTEM TECHNICAL JOURNAL

DEVOTED TO THE SCIENTIFIC AND ENGINEERING
ASPECTS OF ELECTRICAL COMMUNICATION

Volume 51

October 1972

Number 8

Copyright © 1972, American Telephone and Telegraph Company. Printed in U.S.A.

D2 Channel Bank:

System Aspects

By H. H. HENNING and J. W. PAN

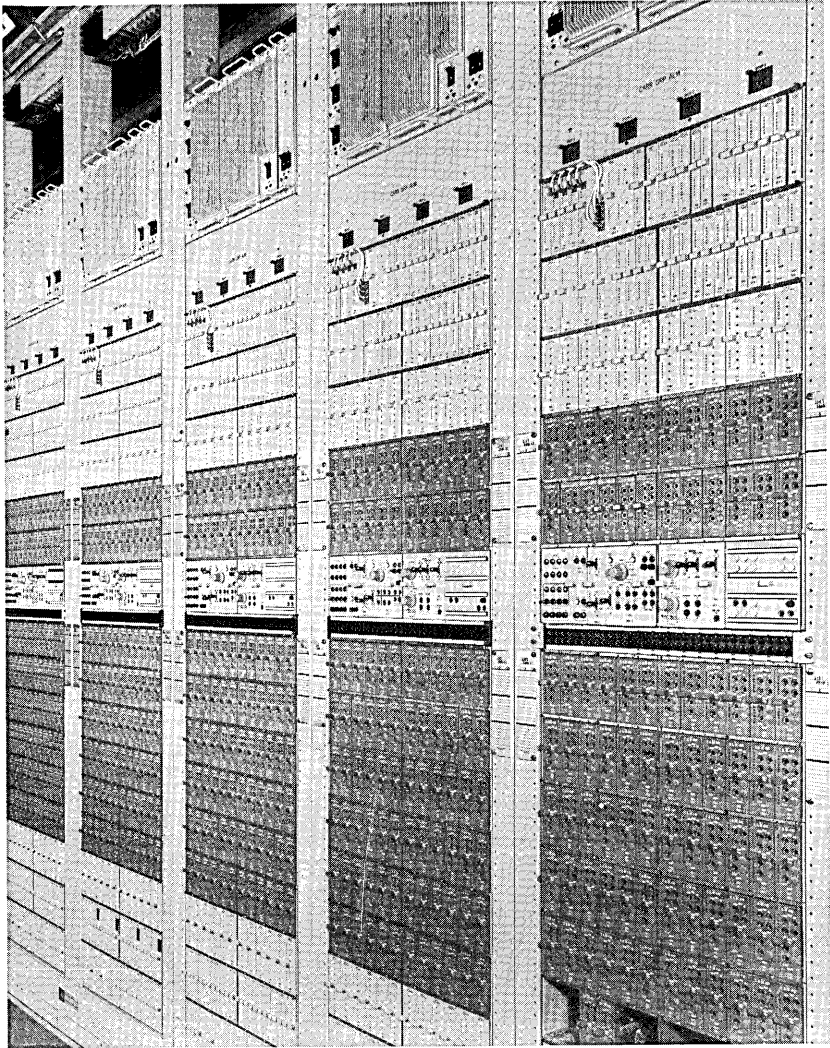
(Manuscript received June 22, 1972)

This is the first of a series of articles describing the D2 Channel Bank—from initial conception, system design, circuit development, physical design, through manufacture, installation and service. Our objective is to provide a complete story of how one product progressed from identification of need through various phases leading from early planning to operating company application.

In this introductory article, the motivation for undertaking the development is pointed out, and the reasons for the choice of the various system parameters are discussed. The D1 Channel Bank, which was designed for exchange application, was the pioneer in digital channel banks. As a second generation channel bank developed for toll application, D2 could be expected to show significant advances over its predecessor. Of the many possible improvements, some would result in incompatibilities with portions of the existing plant. In such cases, engineering judgements were necessary to determine which of these were warranted by the performance improvements they allowed. This article documents the historical evolution of the D2 Channel Bank system parameters.

I. INTRODUCTION

The first digital transmission system used for commercial telephone service was introduced by the Bell System in 1962. This system consists



D2 Channel Bank

of the D1 Channel Bank and the T1 Repeated Line.¹ The D1 Bank uses pulse code modulation to convert 24 voice-frequency signals and the associated signaling information into a 1.544 megabits per second digital stream for transmission over the T1 Repeated Line. Because it has met its performance and cost expectations, the D1/T1 system, which is also known as "the T-Carrier," has been favorably received

by the telephone companies. Today it is the fastest growing carrier facility in the Bell System supplying nearly a million voice channels throughout the country. Its success is a direct result of one of the characteristics of digital systems—low terminal cost.

The success of D1/T1 has stimulated planning of a digital communications network. A fully digital network includes both digital switching and digital transmission. Most of the effort since 1962, however, has been directed towards the development of transmission systems for this network.^{2,3} In this network there will be a hierarchy of digital transmission facilities that provide for long-haul transmission at high bit rates,⁴ digital terminals that convert a variety of signals into a suitable digital form, and digital multiplexers that can derive several smaller capacity digital facilities from a large capacity facility. The existing T1 lines will become part of this hierarchy.

Since the D1 Channel Bank was designed primarily for short inter-office trunks, it does not have all the transmission performance and operating features required for toll service. As a result, there was appreciable motivation for the development of a new digital terminal, the D2 Channel Bank, which is designed to provide economical voice trunks for intertoll service, and which can also be used for exchange trunks. This article describes the system aspects of the D2 Channel Bank. Companion articles describe the various circuits of the D2 Bank, production, installation, and continued improvement after initial service.

II. SYSTEM CHARACTERISTICS

Many objectives and constraints play a significant role in the system design of the D2 Channel Bank. First, a channel bank suitable for intertoll use must exhibit a performance level superior to that of the D1 Bank. This is because toll calls may include many digital trunks interconnected by switching machines which at present can handle signals only in voice-frequency form. In such situations, quantizing noise will accumulate because of repeated analog-digital and digital-analog conversions. Second, in the digital network envisioned for the future, the trunks may be switched in digital form. This means that a signal converted into digital form by one channel bank is expected to be reconstructed by any other channel bank. A high degree of standardization and uniformity is thus required of all such channel banks. Third, the D2 Channel Bank must be able to utilize the T1 line which is an existing transmission facility in the planned digital hierarchy.

Compatibility considerations with the D1 Bank and with other

domestic and foreign digital channel banks in existence or in development have also influenced the system design of D2. Incompatibility will affect the problem of digital interconnection between different telephone administrations in the future. It will also affect future channel bank designs. To be fully compatible, two channel banks must have these same attributes:

- (i) number of voice channels
- (ii) sampling rate
- (iii) companding law
- (iv) code format
- (v) overload point
- (vi) signaling format
- (vii) framing format
- (viii) output bit rate.

If any one of these attributes is different for two channel banks, then digital processing becomes necessary before they can be interconnected digitally. Each attribute requires varying degrees of digital processing to convert from one standard to another.

It was recognized in the system design of D2 that it would not be feasible to strive for complete compatibility with existing digital channel banks for performance reasons, or with future banks because international standards were not yet determined. The choices made in the design of the D2 Bank were such that complex digital processing could be avoided, but simple processing would be permissible for possible interconnection. The use of simple processing can readily provide changes in the code format, signaling format, and framing format. The number of voice channels and the output bit rate can be changed easily and efficiently if the numbers used by two channel banks form simple fractions. This allows an integral number of channel banks on either side to be interconnected. Both companding law and overload point can be changed by digital processing of low complexity. The most difficult parameter to convert is the sampling rate. To change this parameter, interpolation between samples is necessary. The complexity is equivalent to digital filtering.

As a result of these considerations, the number of channels, the sampling rate, and the output bit rate for D2 are chosen to match those of the D1 Channel Bank. This will allow the use of the existing T1 digital transmission line to interconnect D2 Banks and, with simple to moderate digital processing, allow digital interconnection between D2 and other channel banks with 8-kHz sampling.

All other attributes of D2 are not compatible with D1 primarily for reasons of improved performance. The most important differences are the number of digits used per coded sample and the companding law. These differences in turn cause other attributes, such as signaling format and framing format, to be different.

2.1 *Consequences of Eight Digits Per Coded Sample*

As mentioned earlier, in order to meet noise and distortion requirements for toll service, eight-digit PCM is used as compared to seven-digit PCM used in D1. But to increase the number of digits per coded sample and still maintain the same sampling rate and output rate, either the number of voice channels must be reduced or the portion of output bit rate devoted to signaling must be reduced. The latter approach is taken because a reduction in signaling rate to one-sixth of that in D1 could be tolerated without sacrificing the capability of the D2 Bank to operate with all the signaling systems that are presently handled by the D1 Bank. In this way, close to eight-digit PCM performance is achieved. The actual format used in D2 is to code each sample into eight-digit PCM in five out of six frames and into seven-digit PCM in the sixth frame when the eighth digit is used for signaling. The resultant quantization noise is about 2 dB above that of full eight-digit PCM.

The signaling capacity derived in this way is adequate for all signaling systems presently served by the D1 Channel Bank. This includes, for example, dial pulse and reverive pulse signaling. Switching systems in the future are likely to convert to Common Channel Interoffice Signaling (CCIS) where signaling information for a collection of voice channels are sent as a separate digital stream. With CCIS, the signaling circuit packs are simply removed, and the performance of full eight-digit coding can be realized.

2.2 *Frame Format*

The format of D2 allows for signaling in the present plant and future CCIS type signaling. This is accomplished as follows (see Fig. 1). At the frame rate of 8 kHz, there are 193 binary digits per frame as determined by the T1 transmission line. Each of the 24 voice channels occupies a time slot of eight digits. This totals 192 digits leaving the 193rd digit available for framing. Because in every sixth frame the "eight" digit in each time slot is devoted to signaling, and furthermore, for some switching systems two signaling paths are required, it is necessary to identify a super frame of 12 frames of which the sixth and twelfth frames contain the two signaling paths. To accomplish this identification and still allow

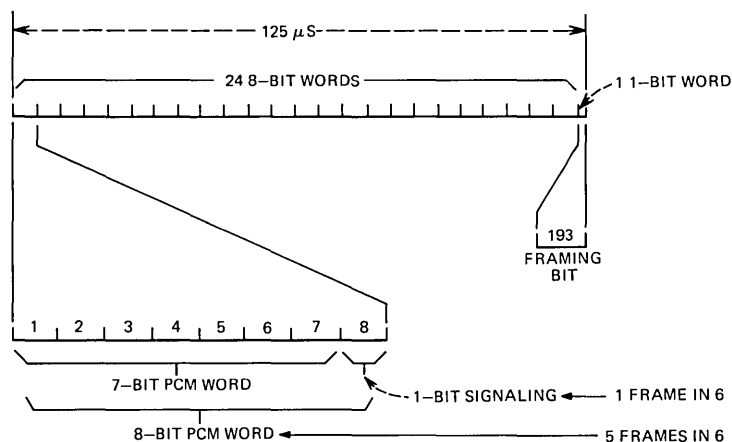


Fig. 1—D2 bit stream format.

rapid synchronization of the receiving framing circuitry, the frames are divided into odd and even frames. In the odd frames, the 193rd digit is made to alternate between 0 and 1. This allows the framing circuit to lock on and maintain synchronism. In the even frames, the 193rd digit is made to follow a 000111000111 . . . sequence. This identifies the sixth and twelfth frames as those that follow an 01 transition or 10 transition of this digit, respectively. When CCIS signaling is used, there is no need to identify the sixth and twelfth frames. The 193rd digits in even frames then become available for CCIS signaling.

2.3. Companding Law and Code Format

The companding law used in the original D1 Channel Bank is based on the nonlinear properties of diodes.* The biasing of the diodes is chosen to approximate a $\mu = 100$ companding law.¹ Careful selection of diodes and stringent temperature control of the diode environments are necessary to maintain matching of the compressor and expander characteristics. To meet the requirement that any two D2 Banks may be interconnected at random under control of a digital switch, even more careful selection of diodes and tighter temperature control would be necessary if diode companding was used. For this reason, it is unlikely that the digital signals from D1 Banks that are based on diode companding can be switched to another similar bank with any assurance of

* Plug-in units will be available for D1 to make its companding law compatible with that of D2.

performance. At the time D2 was developed, techniques were available to achieve nonlinear coding directly. Among these were (i) the techniques of coding linearly and processing the result digitally to produce a nonlinear code, (ii) the technique of using a nonlinear feedback network in a digit-by-digit coding process, and (iii) the technique of stage-by-stage coding where each stage has one digit output and one nonlinear residue output. These techniques are capable of producing different classes of companding laws. In order to provide the greatest flexibility for future channel bank development, the choice was narrowed to two laws that can be implemented by any of the above techniques. The two laws are (i) the 13-segment approximation to the A-law that was actively being considered at that time as the standard for Western Europe, and (ii) the 15-segment approximation to the $\mu = 255$ law.⁵ The $\mu = 255$ law (Table I) was chosen for D2 because it promised better idle circuit noise and crosstalk performance.⁶ Conversion between the two laws is relatively simple.⁷

The most important property of either law is the digitally linearizable property. This allows simple digital processing to convert between companded PCM code words to linear PCM code words. In the linear

TABLE I—CODER THRESHOLD LEVELS AND STEP SIZES
FOR 15-SEGMENT, $\mu = 255$ CODING LAW*

Segment End Points	Step Size	Segment Number
$x_0 = 0$	$\Delta x_0 = 1$	1
$x_1 = 1$	$\Delta x_1 = 2$	
$x_{16} = 31$	$\Delta x_2 = 4$	2
$x_{32} = 95$	$\Delta x_3 = 8$	3
$x_{48} = 223$	$\Delta x_4 = 16$	4
$x_{64} = 479$	$\Delta x_5 = 32$	5
$x_{80} = 991$	$\Delta x_6 = 64$	6
$x_{96} = 2015$	$\Delta x_7 = 128$	7
$x_{112} = 4063$	$\Delta x_8 = 256$	8
$x_{128} = 8159$		

* Points shown are for positive quadrant only; the negative quadrant is symmetrical. With the exception of the first, each segment contains 16 steps of equal step size Δx_n . The output levels are always midway between the threshold points. The table is normalized to 8159 so that all values are represented as integer numbers.

form, signal processing such as gain change, echo suppression, and signaling can be done digitally. Since some processing is expected on international connections or other long circuits, this processing point is then the logical interface between μ -law and A-law regions.

To achieve the $\mu = 255$ companded coding, the stage-by-stage approach was chosen primarily because this approach had the greatest promise of achieving the necessary speed to code 96 voice channels on a time-shared basis.⁸

III. SYSTEM OBJECTIVES

After the major system characteristics have been determined, a set of objectives can be formulated based on reasonable degradations from an ideal channel bank. In an ideal channel bank, the sources of degradations are (*i*) idle circuit noise, (*ii*) quantizing noise, and (*iii*) overload noise. These correspond to small signal, medium signal, and large signal characteristics of a PCM system. In a practical system, idle circuit noise and quantizing noise are expected to be worse than the ideal. Overload noise is not expected to differ from the theoretical value. In addition to the degradation attributed to the quantization process, there is also the degradation due to the sampling process. Some foldover noise (frequency aliasing) is expected due to incomplete removal of signal energy above the half sampling frequency of 4 kHz. Filter characteristics are designed with sufficient out-of-band attenuation so that this source of noise is negligible.

The voice-frequency transmission objectives set for D2 are listed in Table II.

The objectives on idle channel noise and signal-to-distortion can be interpreted as follows. For most of the talker volumes, the objective is

TABLE II—VOICE FREQUENCY TRANSMISSION
OBJECTIVES FOR D2

Overall Idle Channel Noise	<23 dBm0
Interchannel Crosstalk Loss	>65 dB
Signal-to-Distortion Ratio	
C-Message Weighting	
Sine Wave Input	
+3 dBm0 to -30 dBm0	33 dB
-40 dBm0	27 dB
Overload Point	+3 dBm0
Frequency Characteristics	
300 Hz to 3,000 Hz	± 0.25 dB

within 3 dB of the theoretically achievable performance. For very small signals such as those produced by noise and crosstalk, the objective is within 6 dB of theoretical. Overload point can be set arbitrarily, in the sense that every dB of overload results in a dB loss in low signal performance. As in D1, the desired overload point for D2 is near the peak of a +3-dBm0 sine wave.

The frequency characteristics of a channel bank is determined primarily by the filters. Toll transmission objectives for D2 require more complex channel filters than to D1.

IV. SYSTEM DESCRIPTION

The overall block diagram of the D2 Channel Bank is shown in Fig. 2. Per channel equipment includes the channel unit, the transmitting and receiving filters and the multiplex and demultiplex gates. Common equipment performs analog-to-digital conversion at the transmitting terminal and the inverse operation of digital-to-analog conversion at the receiving terminal.

The total cost of a channel bank can be divided into per channel equipment and common equipment. Because of the additional operational and performance requirements necessary for toll operation, the per channel equipment is expected to be more elaborate and thus more expensive than that for D1. To offset this increased cost, the cost of common equipment can be reduced by sharing it over more channels. The most complex common circuits are the coder and decoder, which are shared by all 96 channels. Some common equipment is shared by fewer channels. Engineering judgment was exercised to decide between economy of operation, vulnerability to failure, and effect on working channels during repair whenever it was decided to share common equipment.

The equipment configuration of the D2 Bank also entered into the considerations concerning the maximum number of channels that can be processed on a timeshared basis by common equipment. A major portion of the bay space is taken up by the toll trunk channel units which contain a considerable amount of bulky components such as jacks and relays. Since it was decided to house all the equipment associated with the D2 Bank (including power converter and carrier group alarm) in a factory-assembled bay, 96 channels was the maximum number that could be accommodated on an 11-foot, 6-inch bay. Channel banks for exchange applications with less voluminous channel units could be designed to accommodate a larger number of channels in the same space.

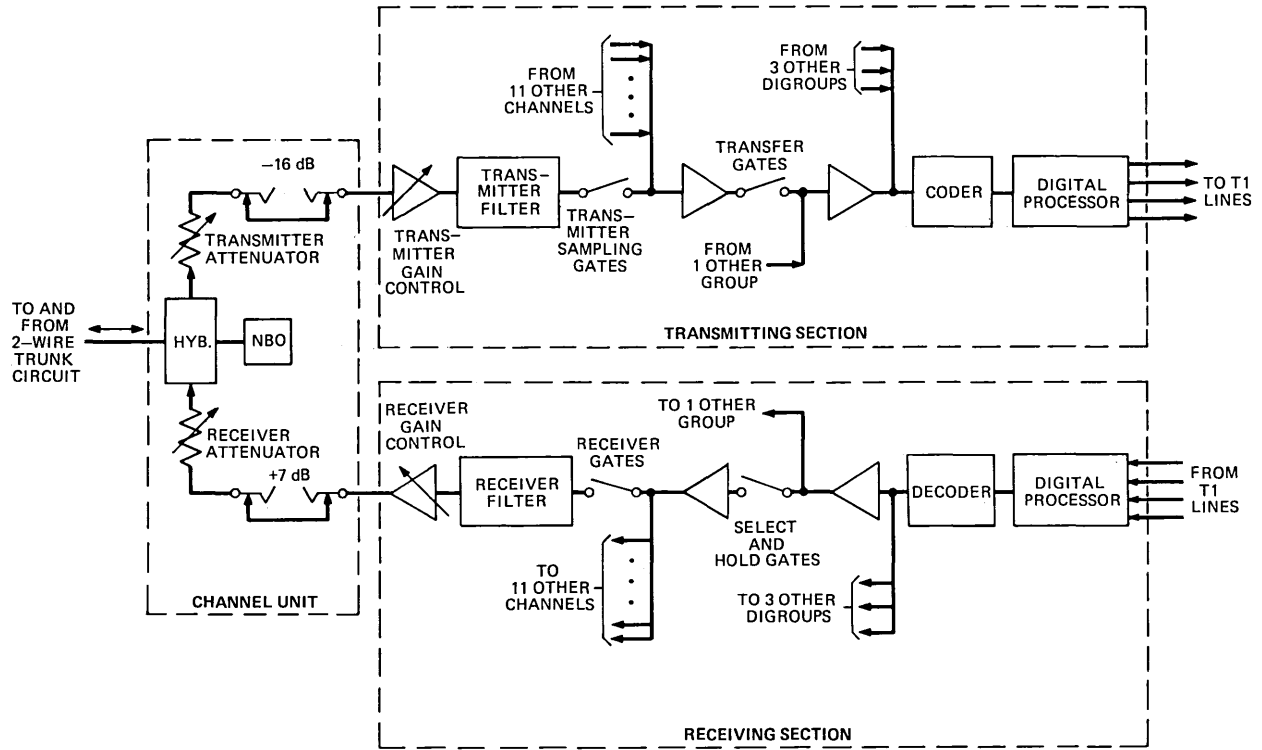


Fig. 2—Block diagram of D2 bank, showing the stages of time-division multiplexing into groups (12 channels), digroups (24 channels), and 96-channel signals.

4.1 *Channel Unit*

The interface equipment between trunk circuits, which are the terminations in a switching machine, and the channel bank, which terminates transmission systems, is contained in the channel unit. The basic functions performed are: (i) converting a 2-wire voice circuit into standard level 4-wire circuit, and (ii) converting the signaling states to digital form. These functions enable the switching machine to treat the voice channels that are derived by PCM techniques as if they were wire pairs. Each channel unit serves one channel, hence a fully equipped D2 Bank contains 96 channel units.

4.2 *Transmitting Terminal*

The transmitting terminal performs analog-to-digital conversion. This includes sampling, multiplexing, coding and digital processing. First, the signal from each channel unit is amplified and bandlimited by a lowpass filter to reject all signal frequencies above 4 kHz. The gain of the amplifier is adjusted to establish the correct level for coding for each channel. The bandlimited signal is then sampled 8,000 times a second by the sampling gate associated with this channel.

Since a single coder is time-shared to code the signals from 96 channels, it is necessary to time multiplex samples from these 96 channels. To achieve this in one step places severe crosstalk and noise limits on the sampling gates. A two-stage multiplexing scheme eases these problems and results in only moderate increase in complexity, since each gate in the second stage of multiplexing is shared by many channels. The number of channels for the first stage of multiplexing is chosen as 12. This number agrees with that of the D1 Bank and most frequency division systems. The eight PAM buses, each containing a *group* of 12 channels, are multiplexed together in the second stage. The resultant multiplex of 96 channels each occupying a time slot of 1.3 microseconds is converted sequentially to digital form by the coder.

The result of coding is a series of pulse-code-modulated digital signals called PCM words. These code words are processed by a coder output processor prior to its application to a digital transmission line. The purpose of digital processing is threefold: (i) the parallel digital output must be converted to serial form, (ii) signaling digits must be inserted at the appropriate times, and (iii) the digital signal must be split into four streams of 1.544 Mb/s each in bipolar format suitable for transmission over T1 lines.

4.3 *Receiving Terminal*

The receiving terminal performs the inverse operation of the transmitting terminal. It has a decoder input processor to bring the incoming digits into a parallel form for decoding by a single decoder, and it uses a two-stage demultiplexing process to distribute the decoded samples to the 96 low-pass filters which reconstruct the original signals.

Whereas it is relatively simple to split the output of the coder into four bit streams, combining four incoming bit streams for decoding by a single decoder is not as simple. The reason is that the four signals can originate from four different channel banks, each with a slightly different clock rate. Thus timing and framing signals are recovered from the four signals individually. Whereas the use of four separate decoders will result in the same apparent complexity as the use of a single decoder which requires digital circuits to combine four asynchronous signals, the single decoder approach results in less analog circuitry. Digital circuits necessary for single decoder operation are simple to build and are more amenable to integrated electronics technology.

The decoded samples undergo two demultiplexing steps to redistribute the signals to the 96 channels. The procedure is the inverse of the one used in the transmitting terminal. The first step demultiplexes the output of the decoder into eight groups of 12 channels each, and the second step demultiplexes each of the eight signals into individual channels for reconstruction by the receiving lowpass filters. In addition to demultiplexing, the first step also removes the timing jitter imparted to the signal by the process of sharing a single decoder. Since the decoder is shared, a digital code word arriving on one of the four lines must queue up for decoding; the waiting time can be as long as 5 microseconds (one word time slot on the line). This waiting time is absorbed by the select and hold circuit, which acts as an analog store as well as a demultiplexer.

V. FEATURES

One important factor in the system design of the D2 Channel Bank is the provision for simplifying the installation and maintenance procedures including gain adjustments and noise measurements. For these purposes, a test panel is built into the D2 bay. This test panel contains standard signal sources, signal and noise test set interface equipment, and signaling test sets. The test panel provides convenient access to standard, permanently installed central office equipment.

Flexible jacking arrangements are also provided so that measurement of transmission, noise, distortion and crosstalk on all voice channels can be accomplished by one man with no portable test equipment.

5.1 *Digital Signal Generator*

To allow D2 to achieve a consistent correspondence between analog signals and their digital representation, a digital signal generator (DSG) is provided as a unique feature of the test panel. This circuit provides an invariant digital reference level for calibration purposes. The digitally derived signal consists of a repetitive sequence of eight PCM words. This signal is defined to be identical to a PCM signal that would be produced if a 0-dBm sinusoidal test tone of 1 kHz, synchronized with the sampling frequency, were encoded by a perfectly gain-adjusted transmitting terminal. The eight binary words that define the signal are listed in Fig. 3. The signal from the DSG can be inserted into any one of the four incoming PCM lines.

The DSG greatly simplifies the installation and maintenance procedures, including gain adjustments and distortion measurements. Prior to D2, the line-up of transmission systems required at least one man at each of the end terminals. One man would connect a calibrated signal at the transmitting end, while another man would perform

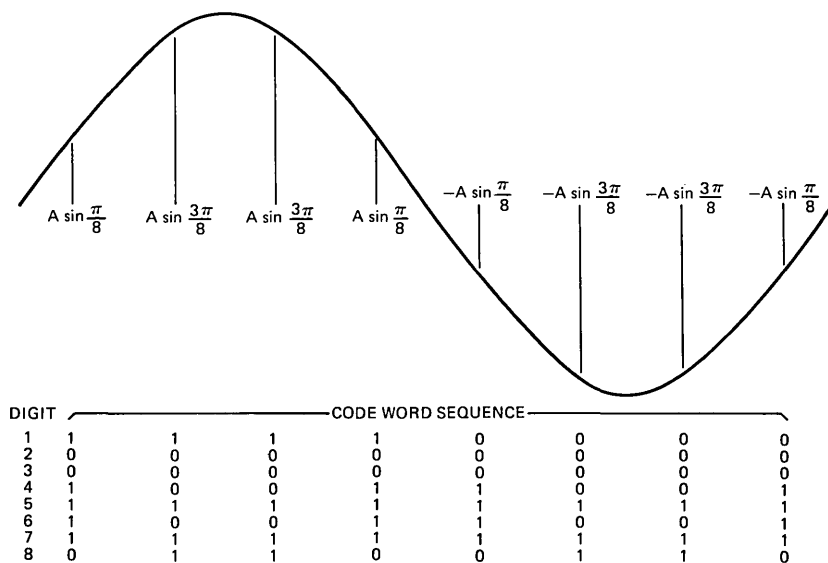


Fig. 3—The digital 1-kHz reference signal.

measurements and adjustments at the receiving end. The D2 Channel Bank takes advantage of the fact that the transmission line is digital. The digital line does not introduce loss or distortion to the coded analog signals. Only the analog circuits of the transmitting terminal and the receiving terminal each require adjustment. A terminal adjusted with the transmitting section looped to its own receiving section has the same transmission performance as that adjusted with distant terminals. This permits independent line-up of each terminal without requiring the assistance of personnel at the remote terminal location.

The first step in this procedure is to insert the digital 0 dBm0 signal into the input line of the receiving terminal. The gains of the voice-frequency amplifiers are adjusted so that the voice-frequency output for each channel is exactly 0 dBm0. Next, the transmitting terminal is looped* to the local receiving terminal in order to line up the transmitting section. A 0-dBm sinusoidal signal is applied to the voice-frequency inputs of the transmitting section and the gains of the transmitting voice-frequency amplifiers are adjusted until 0 dBm0 is detected at the receiving section voice frequency output.

An important result of this method of line-up is that once the gains are aligned with respect to this well-defined digital signal, any one D2 Bank can be digitally switched to any other D2 Bank, or any other bank so adjusted, without changes in signal levels and without variation in gain. This is a significant step towards the possibility of switching the voice signals in their digital form.

When the terminal is adjusted in this manner, the overload point of the D2 Bank is consistently 3.17 dBm0. This is a fraction of a quantizing step from the objective of 3 dBm0.

5.2 Performance

All of the performance objectives for the D2 Channel Bank have been met in the laboratory model during the development of the D2 Channel Bank. Difficulty was encountered for the production models with regards to the idle channel noise performance. All other performance objectives are met by the production models. The original idle channel noise objective was based on the fact that, according to the smallest step size of the coder near the origin, the theoretical noise floor should

* Since the transmitting terminal output line is looped to the receiving terminal input line, looping can only be accomplished on a 24-channel basis. At initial installation procedures, this presents no problem. However, when the bank is in service and a single channel is to be lined up, it would be necessary to busy out the other trunks on that line prior to looping. In such situations, it may be advantageous to use line-up procedures which may involve personnel at both terminals.

be about 18 dBrnC0. The original objective of 23 dBrnC0 allows 5 dB greater noise to account for imperfections in multiplexing and coding. It turns out, however, that although the mean of the idle circuit noise is fairly close to the theoretical noise floor, the standard deviation is much larger than anticipated.

A histogram of idle channel noise of approximately 100 production D2 Channel Banks, is shown in Fig. 4. The mean is 18.6 dBrnC0 which is very close to the theoretical noise floor. The standard deviation is 1.9 dB. For production testing it was found that a 23 dBrnC0 requirement is not realistic in view of this large standard deviation. The average noise contributed by the D2 Bank remains well below the original objective. Because of its large variation, idle channel noise on a few channels of a particular D2 Bank can be as high as 26 dBrnC0. Since the grade of service⁹ is based on probability of an unacceptable circuit consisting of several channel banks in tandem, isolated above-average noise in one channel of one link will not increase this probability.

Signal-to-quantizing noise performance of the D2 Channel Bank is shown in Fig. 5. This figure illustrates the performance of two typical channels; one with very low idle channel noise, and another with relatively high idle channel noise. It is seen that, for medium-to-high signal levels, the quantizing noise performance is very close to the theoretical ideal performance. For a low-level signal, the idle channel noise of the D2 Bank dominates. The small signal performance of the coder by itself tested in its own test set exhibits almost ideal signal-to-noise performance even for small signals.⁸

Because the companding characteristics of the coder and decoder

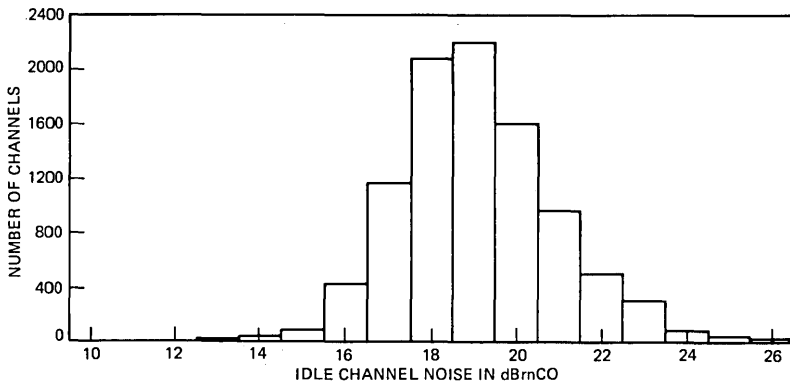


Fig. 4—A histogram of idle channel noise of approximately 100 production D2 Channel Banks.

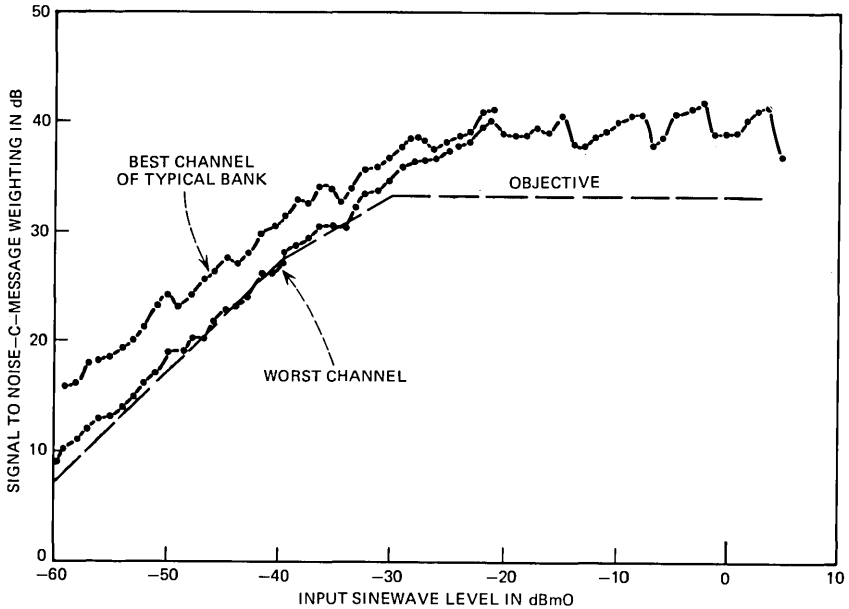


Fig. 5—Signal-to-noise performance of the D2 Channel Bank.

are controlled by highly precise thin-film tantalum resistors, they are expected to be very stable so that field adjustments in any of the common equipment were not necessary. Experience with the gain stability of the voice-frequency amplifiers indicates that analog transmission paths (voice-frequency amplifiers, filters, and gates) have long-term stability to within 0.1 dB after initial alignment.

VI. SUMMARY

This article has summarized the system design considerations of the D2 Channel Bank. The next four articles will cover more detailed circuit aspects of (i) the trunk interfaces which includes channel units and filters, (ii) multiplexing and coding, (iii) digital processing, and (iv) power conversion. The last two articles cover aspects of bringing a product into fruition that are often neglected in the literature. These are physical design, development for manufacture, testing, and continued surveillance and product improvement.

REFERENCES

1. Fultz, K. E., and Penick, D. B., "T1 Carrier System," *B.S.T.J.*, 44, No. 7 (September 1965), pp. 1405-1451.
2. Saal, F. A., "D2: Another Step Toward Nationwide Digital Transmission," *Bell Laboratories Record* (November 1969), pp. 327-329.
3. Davis, J. H., "T2: A 6.3 Mb/s Digital Repeated Line," *IEEE ICC Record*, 1969, pp. 34.9-34.16.
4. Pan, J. W., "Synchronizing and Multiplexing in a Digital Communications Network," *Proceedings of IEEE*, 60, No. 5 (May 1972), pp. 594-601.
5. Kaneko, H., "A Unified Formulation of Segment Companding Laws and Synthesis of Codes and Digital Companders," *B.S.T.J.*, 49, No. 7 (September 1970), pp. 1555-1588.
6. Anderson, E. J., "Considerations in Selection of a Mu-255 Companding Characteristic," *IEEE ICC Record*, 1970, pp. 7-19.
7. Osbourne, P. W., Kaneko, H., and Aaron, M. R., "Synthesis of Code Converters for Segment Companded PCM Codes," *Proceedings of the Conference of Digital Processing of Signals in Communications, Institute of Electronic and Radio Engineers, London, 1972*, pp. 305-308.
8. Dammann, C. L., McDaniel, L. D., and Maddox, C. L., "D2 Channel Bank: Multiplexing and Coding," *B.S.T.J.*, this issue, pp. 1675-1699.
9. Members of the Technical Staff, Bell Telephone Laboratories, Inc., *Transmission Systems for Communications*, Fourth Edition, Winston-Salem, N. C.: Western Electric Company, 1970, pp. 45-48.

D2 Channel Bank:

Per-Channel Equipment

By C. L. MADDOX and D. K. THOVSON

(Manuscript received June 22, 1972)

Traditionally, the design of switching machines has assumed that the transmission plant consists of physical wire pairs. Therefore, appropriate interface equipment is necessary when a multiplex transmission system is substituted for wire pairs. This interface, or per-channel, equipment performs such functions as (i) 2-wire to 4-wire conversion to separate the two directions of transmission, (ii) detection of signaling and supervision information, and (iii) band-limiting and level compensation of the signal prior to sampling. Since failures of the multiplex system can cause massive service interruptions and tie up switching equipment, provision is made to detect such failures and alert the switching machine.

Design of per-channel equipment is complicated by the number of options required to account for the large variety of ways in which switching machines utilize wire-pair plant. This article describes the per-channel equipment of the D2 Channel Bank and some of the considerations that went into its design.

I. INTRODUCTION

In PCM terminals, much of the signal processing takes place in common equipment, where the cost is shared by many or all channels in the system. Although this equipment is more costly than per channel equipment because of higher operating speeds and increased complexity, a net reduction in the per-channel cost is achieved because of time-sharing. However, some functions cannot be performed economically in this way, and per-channel equipment is needed for this purpose.

For the case of the D2 Channel Bank, this equipment includes (i) the lowpass filters and gates, and (ii) the interface equipment between the switching system trunk circuit (or other assigned circuit) and the D2 Channel Bank. The latter circuits are concentrated in

the channel units which contain the signaling scanning and converting circuits in addition to the voice frequency interface circuits.

Looking at a cross section of the telephone plant will reveal a wide variety of switching and transmission equipment in current use. The continued effort to improve the quality of service, the rapid growth of the telephone plant, and the relatively long life of a switching machine have resulted in many switching machines, differing in functions and vintage, existing side by side. In addition, communication between these different machines has also evolved into many schemes. Any new equipment being designed to fit into the existing plant is faced with complex compatibility problems. For the D2 Channel Bank, this meant that many different types of channel units had to be designed and several other types are still in development.

This paper discusses the design of all the per-channel equipment, the manner in which the D2 Channel Bank interfaces with the toll and exchange plant, and how it achieves a measure of compatibility. Included in the general area of interface compatibility are such considerations as VF circuits, signaling and supervision, alarms and service restoration.

II. CHANNEL UNITS

2.1 *General*

Figure 1 shows the major building blocks that constitute one end of a representative 2-wire toll connecting trunk. The location within an office of each of the components will differ, depending on the type carrier system, the means used to signal over it, and the general office layout plan. The D2 Channel Bank was designed as a unitized bay, that is, all equipment shown external to the trunk circuit is included as part of the transmission terminal. To achieve flexibility, all D2 Channel Bank per trunk circuitry is concentrated in a plug-in unit called a channel unit. This permits the channel bank to be adapted to many different trunk types by designing a series of channel units with each channel unit type containing that circuitry unique to a specific trunk.

The major functions of the channel units are (*i*) to provide voice frequency circuit level and impedance conversion and to provide 2-wire to 4-wire conversion for 2-wire trunks, (*ii*) to detect local signaling and/or supervision for transmission to the far end terminal and to reconstruct, from the received digital signal, the far end signaling and/or supervision, (*iii*) to provide jack access to both the voice

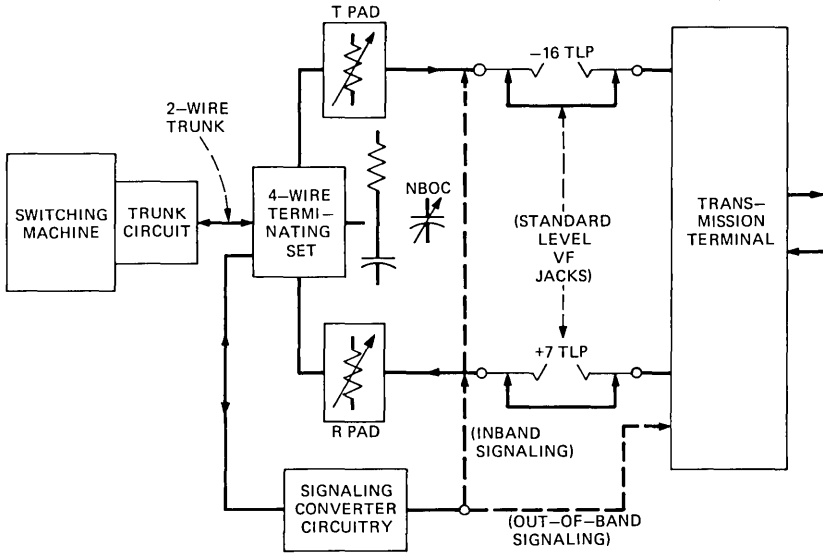


Fig. 1—Typical 2-wire toll connecting trunk.

frequency and the signaling circuitry of the channel bank for alignment and maintenance purposes, and, in some cases, to provide per channel restoration capability. The D2 Channel Bank provides built-in per channel signaling by periodically time sharing the least significant bit of each encoded VF word with a signaling information bit. (A companion article concerning digital functions describes the common signaling circuitry.¹) Thus, referring to (ii), the channel units provide the necessary signaling interface and conversion to the digital format required by the D2 Channel Bank common signaling circuitry.

2.2 Channel Unit Functional Description

Figure 2 shows a simplified block diagram of one of the more complex types of channel units. The channel unit circuitry itself can be functionally separated into two parts; trunk related and channel bank related. The trunk-related circuitry, which is shown to the left of the jacks in the diagram, serves to interface with the specific voice frequency and signaling circuits of the trunk, and to furnish converted signals to the channel bank related circuitry. These circuits account for the primary differences between channel unit types. The channel bank related circuitry shown to the right of the jacks in the diagram has a

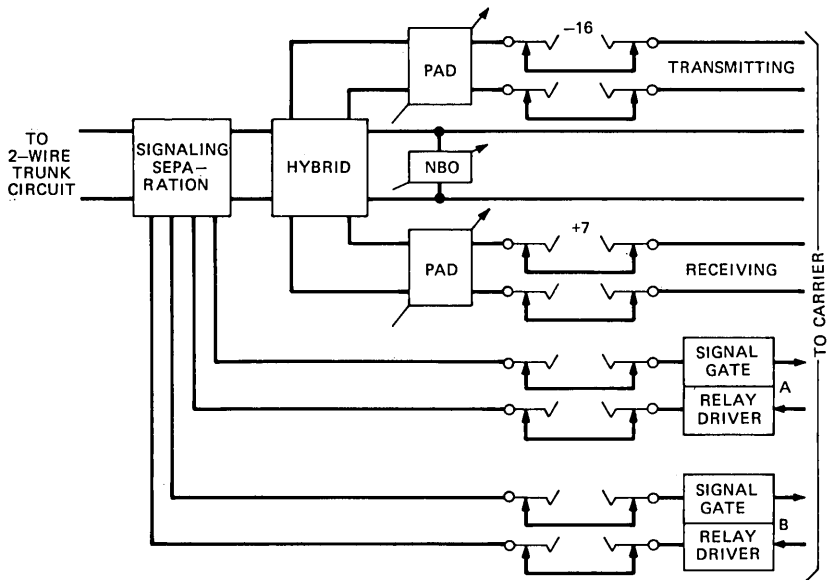


Fig. 2—Simplified diagram of a more complex type of channel unit.

standard output regardless of the channel unit type. Thus there are no restrictions on the physical placement of any channel unit type within the available 96 channel unit positions in the terminal.

For the case of the toll trunk and special service channel units, there are jacks at the interface between the trunk-related and bank-related circuits. These jacks are mounted on the faceplate of the channel unit. At these points, signals (both voice frequency and signaling) appear at standard levels. This permits the restoration of service on an individual channel basis. For example, in case of a failure of channel bank digroup equipment, a trunk may be processed in the trunk-related circuitry in the channel unit associated with the failed channel. The output of these circuits are then patched to the bank-related circuitry of another working channel unit of similar type, either in the same channel bank or in another D2 Channel Bank.

Channel units designed primarily for use with exchange area trunks are not required to provide per channel restoration capability, hence they do not have jack access to standard level signaling circuitry. However, the 4-wire voice-frequency patch jacks with the -16 dBm transmit, and $+7$ dBm receive levels appear on all channel units for maintenance purposes.

2.3 *Voice Frequency Circuits*

Both 2-wire and 4-wire trunks are in common use in the telephone plant. The D2 Channel Bank requires a 4-wire 600-ohm balanced standard level VF circuit at its transmit and receive ports. Therefore, all 2-wire channel units have a built-in 4-wire terminating circuit to provide the necessary 2-wire to 4-wire conversion. The 4-wire terminating circuit consists of a 2-transformer hybrid, transmit and receive loss adjustment pads, a balance network, and a selection of network build-out capacitors. The 4-wire terminating circuit meets all toll transmission requirements.

The 4-wire channel units do not require the hybrid and its associated balance network, but they do provide the transmit and receive loss adjustment pads. These pads, in both the 2-wire and 4-wire channel units, are tantalum thin-film resistor networks controlled by slide switches. The attenuation range is 16.5 dB adjustable in 0.1 dB steps.

As previously mentioned, all channel units provide jack access to the 4-wire standard level VF circuit.

2.4 *Signaling*

2.4.1 *General Considerations*

The term signaling, as used in this paper, is assumed to include both address and supervisory information. Since most switching machines are designed to work with physical wire pairs, signaling systems used by switching machines are based on wire plant. Different switching machines however use the wire plant for signaling in a variety of ways as determined by their trunk circuit. Thus it becomes necessary for carrier systems including the D2 Bank to provide signaling interfaces that imitate the wire plant in a variety of ways. Effort is being made by the telephone industry to separate the signaling information path from the voice frequency path, especially when signaling information is to be transferred from one switching office to another. Economy and efficiency of operation will result with the introduction of the proposed common channel interoffice signaling (CCIS). While the D2 Bank has provisions for carrying a CCIS channel, until CCIS is implemented, signaling interfaces must be designed that look like wire plant to the trunk circuit of the switching machine.

2.4.2 *Signaling System Considerations in the Channel Unit Design*

The signaling systems encountered in the Bell System plant that were considered in the channel unit design are:

- (i) Dial pulse
- (ii) Revertive pulse
- (iii) E and M lead
- (iv) Special access.

Signaling systems (i) and (ii) are used primarily in exchange area trunks. Both are considered 1-way trunks in that a call can only originate from the outgoing end. Different circuitry is required in the channel units that face the outgoing and incoming ends of these trunks. Therefore, different originating and terminating channel units were designed for both systems. Dial pulse and revertive pulse both require 2-state signaling in the forward direction. In the reverse direction, dial pulse requires 2-state and revertive pulse requires 3-state signaling. Figure 3 is a diagram of a representative channel unit for these signaling systems.

Revertive pulse signaling, being a feedback system, is sensitive to round-trip delay. The revertive pulse originate channel unit reduces the round trip time somewhat by taking advantage of the fact that only 2-state signaling is required in the forward direction, and transmits this information in both signaling channels A and B. Because the signaling channels in the D2 Bank are derived serially, transmitting the same signaling information in both channels reduces the time to detect a signaling transition by a factor of two.

E and M lead signaling, (iii), is found primarily in toll applications. It is a symmetrical signaling system that can be used for both one- or two-way trunks. The same type channel unit serves both ends. The signaling takes place on the E and M leads which are separate from the VF circuit. Two-state signaling is required. Figure 4 shows a diagram of a 4-wire E and M channel unit.

Special access, (iv), is a type of service characterized by a trunk that signals toward the station by ringing and toward the office by dial pulsing. It is a nonsymmetrical 2-way signaling system. Ground start PBX trunks are an example of this type of signaling system. The full 4-state signaling capacity of the D2 Bank is used in both directions. Different channel units are required at each end of the trunk. Figure 5 is an example of a channel unit designed for this service.

2.4.3 *Detecting the Local Signaling States*

The standard signal conditions that must be detected and reconstructed by the channel unit fall into three categories:

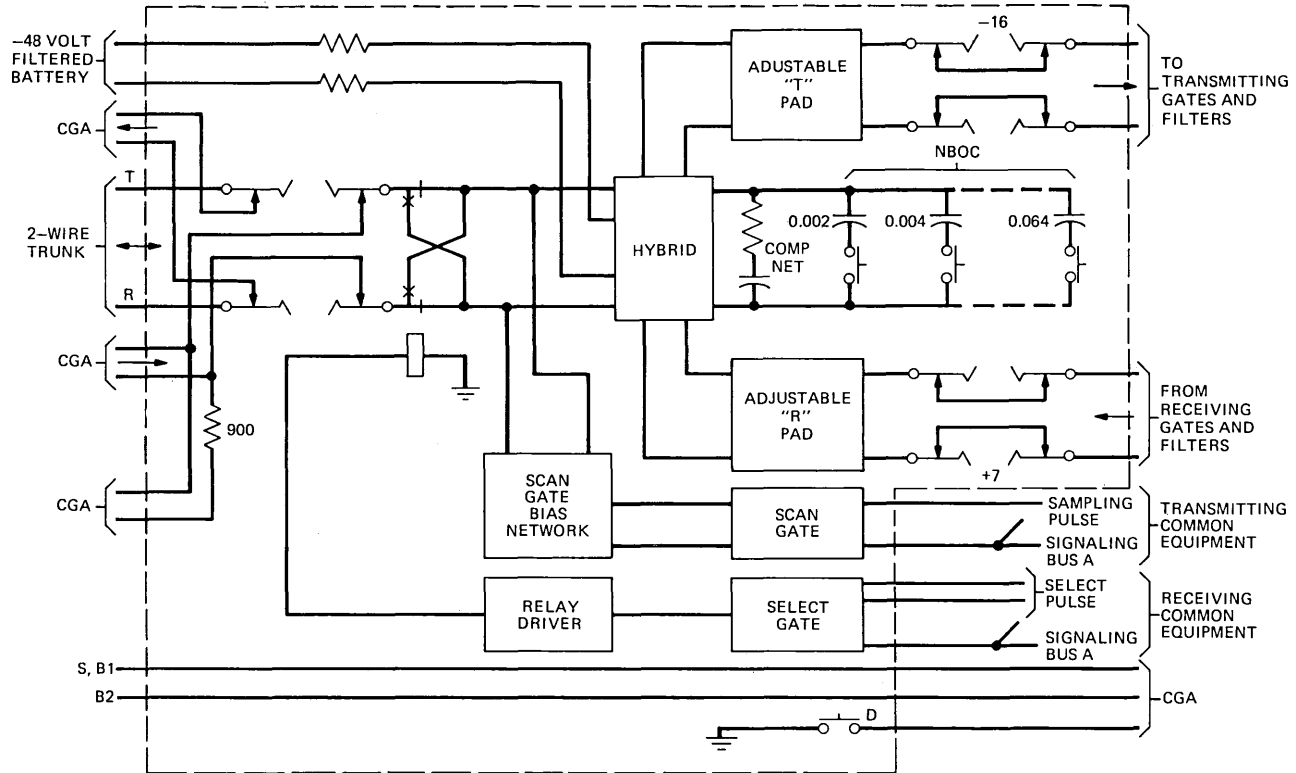


Fig. 3—Dial Pulse Originating channel unit.

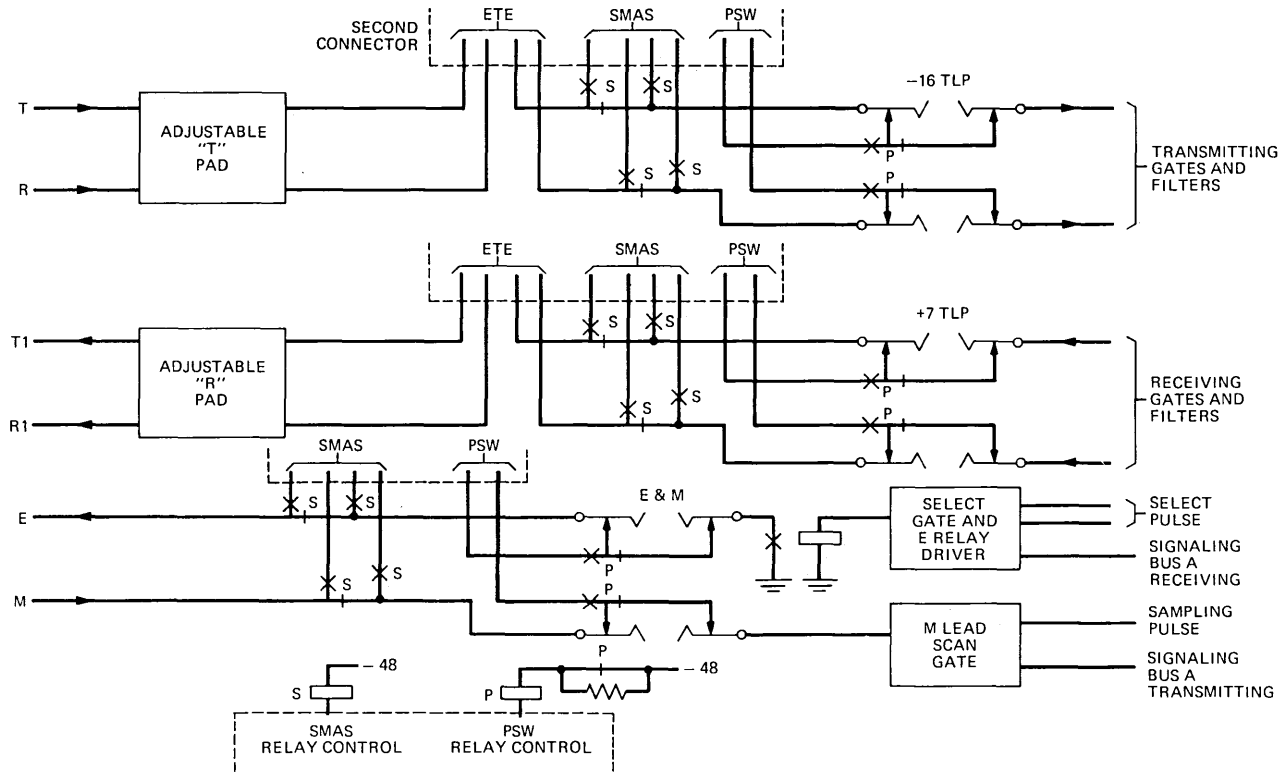


Fig. 4—E & M 4-W + ETE + SMAS + PSW.

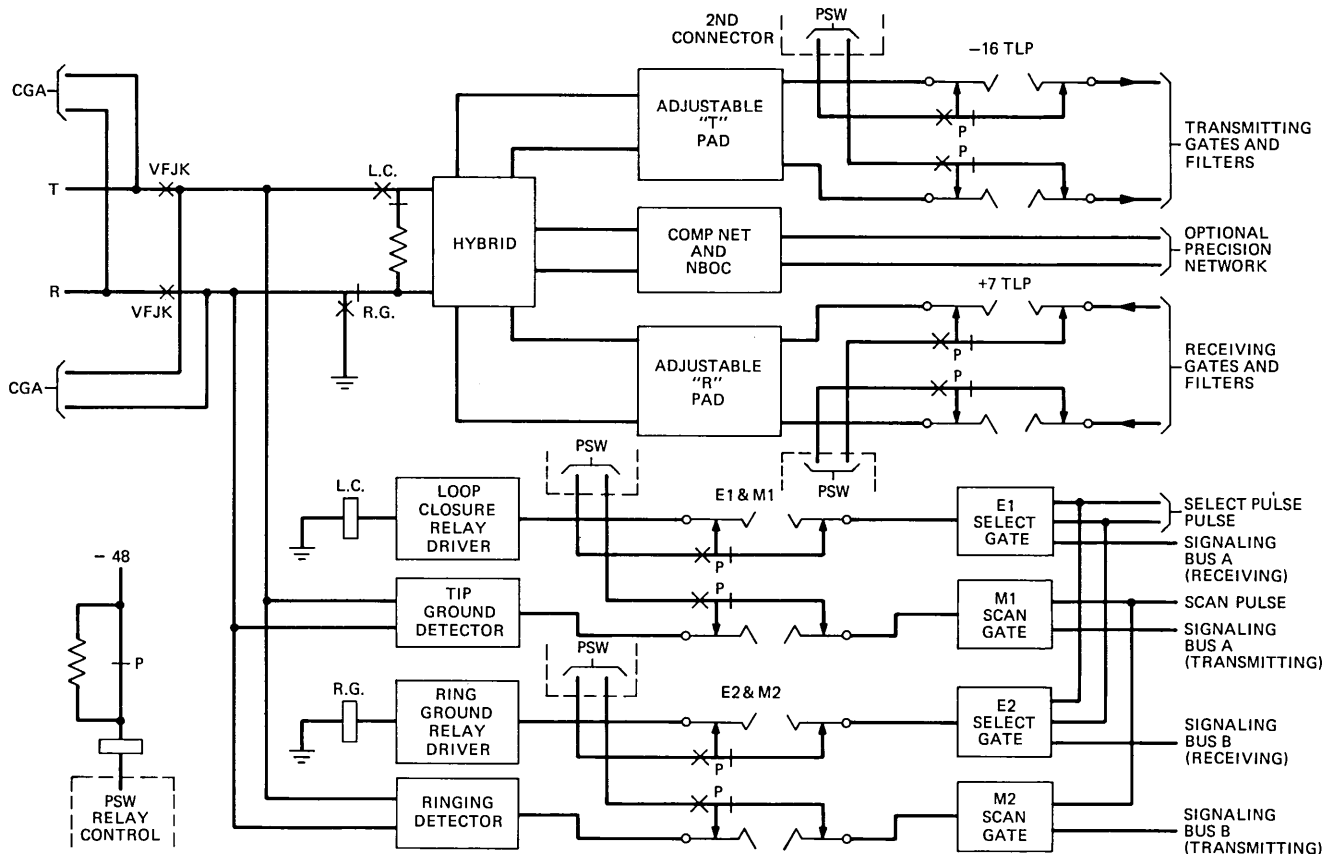


Fig. 5—Special Access-office end + PSW.

- (i) Impedance, high or low
- (ii) Voltage polarity
- (iii) Ringing.

These conditions must be detected either on a 2-wire loop which also carries the voice frequency signal, or on separate leads devoted exclusively for signaling. For the detection of impedance conditions, a conversion to dc voltage conditions is made by supplying dc current to the loop via a resistive network at the center of the hybrid. In many cases, this current also serves as talking battery. The various scanning bias networks shown in the example channel units in Figs. 3 to 5 are high-impedance resistive networks bridged on to the circuit. They convert the dc signaling (and supervision) voltage to standard levels suitable for sampling by the scanning gates. The scanning gate bias circuitry also includes filters to suppress the 8 kHz sampling pulses which would feed back to the voice frequency circuit. The scanning gate circuitry is similar for all channel units, and consists of a sampling diode, a multiplex diode, and two dc isolation capacitors arranged such that the bias network will either forward or reverse bias the sampling diode, depending on the signaling state of the 2-wire circuit. The transmitting channel counter provides a pulse in each channel time slot which interrogates the sampling diode in the channel unit. This channel pulse is either passed or blocked by the sampling diode and its presence or absence is multiplexed on one of the transmitting signaling busses.

Each scanning gate handles two signaling states. When 4-state signaling is required, as shown in the case of Fig. 5, this circuitry is duplicated. For the case of revertive pulse signal, where signaling round trip delay is a critical requirement, one scanning gate drives both signaling channels in the revertive pulse originate channel unit, thus doubling the speed at which signaling information is transmitted in the forward direction.

2.4.4 *Reconstructing the Far End Signaling States*

Signaling information received from the far end is furnished to the channel unit from receiving common signaling on receiving signaling bus A (and bus B when required). A balanced channel pulse from the receiving channel counter interrogates the receiving signaling bus in the select gate. On the basis of the information on the receiving signaling bus, the select gate provides a positive or negative pulse to a flip-flop contained in the relay driver. The flip-flop bridges the time between incoming signaling information samples, and operates or releases a

relay the contacts of which restore the far-end signaling condition to the near-end trunk. In the example shown in Fig. 3, the far-end supervisory state is indicated by the polarity of the battery supplied to the near-end switching equipment. As shown in Fig. 5, two additional signaling states can be detected by adding another select gate operating from receiving signaling bus B analagous to the transmitting end.

2.5 External Trunking Equipment Option

Some of the E and M lead signaling-type channel units are designed with additional access to the 4-wire, voice-frequency circuits. This access permits external trunking equipment (ETE) such as echo suppressors, delay equalizers, voice-frequency filters, etc. to be connected with either transmit or receive voice frequency circuits. Figure 4 shows how this access is implemented. The mass of additional wiring required for ETE necessitates an additional 40-pin connector on those channel units with ETE provision.

2.6 Channel Unit Design Summary

The difficulty in the design of these circuits has been to accommodate large variation of loop length and tolerances of impedance and voltage encountered in the trunk circuit. Also, existing wire trunks were designed with components and voltages which produce circuit conditions hostile to solid state devices. Consequently, relays were used in some cases as the most practical device to reliably perform the required function. In spite of these problems, the per-channel signaling circuitry is extremely simple. Economy in providing signaling contributes considerably to the low cost characteristic of PCM channel banks.

III. FILTERS AND SAMPLING GATES

Lowpass filters are employed for bandlimiting prior to sampling in the transmitting terminal as well as for reconstructing the signal after demultiplexing in the receiving terminal. The transmitting and receiving filters and gates are shown in Figs. 6 and 7. The multiplexing gates will be discussed in more detail in the next article.²

Lumped element LC filters were chosen rather than active RC filters because at the time of introduction of the D2 Channel Bank, the large quantities of active filters required could not be produced at a cost competitive with lumped element LC filters. However, the impedance and signal voltage levels for the demultiplexing filters were chosen so that it will be possible at a later time to convert production of the

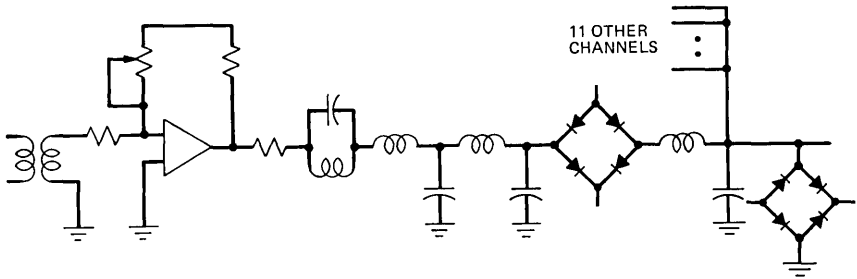


Fig. 6—Transmitting gates and filters.

demultiplex filter to active RC filters if it becomes economically attractive to do so.

The transmission characteristics of the transmitting and receiving terminals are shown in Fig. 8. The voice-frequency transmission characteristics of the terminal are primarily determined by these filters. However, the low frequency cut-offs are due to coupling capacitors and the voice frequency transformers at the input and output terminals.

The transmitting voice-frequency amplifier is a two-transistor circuit. The purpose of this amplifier is to provide a good return loss, and to provide a means of adjusting the level of the signal in the transmitting terminal. The transmitting gain control permits the lineup of signal levels in relation to the digital signal, and so permits the possibility of digital switching.

Two-transistor voice-frequency amplifiers are used in the receiving terminal to raise the output level to +7 dB at the impedance level of

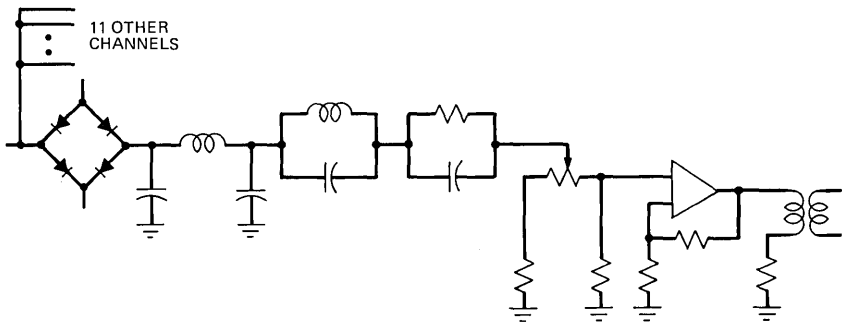


Fig. 7—Receiving gates and filters.

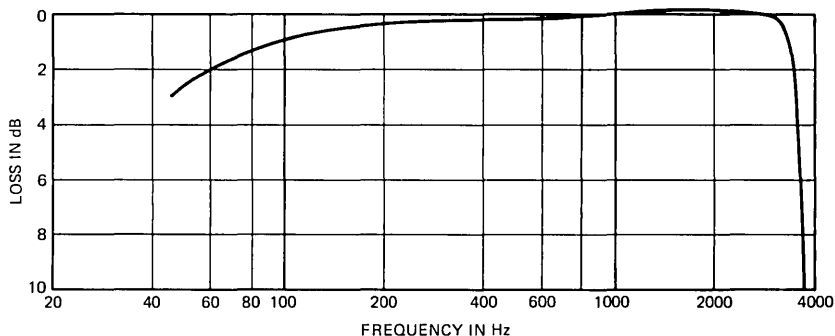


Fig. 8—D2 transmission.

600 ohms. A gain control is included to permit per-channel gain adjustments at the receiving terminal.

IV. MAINTENANCE AND ALARMS

4.1 *Maintenance Access*

In the design of the channel unit, great emphasis was placed on having convenient access to each channel for maintenance purposes. This resulted in the availability of standard 4-wire voice-frequency level and impedance points on all channel units. In the exchange-type channel units using loop signaling, 2-wire jacks were provided for use in conjunction with signaling tests.

Some of the E and M channel units were designed with provisions for switched maintenance access (SMAS). This system permits an operator at a centralized location to detect and condition a particular trunk through the switching machine, and to monitor remotely the 4-wire voice frequency and the signaling leads. For this purpose, a special interface circuit for the D2 Bank was designed. This circuit energizes the SMAS relay (Fig. 4) in the channel unit and routes the voice frequency and signaling points via a switching network to the maintenance center.

4.2 *Alarm Control*

The D2 Channel Bank is equipped with alarm control circuits. These circuits monitor power supply voltages, fuses, and coder zero set voltages. They also monitor the receiving terminal and its ability to maintain synchronization with the incoming bit stream. When a failure occurs, the alarm circuit initiates the appropriate action. This includes

audible and visual alarms, operation of the carrier group alarm, and transmission of the alarm indication to the remote terminal. No action is taken by the alarm control circuit unless the trouble condition persists for a period greater than 300 milliseconds. This delay time is built into the system in order to prevent a complete alarm cycle to be initiated by short hits on the repeatered line, which under normal conditions will recover well within this period. The transmission of the alarm condition to the remote terminal is accomplished by forcing the second digit in all words to zero. Under normal conditions, a long string of zeros is extremely unlikely.

4.3 *Carrier Group Alarm*

The carrier group alarm is a trunk processing unit the function of which is to disconnect existing calls and make all trunks busy when an alarm condition occurs. Activation of the CGA starts a timing sequence which provides each trunk with on-hook supervision to disconnect active channels and to stop charges and then provides the appropriate signal to make the trunk appear busy to the switching machine. The timing sequence lasts for a minimum of 20 seconds upon the initial receipt of an alarm signal, or for 10 seconds after the alarm signal disappears. This timing sequence minimizes the susceptibility of the CGA to intermittent operation.

As is the case with the channel unit circuits, the CGA has to be matched to the processing requirements of the particular trunk circuit to which it is connected. This selection is done in the channel unit. Jumper wires in the channel unit connect together those leads which provide the proper sequence and final state necessary to remove that trunk from service. However, most channel units are used with different trunks or switching machines even though they use the same type of signaling and supervision method. Therefore, some channel units have CGA options which are selected at the time of initial installation.

V. RESTORATION

In addition to manual restoration patching, provision for manually-initiated protection switching has been made.

High-priority circuits, remotely monitored unattended offices, or other special services sometimes require rapid automatic means of restoring service on a failed facility. A special standby D2 Channel Bank has been designed which can be used to protect up to ten fully equipped D2 Channel Banks. When it is determined that a failure has occurred

in one or more digroups of a protected D2 Channel Bank, all of the circuits being served by that digroup can be switched en masse to the standby bank. The digital line output of the standby bank is simultaneously switched to the digital transmission facility originally serving the failed channel bank.

As shown in Fig. 4, the per channel switching occurs at the standard level 4-wire VF point, and at a defined E and M like signaling point in the signaling circuitry. Provisions for 4-state signaling have been included. The actual per-channel switch is performed by a relay included in those channel units designed for protection switching and designated by the suffix PSW. As with external trunking equipment, the PSW channel units also make use of the second connector on the channel unit.

VI. SUMMARY

This article has described the per-channel equipment of the D2 Channel Bank. This includes (i) equipment that is used as the interface between the switching machine and the carrier transmission equipment, (ii) the filters and gates that condition the signal for sampling and multiplexing, and (iii) provision for maintenance and restoration of service in case of failure. This equipment constitutes the bulk of the physical space occupied by the D2 Bank and shares about half of the total cost. The design becomes complicated due to the large variety of trunk circuits and options encountered in the field.

REFERENCES

1. Cirillo, A. J., and Thovson, D. K., "D2 Channel Bank: Digital Functions," B.S.T.J., this issue, pp. 1701-1712.
2. Dammann, C. L., McDaniel, L. D., and Maddox, C. L., "D2 Channel Bank: Multiplexing and Coding," B.S.T.J., this issue, pp. 1675-1699.

D2 Channel Bank:

Multiplexing and Coding

By C. L. DAMMANN, L. D. McDANIEL, and C. L. MADDUX

(Manuscript received June 22, 1972)

Analog multiplexing and coding in the D2 Channel Bank is discussed in this article. Multiplexing of the message signals is accomplished in two stages. In the first stage, groups of 12 channels are multiplexed together using resonant transfer gates. The resulting eight buses, each carrying pulse-amplitude-modulated signals of 12 channels, are then multiplexed in the second stage. The samples of all 96 channels are presented to a single coder. The demultiplexing plan follows the inverse of the multiplexing plan. The output of the decoder is first divided into eight buses, and the final demultiplexing is accomplished in groups of 12 channels. Because the decoding is accomplished by an asynchronous time-shared decoder, storing and stretching of the analog samples is necessary to permit removal of the time jitter due to the queuing process.

The coder used in D2 is a nonlinear coder using a compression characteristic called the 15-segment approximation to the $\mu = 255$ law. To ensure the success of the coder development, a stage-by-stage binary coding plan was chosen. The first stage determines the polarity of the signal, and the succeeding binary stages determine the amplitude of the compressed signal one digit at a time. To achieve accuracy in the coder with available devices, automatic zero-setting circuits are used in a feedback loop to control offset deviations. This is in addition to the use of precision resistors and precision power supplies for the remaining critical parts of the coder. In order to achieve comparable accuracy in the decoder, the same stage-by-stage arrangement is also used. Again, automatic zero-set feedback loops are used to control drifts. The performance of the coder/decoder combination has met the objectives.

I. INTRODUCTION

This paper is concerned with the analog multiplexing and the coding processes in the D2 Channel Bank. Economies in a digital channel bank

are mostly due to the ability to share complex equipment such as the coder among all of the channels. To permit such sharing to take place, it is necessary to time-division multiplex the message signals from all the channels. This is accomplished by the multiplexing circuits. Because of the fragile nature of pulse-amplitude-modulated signals, it is here that most of the degradations, other than quantizing noise, are introduced. Great care must be exercised in the electrical and mechanical design of the multiplexing and demultiplexing circuits.

The bulk of this paper is devoted to coding. The presence of a coder is characteristic of digital channel banks, since it converts analog signals into digital form. The choice of the compression law will be discussed here. The method used to achieve the chosen compression law will be described. The speed and accuracy achieved by the D2 coder represented a significant technical advance in the art of analog-to-digital conversion. The success in developing and manufacturing this coder to the required speed and accuracy proved the basic soundness of the approach taken.

II. MULTIPLEXING

It was recognized early in the system design of the D2 Channel Bank that time-division multiplexing of the analog signals of all 96 channels in one step would present a difficult electrical and physical design problem. This is because it would be very difficult to control crosstalk and signal interference with a fan-in of 96 to 1. With two stages of multiplexing, the fan-in is reduced for the first stage. Greater care can then be exercised in the circuits used for the second stage of multiplexing since these will be fewer in number and their cost will be shared over many channels. Although the number 96 can be factored into many different combinations, it was decided to use the 12-8 multiplexing plan because the Bell System has traditionally used 12 channels to form a basic group in the frequency division multiplexing plan. (See Fig. 1.) Two such groups consisting of 24 channels are called a digroup meaning two groups. Digroup has also been construed to mean digital groups. The second meaning gained popularity because digroups appear predominately in digital channel banks.

2.1 *The First Stage of Multiplexing*

For the first stage of multiplexing the channels must be sampled as well as multiplexed. The sampling rate is 8 kHz. The basic cycle for sampling one of the 12 channels is 10.4 microseconds. The resonant

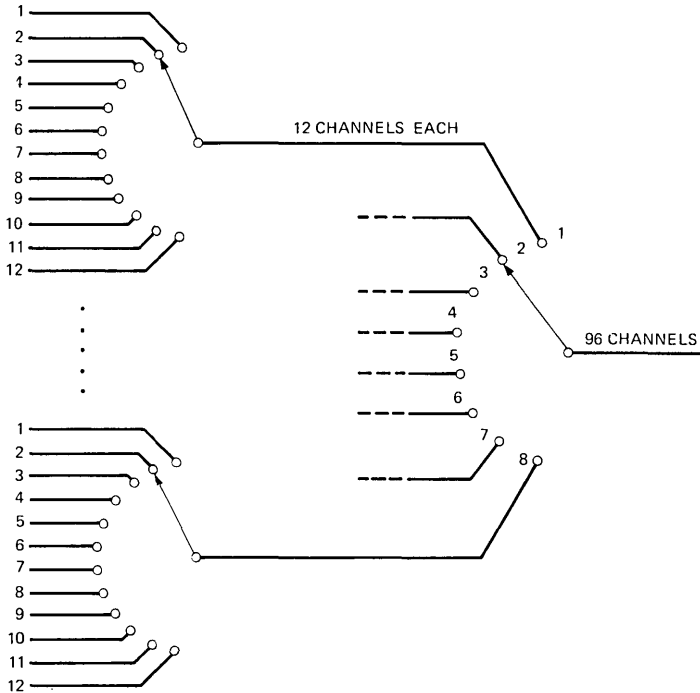


Fig. 1—Multiplexing plan.

transfer technique is used to perform the sampling. This transfer is accomplished in 2.2 microseconds (Fig. 2a). The result of resonant transfer is stored in the hold capacitor of the multiplexing bus for 5.2 microseconds (Fig. 2b). During this time the held sample is available for the second stage of multiplexing. After the hold period, clamping takes place for the rest of the 10.4-microsecond time to prepare the hold capacitor for the next resonant transfer. It is seen that on each multiplexed bus of 12 channels the samples are available only half the time. It is thus necessary to stagger the operation of sampling in pairs so that when a sample is available on one bus ready for the next stage of multiplexing, clamping followed by resonant transfer takes place on the other bus.

In a fully-equipped channel bank, there are eight buses, each with a holding capacitor, one for each group of 12 channels. Resonant transfer and the clamping occur in four of the eight buses simultaneously. This permits the use of a single channel counter to drive four resonant

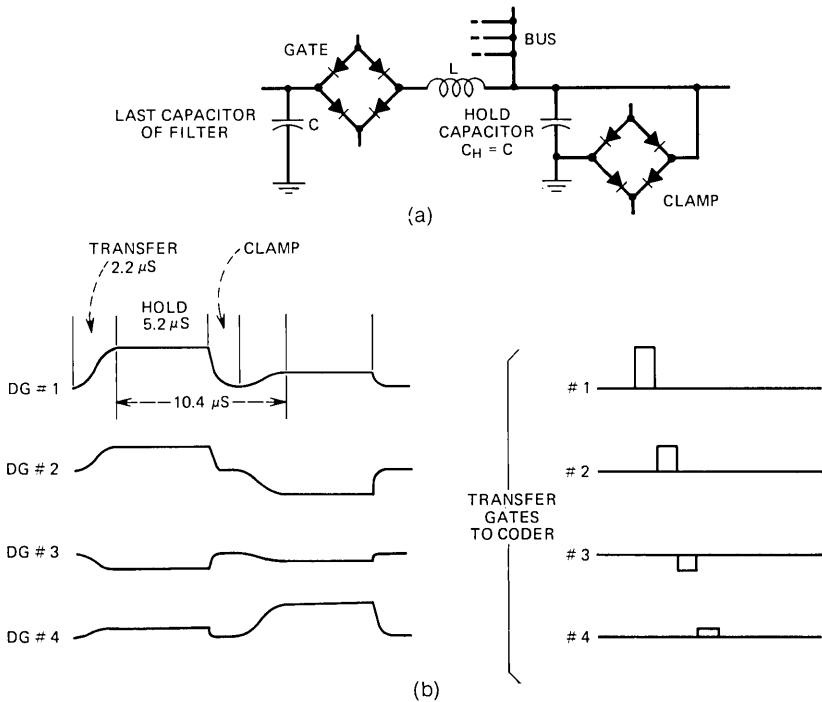


Fig. 2a—Resonant transfer.
 Fig. 2b—First stage sampling and multiplexing.

transfer gates. A shift register of 24 stages is thus shared by 96 channels. Resonant transfer is used for this first stage of multiplexing because the resultant higher signal levels are more immune to noise and interference. This is particularly important because the multiplexing bus is physically spread over a large area on the D2 bay.

2.2 Second Stage of Multiplexing

The second stage of multiplexing combines the sample values from each of the eight holding capacitors, and feeds them to the coder for conversion into digital form. (See Fig. 3.) This multiplexing is accomplished by the use of balanced diode gates that follow the holding amplifiers associated with each of the eight holding capacitors. An operational amplifier is used for each pair of the diode gates. Thus four operational amplifiers are required. The outputs of these operational amplifiers are summed at the input of the coder. This arrangement permits equipping a channel bank with one digroup at a time and has

the further advantage that in the event of a failure of a single amplifier or gate, all channels do not have to be disabled to replace the defective one. The holding time of the first stage of multiplexing is divided into four 1.3-microsecond intervals and four samples are read sequentially into the coder during these intervals of coding. As will be explained in Section V, in order to provide a full 1.3-microsecond interval for zero setting the coder, the transfer gate operation has a timing offset every other frame. To permit this timing offset, the holding time at the first stage of multiplexing is made slightly longer than 5.2 microseconds. This delays the clamping operation by 0.65 microsecond.

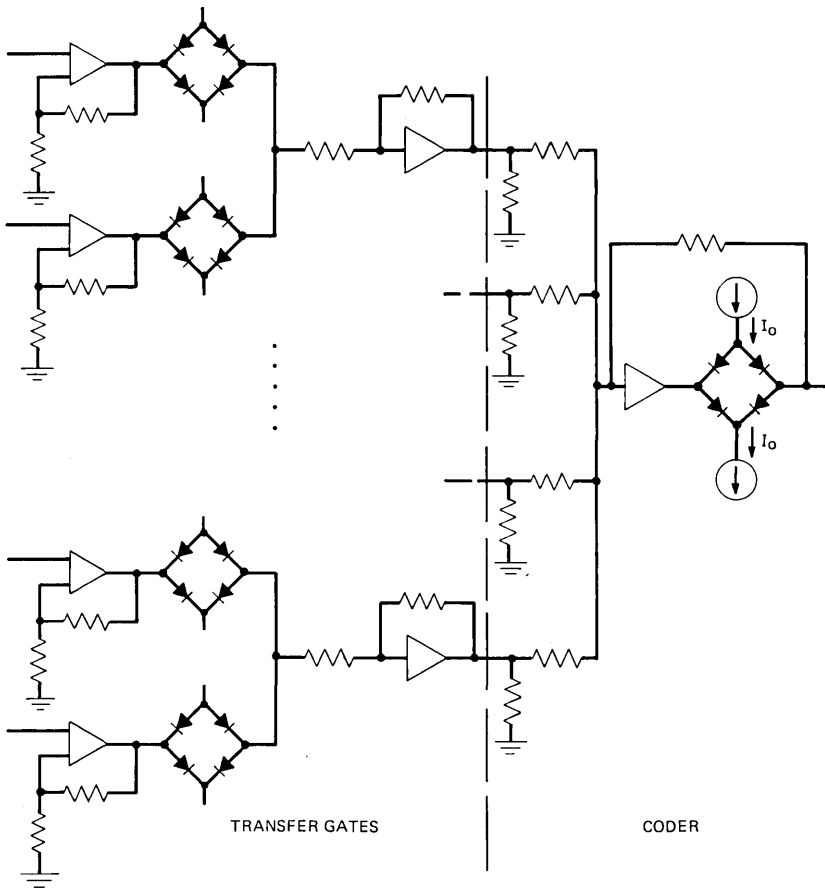


Fig. 3—Transfer gates.

III. DEMULTIPLEXING

On the receiving side after the incoming PCM code words are decoded by a common decoder, the PAM samples are demultiplexed in a two-stage operation following the inverse pattern of the multiplexing operation. First, the PAM samples are demultiplexed into eight groups of 12 channels each. Second, the groups of 12 channels are broken down into individual voice-frequency channels.

3.1 *Select and Hold*

Because the single decoder in a D2 Channel Bank is shared by four asynchronous incoming digroups, queuing takes place at the digital input to the decoder.¹ Consequently, the PAM samples, as delivered by the decoder, are jittered in time with respect to the derived incoming clock. This jitter can be as large as 5 microseconds. Removal of this jitter is accomplished in the select-and-hold circuit which also acts as the first stage of demultiplexing. The details of queuing are discussed in the next article.

At the same time that a digital word is transmitted to the decoder, the appropriate selection gate is turned on to steer the decoded PAM sample to one of the eight holding capacitors (Fig. 4). The samples are then ready for the channel pulses to steer them for the second demultiplexing stage. The time constant of this capacitor in parallel with the input impedance of the holding amplifier is made long so that the result of variable holding time will not add any detectable noise to the signal.

In each of the holding capacitors there is a sequence of PAM samples representing signals from the group of 12 voice frequency channels. Crosstalk can be caused by residual charge from the sample from a previous channel. Clamping the capacitor is one way to reduce this crosstalk. It is obvious that the same crosstalk performance can be achieved by precisely charging the hold capacitor to the new sample value. This is accomplished by making the time constant of the hold capacitor in combination with the output impedance of the driving amplifier very small. To protect the amplifier from the high currents that would result for such short-time constants, a current limiting diode bridge is interposed between the amplifier and the selection gates.

3.2 *Final Demultiplexing*

Each select-and-hold circuit is connected to the 12 voice frequency receiving filters through balanced diode transmission gates. These

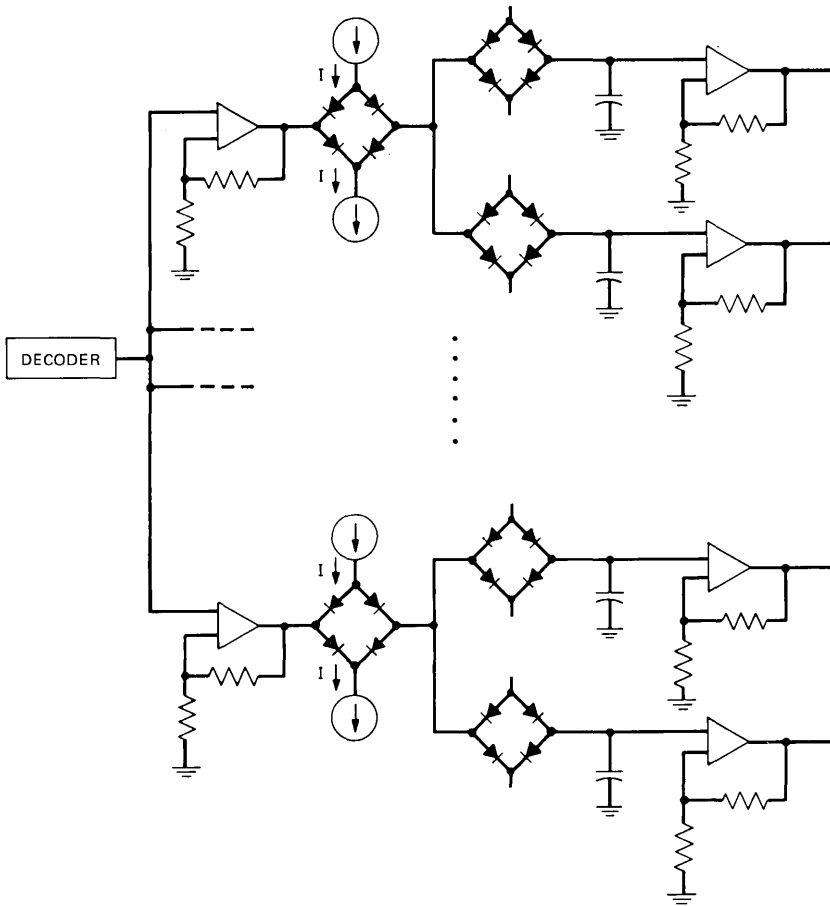


Fig. 4—Select and hold.

gates are under the control of the receiving channel counters and they demultiplex the 12 sequential PAM samples held by the select-and-hold. Voice frequency filters then reconstruct the signal from the PAM sample. Since the four incoming digroups are expected to be asynchronous, there are four sets of receiving channel counters of 24 stages each—one set for each digroup. The receiving channel pulses generated by the channel counters are derived from the incoming line. Thus they are not influenced by the asynchronous queuing logic of the decoder.

IV. CODING PLAN

For PCM coding of message signals, it is well known that a non-uniform assignment of code words to various amplitudes is necessary. This is because message signals can be expected to have a dynamic range of over 40 dB.² In the D1 Bank, this nonuniform coding step size is achieved by compressing the signal with a nonlinear circuit before coding and expanding the signal after decoding. The nonlinear circuit that does the compression and expansion is called the compandor and the transfer characteristic of the compressor is called the compression (or companding) characteristic. For uniform signal-to-distortion ratio over a wide dynamic range, a logarithmic compression characteristic is required. A pure logarithmic function cannot be used since the function approaches minus infinity at zero. Various approximations to the logarithmic characteristic have been proposed. The one proposed by the Bell System is called the μ -law,² and the one proposed by the British is called the A-law.^{3,4} Both of these laws become linear at the origin. The D1 Channel Bank uses the nonlinear properties of diodes to approximate the μ -law characteristic. The resultant compression characteristic falls somewhere between the μ - and A-laws. During the initial planning of the D2 Channel Bank a new family of compression laws was proposed. These are called the digitally linearizable compression laws. These laws are piecewise linear approximations to the logarithmic laws. Furthermore, coders using these transfer characteristics have step sizes that are related to each other in powers of

TABLE I—CODER SEGMENTS

Coder Input	Segment	Step Size
8159	000	256
4063	001	128
2015	010	64
991	011	32
479	100	16
223	101	8
95	110	4
31	111	2
0		

two. It is thus possible to take the resultant binary code words representing the compressed signal amplitudes, to translate them easily into binary code words representing the linear or uncompressed signal, and to do so without incurring any additional quantizing noise degradations and with the linear step sizes no smaller than the smallest step sizes of the compressed code words.

There was some disagreement as to the exact detail of these digitally linearizable laws. The Bell System proposed a 16-segment law, symmetrical about the origin, with the center two segments having the same step size, and each succeeding outer pair of segments having step size double that of the previous pair of segments. This is called the 15-segment approximation of the $\mu = 255$ compression law. The definition of the 15-segment compression law is shown in Table I and also illustrated in Fig. 5. At the same time CCITT in Europe favored a similar 16-segment approximation where the inner *four* segments are made to have the same step size and the next succeeding pairs of outer segments doubling in step size. This is called the 13-segment approximation to the A-law. In comparing the two compression laws, one can find that the smallest step size used by the 15-segment approximation is about one half the step size used by the 13-segment approximation. This tends to give the 15-segment law better idle channel noise and crosstalk performance. However, for the same number of digits, the 15-segment law necessarily has slightly larger step sizes in the outer segments. It has been argued that the potential performance

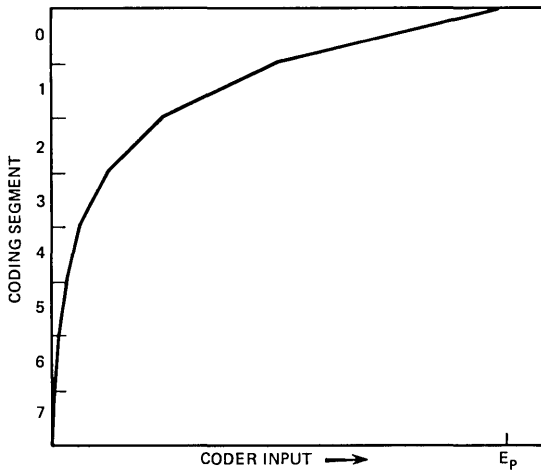


Fig. 5—15-segment $\mu = 255$ compression characteristic.

capabilities of the 15-segment law cannot be achieved with present technology. This is largely true as can be seen by the performance of the production D2 Channel Banks. The 15-segment law was chosen for the D2 coder because (i) the slightly lower signal-to-distortion performance at mid-to-upper signal levels as compared to the 13-segment law is not significant, and (ii) it allows the potential for future improved performance so that the 15-segment law can be made to be the standard for future digital channel banks as well.⁵

4.2 *Output Code Format*

The choice of the output code format is important because of the desire to make this the standard for future channel banks also. The D1 Channel Bank used the ordinary binary code where negative signals are expressed as the complement of the corresponding positive signals. An alternative to the ordinary binary code is the signed binary code or folded binary code where the first digit indicates the polarity of the signal, and the following digits indicate the magnitude. It has been shown that the folded binary code is superior to the ordinary binary code in masking transmission errors when a speech signal is carried.⁶ This is because, with the ordinary binary code, an error in the first or most significant digit will always cause an amplitude error of half range, whereas with the folded binary code, a transmission error in the first or sign digit causes an error which is proportional to the signal. Since speech signals have the highest probability at zero, the average error voltage caused by a transmission error would also tend to be small.

A second choice concerning the output code format is whether magnitude should be transmitted in straight binary or its complements. For the folded binary code, the two choices result in different densities of pulses on a transmission line. When the magnitude is transmitted straight, small signals or no signal results in a very low density of ones being transmitted on the line. This could cause timing problems in the repeatered line. When the inverted binary code is transmitted, small or no signal would cause a very high density of ones on the transmission line. This would result in a strong timing signal in the repeater line. However, in a cable containing many digital transmission systems, a dense pattern of pulses would result in greater amounts of crosstalk from one system to another. For the ordinary binary code, there is no reason to make a choice since negative values have code words which are the exact complement of positive values.

The inverted folded binary code was chosen as the output code

format because it is thought that maintaining good timing performance in the repeatered line is important. The reduced timing error in the repeatered line would then leave more margin for amplitude degradations due to crosstalk. Compromise alternatives have been considered. For example, it is possible to invert every other bit at the output. This was proposed for the 13-segment approximation to the A-law. Compared with the inverted folded binary code it has both advantages and disadvantages, and the net difference was not considered significant.

V. THE CODER DESCRIPTION

Many methods are available that can be used to code an analog sample into a compressed binary code according to the 15-segment approximation to the $\mu = 255$ compression law. The nonlinear characteristics of diodes used in the D1 Bank are not suitable because they do not yield a piecewise-linear approximation. The digitally linearizable property permits the sample to be initially coded into a linear code which can then be followed by digital processing to achieve the compressed code. A desirable property designed into either the 15-segment or 13-segment code is that it permits a coder design using either the feedback arrangement or the stage-by-stage coding arrangement. For the feedback arrangement a local decoder is used, and the digits are determined sequentially by a single comparator. The local decoder consists of a linear binary decoder which determines the steps within each linear segment, followed by a ladder attenuator which determines the slope of the compression law.* The stage-by-stage coding method was chosen for D2 primarily because transmission gates and control logic in integrated-circuit form were not available at the time of development. Stage-by-stage coding does not use transmission gates. It is more complex, but within a D2 Bank the coder constitutes a very small fraction of the total cost.

5.1 Coding Method

The nonlinear coder consists of an arrangement of tandem stages, one stage for each of the eight digits.⁸ Each coding stage produces two outputs, a digit output and a residue output. The residue is the input to the next tandem stage. This arrangement is illustrated in Fig. 6. The tandem stage method can be applied to any logarithmic compression coding characteristic including linear coding and piecewise-linear approximation to logarithmic compression laws.

* A more detailed discussion of coder types can be found in Ref. 7, pp. 583-592.

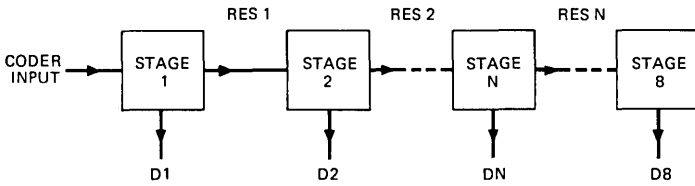


Fig. 6—Tandem stage coder.

Residue and digit output characteristics for the first three coding stages are plotted versus the amplitude of the PAM sample, E_{in} , in Fig. 7. Stage 1 determines the polarity of the sample, D1. Its residue is proportional to the magnitude of the sample. Each of the remaining

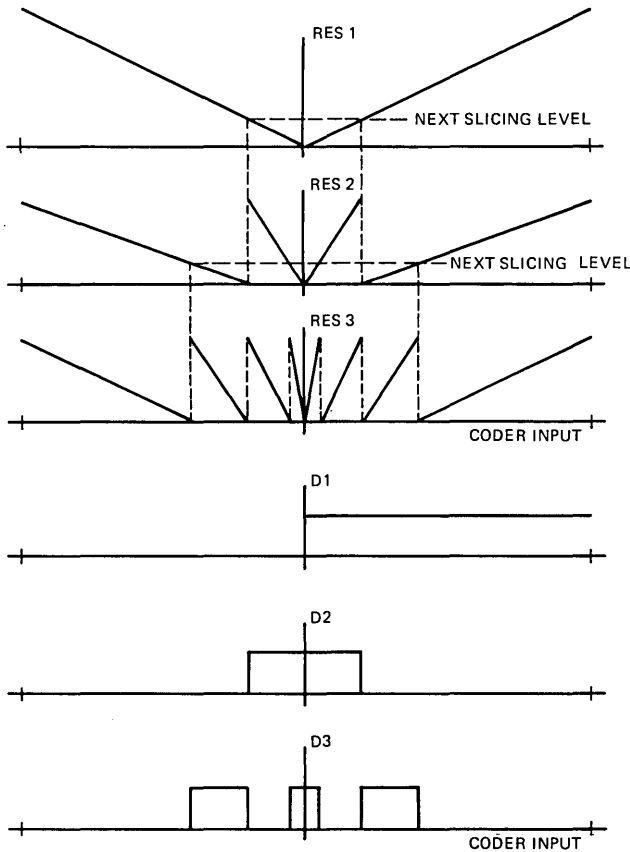


Fig. 7—Residue and digit characteristics.

stages produces one digit of the binary code representing the magnitude of the sample. The compression characteristic is generated by controlling the break-points and relative slopes of the residue segments.

Logarithmic compression coding is possible with this technique because the logarithmic curve has the property that it is congruent to itself by magnification and translation.

5.2 Polarity Stage

The first coding stage, the polarity stage, determines the polarity of the sample to produce the sign digit D_1 , and full-wave rectifies the PAM sample amplitude to produce the residue. The configuration of this stage is shown as Fig. 8. The stage consists of a high-gain inverting amplifier with nonlinear feedback and a digit detector. It is a combination half-wave rectifier and precision slicer. Assuming the high-gain inverting amplifier and no reverse conduction in the diodes, any input current, $I_{in} = E_{in}/R$, at the summing node of the amplifier, must result in conduction through one of the two feedback paths which contains either diode D_1 or D_2 . Any input voltage E_{in} greater than zero results in conduction through D_2 while an input voltage less than zero results in conduction through D_1 . The voltage and current relationships for the stage are shown in Fig. 8. The very steep slope of E_d in the vicinity of zero E_{in} results from the $V-I$ characteristic of the feedback diodes. For the silicon diodes used in the

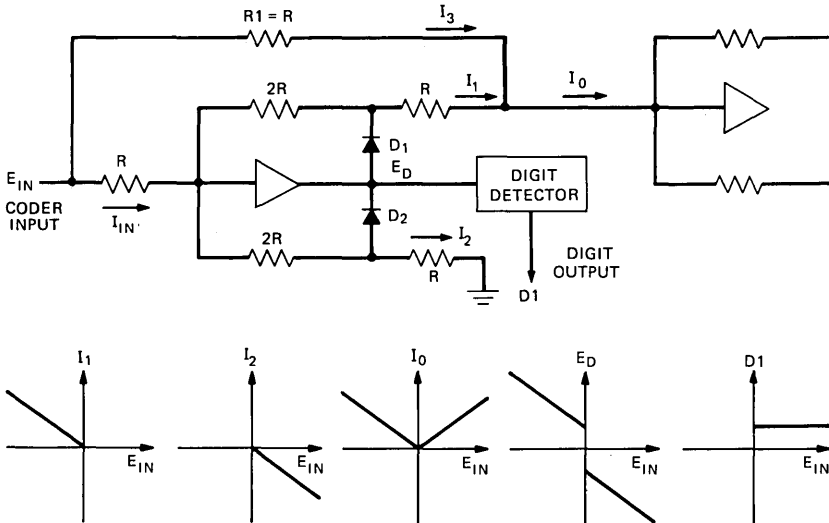


Fig. 8—Polarity stage.

coder, this slope is greater than 200 mv per microamp of diode current. The digit detector, which is a coarse threshold detector, senses the polarity of the voltage E_d and produces the one or zero code output for D1. The full-wave rectifier characteristic is produced by combining the current I_3 transmitted through the path of resistor R1 with the half-wave rectified current I_1 at the summing node of the next coding stage.

5.3 Binary Stage

A typical binary coding stage is shown in Fig. 9. The input current I_{in} (I_{in} is the output current or residue of the previous tandem stage) is combined with a reference current I_{ref} . I_{ref} defines the slicing level of the half-wave rectifier stage. As in the polarity stage, any net current at the summing node must result in conduction through one of the two feedback paths. Input currents greater than the reference current result in conduction through D_2 while input currents less than the reference result in conduction through D_1 . Voltage and current relationships for the stage are shown in Fig. 9. As in the polarity stage, a digit detector senses the polarity of the voltage E_d , and produces the one or zero code output for the digit. The digit detector also controls the switching of a reference current, I_p , which, when added to I_1 and I_2 at the summing node of the next stage, produces the required residue characteristic for binary coding.

The binary coding stage is a precision slicer producing a binary digit with one-zero transitions at $I_{in} = I_{ref}$, and producing a sawtooth

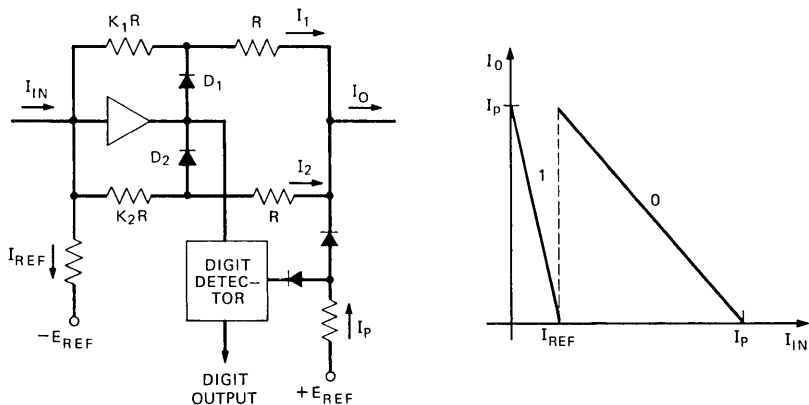


Fig. 9—Binary coding stage.

residue with different slopes to give the desired compression characteristic. For an arbitrary μ , the gain ratio and reference current for the first stage after the polarity stage are given by:

$$\text{gain ratio} = \sqrt{1 + \mu}$$

$$I_{\text{ref}} = I_p \frac{1}{1 + [\text{gain ratio}]}$$

where I_p is the peak current.

For the next stage, the gain ratio is the square root of that of the previous stage or

$$(1 + \mu)^{\frac{1}{2}}.$$

and

$$I_{\text{ref}} = I_p \frac{1}{1 + (1 + \mu)^{\frac{1}{2}}}$$

and so on. The gain ratios and reference currents for the 15-segment $\mu = 255$ compression characteristic are summarized in Table II.

The recurrence relationship between each succeeding gain ratio described above is discontinued at the fifth stage. Had this been continued, an exact $\mu = 255$ compression law would result. By forcing the gain ratios to be one for stages five through eight, a piecewise linear approximation to the μ -law results. Since the only gain ratios used are powers of two, the resulting coder step sizes are related to each other also by factors that are powers of two. This is due to the choice of μ such that

$$1 + \mu = 2^{2^N}. \quad \text{For } \mu = 255, N = 4.$$

5.4 Coder Timing

The step discontinuities in the residue output of each binary coding stage correspond to transitions of the binary digits. These discontinuities

TABLE II—CODING STAGE PARAMETERS

Stage	Reference I_{ref}/I_p	Gain Ratio K_1/K_2
2	1/17	16
3	1/5	4
4	1/3	2
5-8	1/2	1

are produced by the digit detectors switching reference currents I_r during the coding interval. Sample amplitudes near a transition are likely to cause a change in digit output, and a switching of the reference current at the last coding instant. All succeeding stages must then respond to this transient and settle to new outputs. Use of the resultant code word during this transient interval would result in coding errors. It is this response time that limits coding speed.

These errors are prevented by sequentially clocking digit detectors for the binary stages to define the time during the coding interval when each digit output can change state. These times are indicated in Fig. 10. The $1.3 \mu\text{s}$ coding interval is divided into eight 163-ns phases defined by the 6.176-MHz clock. Coder control (CC) pulses of various widths clock the digit detectors: when a coder control pulse is high, the digit output can change state, and when it is low, the digit output is inhibited from changing state.

The detector for D2 is allowed two clock phases to make its digit decision and is then inhibited by CC2 from making further changes, the detector for D3 is given three clock phases to make its digit decision and is then inhibited by CC3 from making further changes, the detector for D4 is given four clock phases to make its digit decision and is then inhibited by CC4, etc. This clocking sequence is followed through D6. D7 and D8 are inhibited simultaneously at Phase 7. Since the polarity

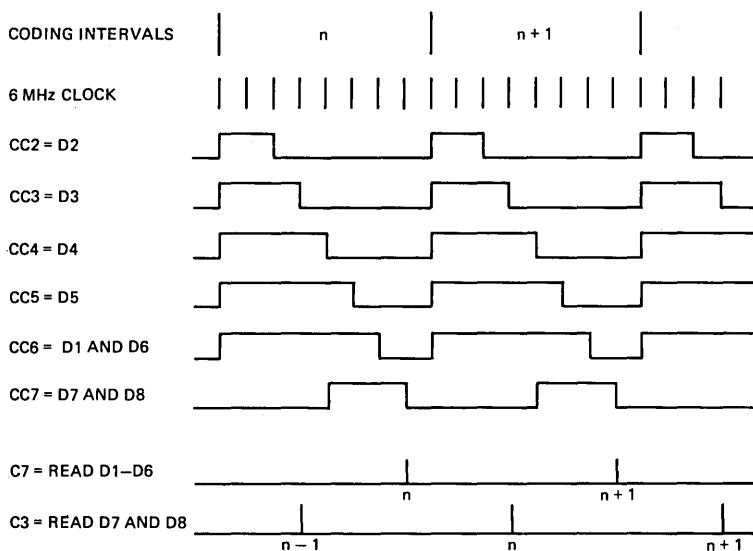


Fig. 10—Coder timing diagram.

digit, D1, has no effect on the residue passed forward to the binary coding stages, it is given until Phase 6 to make its digit decision and is clocked by CC6.

Digits are read into the appropriate coder output processor by phases of the 6-mHz clock. D1 through D6 are read out at Phase 7 by C7. At that time, D7 and D8 have just been inhibited by CC7 and may not have settled to a solid one or zero logic status. Read out of these digits is delayed until the next available processor read time, C3 (each of the four processors operates at one fourth of the 6-mHz clock rate). D7 and D8 are stored in the digit detectors until Phase 4 of the next coding interval by inhibiting them with CC7.

5.5 Coder Accuracy

Numerous coder parameters determine the accuracy to which input samples are coded. Among these are:

- (i) reference voltage supplies,
- (ii) reference resistors,
- (iii) ratio of amplifier gain resistors,
- (iv) summing node bias current and offset voltage.

The significance of the coding error introduced by these parameters depends on which coding stages are affected and the amplitude of the sample being coded. Obviously, errors introduced in the input coding stages are more significant than comparable errors in the latter stages, and a given error is most significant when the coding step size is smallest.

The effect of coding errors can be illustrated by a simple example. Current levels in the coder are chosen such that the peak input voltage to the coder, E_p , results in a peak residue current $I_p = 7$ ma. At the input to the second coding stage, the smallest step sizes, those on the inner segment, correspond to a residue current of approximately $1.7 \mu\text{a}$. The effect of a current error of $3.4 \mu\text{a}$, which could be introduced by summing node offset at the input of this stage, is illustrated in Fig. 11. This error has resulted in the omission of four code words at the origin of the transfer characteristic. This abrupt step in the characteristic would result in excessive idle channel noise and crosstalk as well as degraded gain tracking and distortion performance in message channels.

The magnitude of coding errors could be controlled by placing stringent and costly stability requirements on reference voltage supplies, resistors, and summing node offset. Instead, two simple automatic zero-set loops are used to maintain alignment of the coding segments adjacent to the origin.

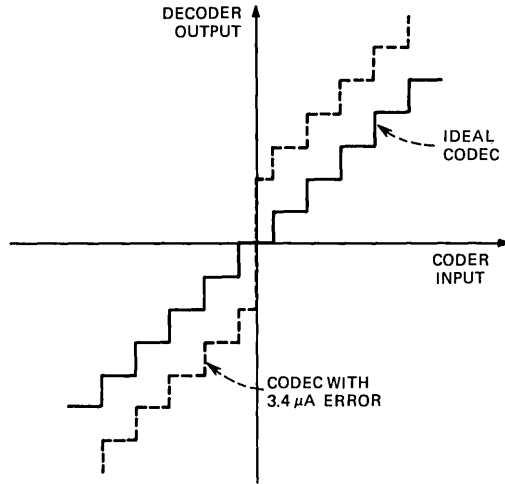


Fig. 11—Coding errors.

Automatic zero-set loops are operated in a housekeeping coding interval set aside at the end of every two frames, or $250 \mu\text{s}$. The coder input during this interval is taken to be zero (although the input from the transfer gates need not be zero volt). The first loop, AZS1 in Fig. 12, samples D1. If the digit is a one, the zero set current into the polarity stage is adjusted to move the digit toward a zero. When the digit output

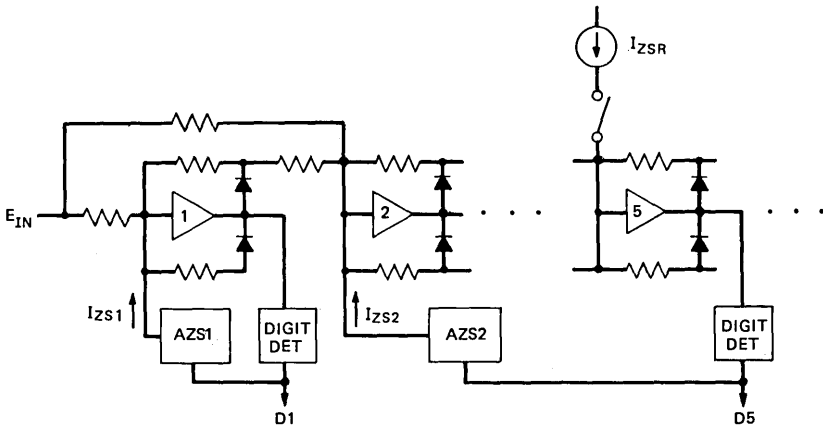


Fig. 12—Automatic zero-set loops.

is zero in the zero set interval, the current is adjusted to move the digit toward a one. The change in I_{zs1} between sampling intervals (250 μ s) is held to a fraction of the smallest coding step. Thus, this loop dithers the polarity stage at its one-zero transition.

Since the polarity stage is at its transition, the residue is at the tip or folding point of the full-wave rectifier "vee" characteristic. This information is vital to the operation of the second loop, AZS2. During the zero-set interval a reference current, I_{zsr} , equal to 7-1/2 steps, corresponding to the transition of D5, is switched into the node of stage 5. This current is not present during normal message channel coding intervals. With this input current and the polarity stage at its transition, stage 5 *should* be at its one-zero transition. By injecting a correction current, I_{zs2} , at the second stage the AZS2 loop, sampling D5 and functioning similar to AZS1 dithers stage 5 around this transition.

During the zero-set interval, the coder is forced to the digit-five transition at the middle of the innermost segment independent of any errors introduced by stages two through four. When I_{zsr} is removed during normal coding intervals, the segments adjacent to zero input are aligned. Any error in coding on these segments is due to stages five through eight. These stages resolve the residue input to stage five into 16 uniform steps, a task requiring only modest accuracy. Coding errors that would have resulted from static imperfections in stages one through four without the zero-set loops have been shifted from the inner segments to the segments where the step sizes are larger and a fixed coding error is less significant relative to the step size. The current I_{zs1} serves the function of the slicing reference for stage 2 and precisely cancels any summing node bias current of this stage.

Each zero-set loop (Fig. 13) is a simple bang-bang servo with an integrator. The state of the digit is clocked as the input to a one-shot multivibrator during the zero-set interval. If the one-shot does not

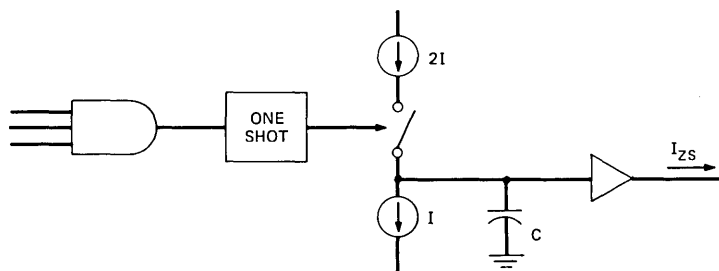


Fig. 13—Zero set.

operate, current source I discharges the voltage on capacitor C . When the one-shot does operate, current source $2I$ is connected to charge C with a net current equal to I . Loop parameters I and C are chosen such that the change in I_{zs} between sampling intervals is a fraction of the smallest coding step.

Precision resistors determine the gains of the coding stages and, with reference voltage supplies, determine fixed and switched reference currents. These resistors define the coder compression characteristic. Thin-film tantalum nitride resistor networks are used to ensure coding accuracy. The ratio of two resistors determines the gain of each path of a coding stage. These gain resistors are fabricated in networks consisting of four resistors, two for each gain path of the stage. A typical resistor network is shown in Fig. 14. Gain resistors are specified to have a ratio tolerance of ± 0.1 percent for the input stages to 2 percent for the latter stages. Reference resistors are fabricated as networks consisting of the fixed and switched reference for each summing node. Absolute tolerance requirements range from ± 0.25 percent to 2 percent. Computer simulations and actual measurements have confirmed the necessity and the adequacy of these requirements.

5.6 Zero Code Suppression

In order to insure adequate timing information for the T1 line repeaters, the D1 Channel Bank suppresses the all zeros code. In so doing, the density of ones is at least one eighth, and the longest run of zeros is fourteen. To comply with the same constraint, the all zeros code in D2 is also suppressed and replaced by the code word 00000010.

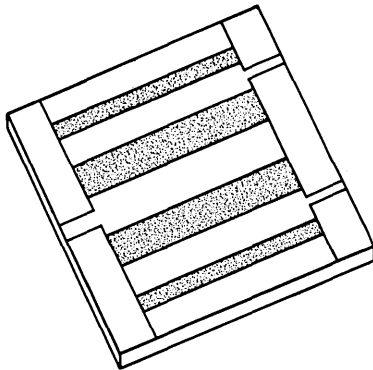


Fig. 14—Thin-film gain-resistor network.

The next to the last digit is changed to a one because the last digit is reserved for signaling in one sixth of the frames.

It has been argued that, since the average density of ones in D2 is higher than that in D1, and since the probability of code words with sparse ones following one another is much smaller than that in D1, zero code suppression is not really necessary. It is included because the additional complexity is very small.

VI. DECODER

The nonlinear decoder has an expansion characteristic that is the inverse of the 15-segment approximation to the $\mu = 255$ compression characteristic of the coder. The decoder is time-shared over four asynchronous digroup inputs. Words to be decoded are written in parallel into a register in the decoder by one of four decoder input processors. The word remains in the store for the approximately $1.2 \mu\text{s}$ decoding interval until reset by the decoder clock. Digits D1 through D5 are decoded in tandem stages very similar to those used in the coder. Because these digits represent linear divisions within segments, and because accuracy is no longer critical, digits D6 through D8 are decoded by summing binary-weighted currents at the node of the first tandem stage, which is for D5.

A typical decoding stage is shown in Fig. 15. The binary state of the digit controls a switched reference current, I_p , at the input of the

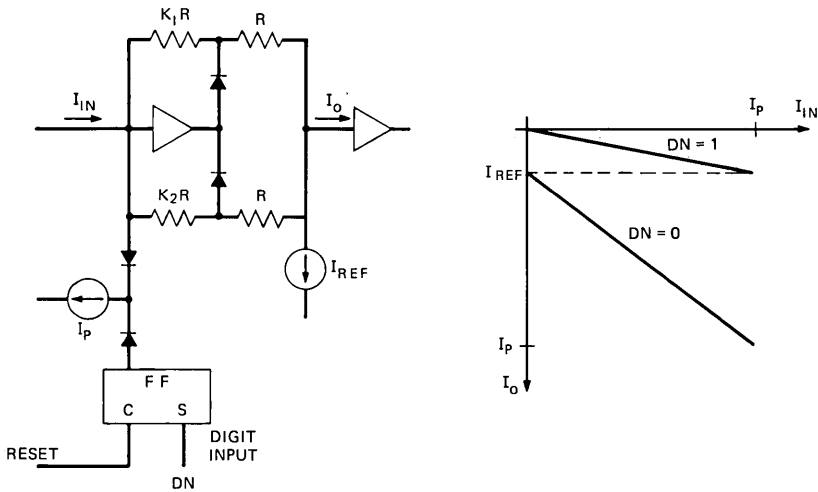


Fig. 15—Typical decoder stage.

stage and determines which gain segment is used for the stage. The gain ratios of these stages are the inverse of the gain ratios for the corresponding coder stages given in Table II. A fixed reference current, I_{ref} , shifts the output to produce a residue of constant polarity. This output, I_0 , is the input to the next decoding stage.

Zero set is used in the final decoder stages in a manner similar to that for the coder. A known signal is used as input to the last two stages. Although the decoder does not make decisions, the equivalent digit output point, which is the junction of the output of the operational amplifier and the steering diodes, is tested for polarity and the zero set correction currents are adjusted according to that result.

VII. SIGNALING

During signaling frames (every sixth frame of PCM words) only seven data digits contain information about the amplitude of the sample: The eighth data digit contains signaling information. Different values are produced by the decoder for 7-digit and 8-digit words in order to minimize quantizing error. This is illustrated in Fig. 16 which is

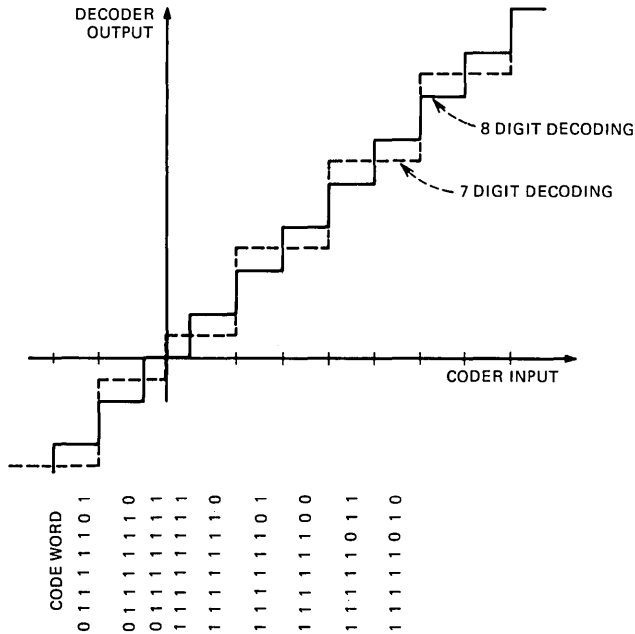


Fig. 16—7- and 8-digit decoding.

an overall CODEC transfer characteristic in the vicinity of zero input for 7-digit and 8-digit decoding.

VIII. PERFORMANCE

The measured performance of a production coder-decoder combination (CODEC) is shown in Fig. 17. Sine waves are generally used as input to measure the performance of a typical production CODEC rather than Gaussian noise which has a distribution that resembles speech more closely. Sine waves are not only easier to generate and measure but they also have the property that localized errors can be detected in the coding and decoding process with greater sensitivity and accuracy. A Gaussian-distributed signal would exhibit a smoother curve over the same defects. This difference is particularly evident in the theoretical signal-to-noise performance of an ideal CODEC when sine waves are used as inputs. This is illustrated in Fig. 18, where each step of the CODEC transfer characteristic manifests itself in a cyclic oscillation in the signal-to-distortion curve, and each segment of the piecewise linear approximation manifests itself in a cyclic oscillation

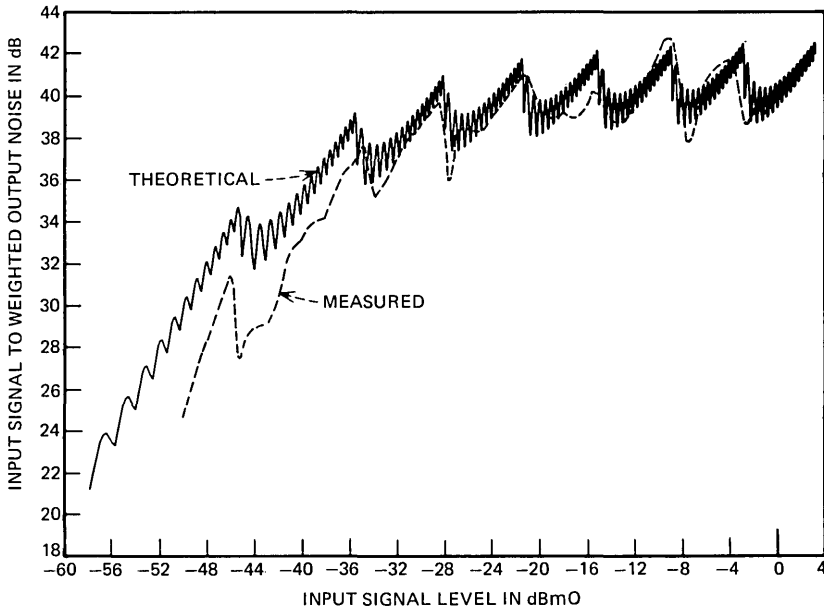


Fig. 17—8-digit, 15-segment, μ -law, non-uniform CODEC. Comparison of theoretical and measured signal-to-noise ratio for sine wave input.

of the envelope. With a Gaussian signal such fine structures are not evident. Furthermore, the Gaussian signal, because of its high peak factor, exhibits the overload characteristic early so that the accuracy of the outermost segment cannot be tested.

The theoretical signal-to-distortion curve is indicated in Fig. 17 for comparison with the measured performance. It can be seen that the typical performance of a production CODEC is very close to that of an ideal CODEC. Due to noise and crosstalk in the D2 bay the measured performance of the CODEC in the D2 Bank is not as good as that of the CODEC alone.

A photograph of the transfer characteristic of the CODEC is shown in Fig. 19. Only the transfer characteristic of the innermost segment is shown since this is the most critical area of the CODEC performance, where offsets will be most evident and errors in step size are most pronounced. If the entire transfer characteristic is plotted, only the step size for the outer segments will be obvious. The inner segments, being only one-128th the length of the outermost segment, will not exhibit visible steps at all.

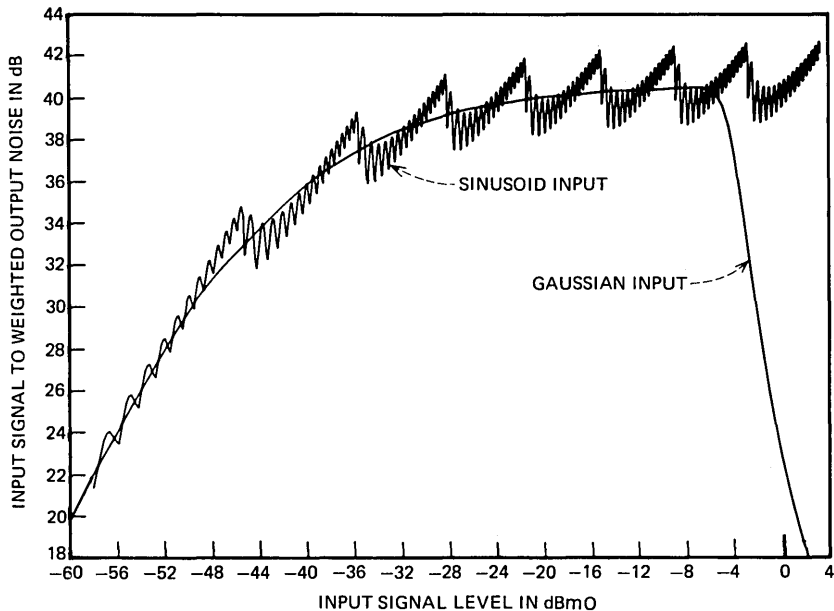


Fig. 18—8-digit, 15-segment, μ -law, non-uniform CODEC. Comparison of theoretical signal-to-noise ratios for sine wave and Gaussian inputs.

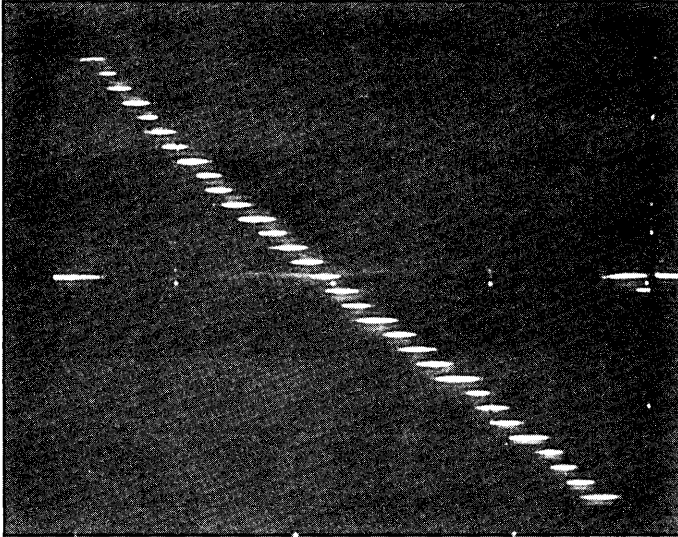


Fig. 19—CODEC transfer characteristic—only the inner two segments are shown representing $1/255$ or 0.4 percent of the total range.

IX. SUMMARY

This article has discussed the analog multiplexing and coding aspects of the D2 Channel Bank. Multiplexing and demultiplexing are accomplished in two stages to ease the timing and crosstalk problems. Coding is accomplished by a digit at a time stage-by-stage coder. Its performance is very close to that expected of an ideal CODEC.

REFERENCES

1. Cirillo, A. J., and Thovson, D. K., "D2 Channel Bank: Digital Functions," B.S.T.J., this issue, pp. 1701-1712.
2. Smith, B., "Instantaneous Companding of Quantized Signals," B.S.T.J., 36, No. 3 (May 1957), pp. 653-709.
3. Purton, R. F., "Survey of Telephone Speech-Signal Statistics and their Significance in the Choice of a PCM Companding Law," Proc. IEEE, B109 (January 1962), pp. 60-66.
4. Cattermole, K. W., *ibid.*, pp. 485-487.
5. Anderson, E. J., "Considerations in Selection of $\mu = 255$ Companding Characteristic," IEEE ICC Record, 1970, pp. 7-19.
6. Dostis, I., "The Effects of Digital Errors on PCM Transmission of Companded Speech," B.S.T.J., 44, No. 10 (December 1965), pp. 2227-2243.
7. Members of the Technical Staff, Bell Telephone Laboratories, *Transmission System Design*, Fourth Edition, Winston-Salem, N. C.: Western Electric Company, 1970.
8. Dammann, C. L., "An Approach to Logarithmic Coders and Decoders," NEREM Record, 1966, pp. 196-197.

D2 Channel Bank:

Digital Functions

By A. J. CIRILLO and D. K. THOVSON

(Manuscript received June 22, 1972)

This article describes the generation of timing signals and processing of digital information in the D2 Channel Bank. The transmitting portion of the D2 Channel Bank, which has four digital outputs, operates under the control of a single synchronous timing circuit. Because the four digital inputs to the receiving section of the D2 Bank are generally asynchronous, independent timing circuits are used for each of the four inputs. In addition, a separate clock is used in the receiver section to operate a single decoder shared by the four digital inputs.

Digital processing of the transmitting terminal includes the serializing of the coder output, inserting of signaling and framing information, and converting the binary code into a bipolar format for transmission over the T1 digital line. To perform the inverse operation just mentioned, the receiving portion of the D2 Bank must extract timing information from the received digital signal and recover the framing information so that the decoded PCM words can be properly demultiplexed. In addition, queuing logic must be performed to permit sharing a single decoder among four asynchronous inputs.

I. INTRODUCTION

The digital functions necessary in the D2 Channel Bank include the generation of control signals and the processing of digital information. The four digital outputs, called digroups, of the D2 Channel Bank are synchronous. A single crystal oscillator followed by a countdown chain is used to provide timing information for the operation of the entire transmitting section. The four incoming digital signals, however, are not expected to be synchronous. Independent extractions of the line clock for each of the four incoming digroups must be provided. This is followed by four independent countdown chains. Since a single decoder will be shared among the four incoming digroups, an independent timing supply is used for decoder operation.

Signals to be processed in the transmitting section of the D2 Bank include the PCM words from the coder which must be serialized and separated into four outgoing digital streams, and signaling information from the channel units which must be multiplexed along with the PCM code words. To provide for the proper demultiplexing of the information at the distant receiver, framing information must also be added to the signal. Automatic search and verification of the framing pulse is performed by each of the countdown chains to insure the proper demultiplexing of the digital signals. Asynchronous sharing of a single decoder and a fast framing search procedure represent unconventional approaches taken in the system design of the D2 Channel Bank. This article will describe circuits that perform the above two functions as well as the more straightforward operation of timing and digital processing.

II. TRANSMITTING SECTION TIMING

2.1 *Requirements*

Operationally, the 96 channels are divided into four independent digroups of 24 channels each. However, the four digroups of the D2 Channel Bank are treated synchronously on the transmitting side. They share a single timing generation circuit. The requirements for transmitting timing are dictated by both the output format and circuitry requirements. First, channel pulses occurring at an 8-kHz rate are necessary for the operation of the sampling gates. Four channels, one from each digroup, are sampled simultaneously. Thus 24 channel pulses, staggered in time and placed on separate leads, are required for the first stage of multiplexing. The first multiplexing stage results in eight groups of 12 channels each. The samples of four of these groups are staggered with those of the other four. Second, these samples are sequentially transferred to the coder by transfer gates operated by group pulses. Thus eight group pulses are needed for this final multiplexing with each of the eight group pulses repeating at a 96-kHz rate. Third, the digit-at-a-time sequential nature of coder operation requires control pulses. Six coder control pulses, each on separate leads, and each repeating at 772 kHz are necessary. In order to provide for proper phasing of the coder control pulses, a crystal clock of 6176 kHz is used at the head of the countdown chain. Fourth, the output format requires the insertion of the framing pulses at a 2-kHz rate, and requires the insertion of signaling framing pulses repeating at 667 Hz. Finally, the output bit streams applied to the four transmission lines require a clock of 1544 kHz.

2.2 *Timing Circuitry*

Figure 1 is a block diagram of the countdown chain used to produce the clocks in the transmitting section. The major use of each frequency is also indicated in the figure. Combinations of pulse trains produced by this chain are used wherever necessary to provide proper phases and duty cycles. The lowest frequency of 667 Hz provides for the

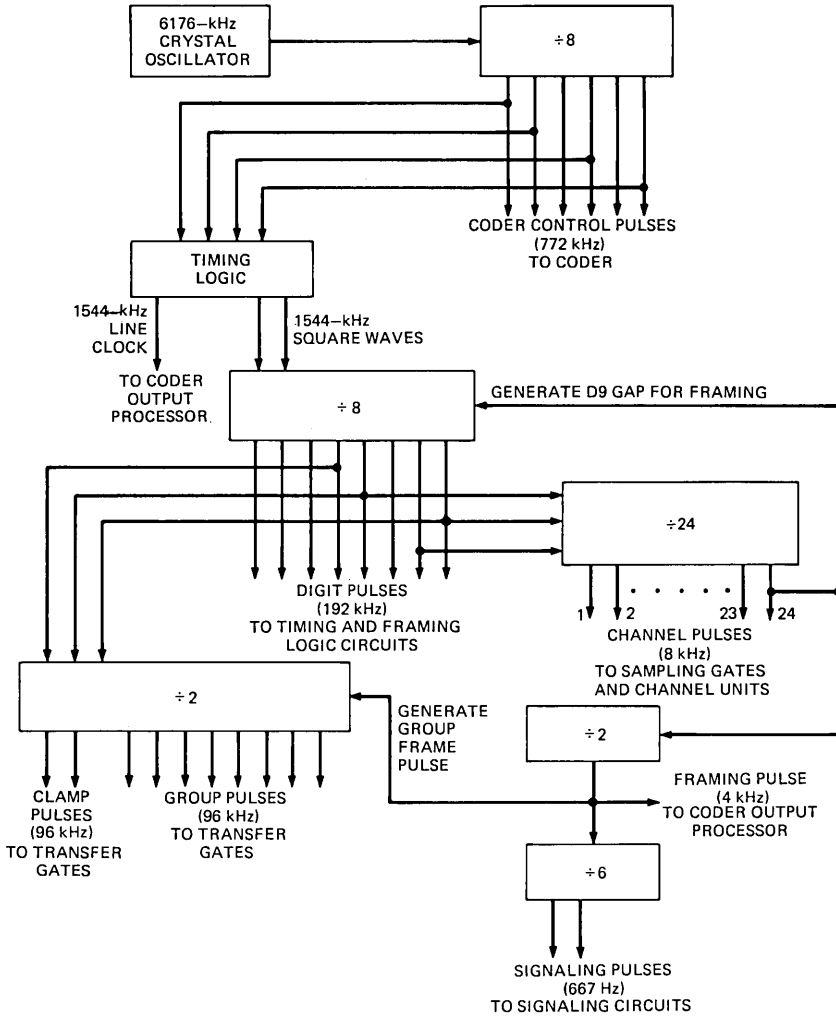


Fig. 1—Transmitting timing: a block diagram of the countdown chain used to produce the clock in the transmitting section.

insertion of the signaling framing bit, which is a repeating pattern of 111000, and controls the insertion of the signaling bits for all 96 channels. The 4-kHz signal controls the insertion of the main framing pulses which have the alternating pattern of 1010. The main framing bit and the signaling framing bit alternately occupy the 193rd bit in each 125-microsecond frame of each of the four digroups.

Even though the first stage of multiplexing is accomplished in groups of 12, 24 channel pulses are necessary because while a sample from one channel is obtained and held for further multiplexing, the sample from another channel must be clamped in preparation for a new sample. This requires that two staggered sets of 12 channel pulses be generated. Figure 2 indicates a portion of the timing near the framing pulse. The straightforward counting down of the frequency is modified by the insertion of a gap for zero setting of the coder. As mentioned in a companion article,¹ accuracy of the coding is achieved in part by a zero-setting circuit operating on the output of the coder when the coder is connected to a known reference signal. The time for this house-keeping chore is the gap which is present in the output frame format when the framing pulses are being inserted. This time, which is about 650 nanoseconds, is only half the interval that the coder normally takes to code a sample. To provide for a full 1300-nanosecond gap for the zero-setting circuit, the coding operation is offset in time by 650 nanoseconds every other frame so that two framing intervals can be lumped together. This also requires that the timing offset must be removed at the output so that framing pulses can be properly inserted for transmission.

III. RECEIVING SECTION TIMING

Much of the receiving section timing is similar to the transmitting section timing to the extent that channel pulses, framing pulses, and signaling framing time must all be derived. Since the receiving timing chain should be synchronous with the incoming signal, and since four independent incoming digroups are expected, four independent count-down chains are used in the receiver. Each of these four countdown chains is driven by a 1544-kHz clock extracted from the incoming line (Fig. 3). Unusual aspects of the receiving timing circuitry are the framing search procedure, and the operation of a shared decoder.

3.1 Framing

Framing is necessary at the receiver to bring the timing generation at the receiver into phase with respect to the incoming line, so that the digits can be properly identified for decoding and demultiplexing.

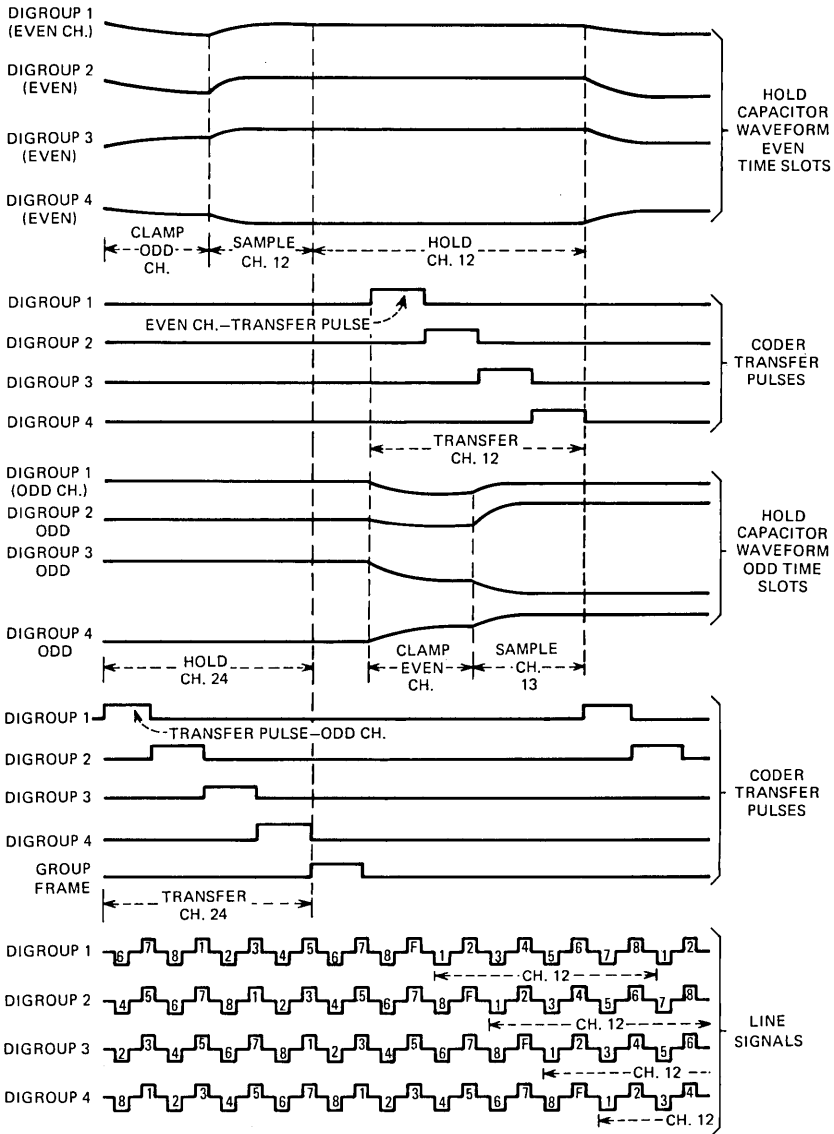


Fig. 2—Transmitting timing; a portion of the timing near the framing pulse.

Framing is accomplished by searching and verifying the framing pattern which was inserted at the transmitting end. The requirement on a framing circuit is that it should accomplish the search for the framing pulse within a certain time. Whenever the synchronization

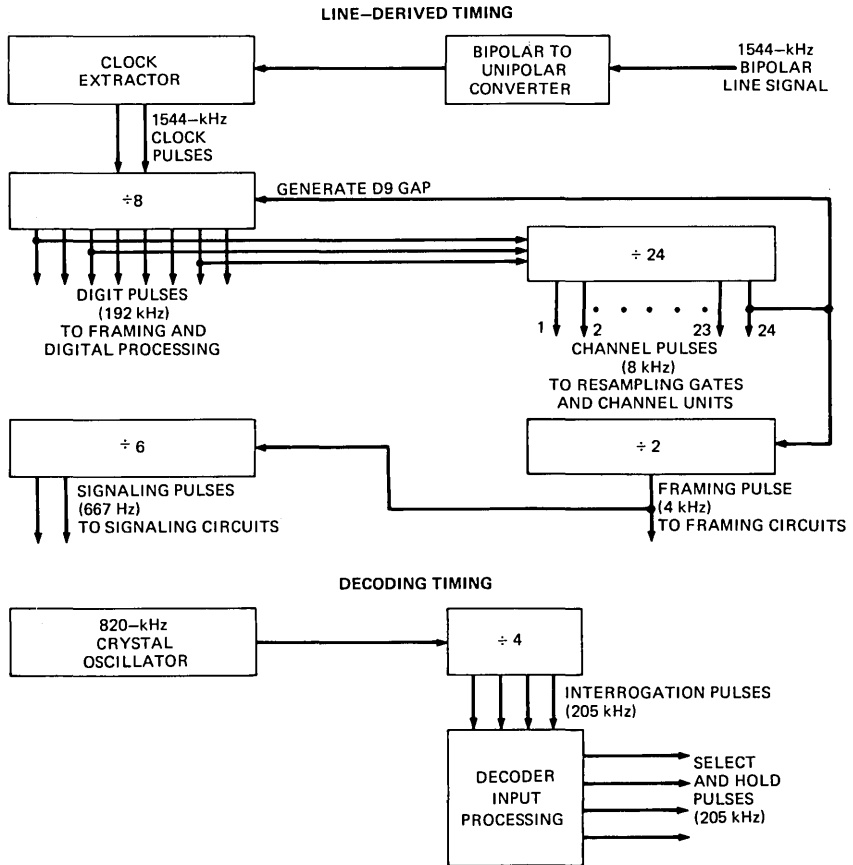


Fig. 3—Receiving timing: each of four countdown chains is driven by a 1544-kHz clock extracted from the incoming line.

of the receiver with the incoming signal is inadvertently lost due to error conditions, the signaling information will be incorrect. The requirement that reframes be accomplished within 50 milliseconds insures that although the telephone customer will receive no useful signal during this 50 milliseconds, at least he will not be disconnected. Because framing pulses in D2 are inserted every 250 microseconds as compared to every 125 microseconds for D1, a simple bit-by-bit search strategy, such as that used in D1, will result in a reframe time up to 200 milliseconds long. Therefore, a more sophisticated reframe procedure is used in D2.

As is usual in framing circuits, there are two modes of operation. The normal mode is the in-frame mode where the framing bit is checked for the alternating 1010 pattern. Occasional deviations from this pattern are ignored so that isolated line errors will not cause the framing circuit to initiate a false search. This flywheel action is achieved in the form of a capacitor store where about four closely-spaced errors are necessary to cause the framing circuit to enter the out-of-frame mode.

When a framing circuit enters the out-of-frame mode, the search procedure is initiated. Two 8-bit registers are used by the framing circuit in order that eight bits of the received stream may be examined at a time for possible candidates for the framing pulses (Fig. 4). Normally, these two 8-bit registers are used by the decoder input processor for the queuing operation. This will be discussed later in this article. During the out-of-frame mode, these registers are transferred for use by the reframe circuit. One register is used to store the incoming information bits. This is called the I register. The second register is used to note which of the bits in the I register are still suitable candidates for the framing pulse. This is called the S register. At the start of the reframe procedure, eight consecutive incoming bits are stored in the I register. These saved bits are then compared with eight consecutive incoming bits occurring two frames (0.25 milliseconds) later. A true framing bit should exhibit the alternating 1010 pattern. Any other bit position is assumed not to have this property on the long term but may or may not exhibit this property on the short term. The new eight bits replace the old eight bits in the I register. A comparison is made during this replacement between the old and the new to see if

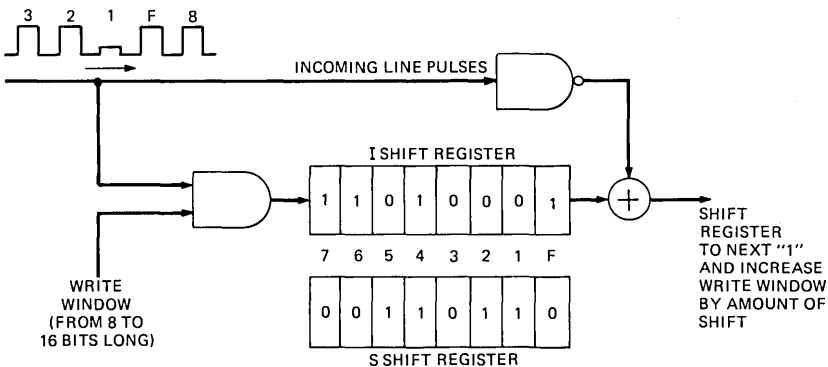


Fig. 4—Two 8-bit registers which are used by the framing circuit so that eight bits of the received stream may be examined at a time for possible candidates for the framing pulses.

any of these pulses are still candidates for the framing pulse. The result is stored in the S register. If the first of the eight pulses fails to qualify as the framing pulse, both the I and the S registers will shift one position allowing the second pulse to occupy the first position and making room for the next pulse at the end. It is thus possible for the registers to shift up to eight bit positions if all eight bits fail to meet the test. No shift takes place, however, if the first position exhibits the 10 alternation, even though the succeeding positions do not. This fact, though, is noted in the S register. In so doing, as soon as the first position fails the test, the succeeding bit positions could be rapidly passed over.

During the out-of-frame mode, a single detected violation of the alternating pattern will disqualify a particular bit position as a possible framing pattern. Thus line errors could cause a legitimate framing pulse to be passed over. The reframe time will exceed 50 milliseconds when this happens. With the expected line error rates of 10^{-6} or better, such incidences will not be frequent enough to cause concern. When an alternating pattern has persisted for about 2.5 milliseconds during its examination, the corresponding bit position will be considered as the framing pulse. The in-frame mode will then be entered. Again isolated errors will be ignored, and the I and S registers returned to their former function in the decoder input processor. By testing the framing position for 2.5 milliseconds in which time ten checks are made, the incidence of falsely returning to the in-frame mode due to information pulses exhibiting the alternating pattern on the short term will be small. When it happens, the penalty is a slightly increased reframe time since the out-of-frame mode must be re-entered.

By using a window of eight bits, the reframe time of about 43 milliseconds is achieved. The standard deviation is about 5 milliseconds. Increasing the number of bits stored and compared at one time would, of course, reduce the reframe time further. For example, with a pair of 16-bit registers, the reframe time is calculated to be 26 milliseconds. However, the 43 milliseconds of reframe time is adequate, and permits the sharing of these registers with the decoder processor.

IV. CODER OUTPUT AND DECODER INPUT PROCESSING

4.1 *Coder Output Processor*

The coder output processor is basically a parallel to the serial converter. In the D2 Channel Bank the coder output must first be split into four outgoing streams. Thus, there is a coder output processor for each of the four digroups. Parallel outputs from the coder are written

sequentially into one of four shift registers from which the serial information stream is read out. In addition, after 24 channels have been processed, the framing bit is added to the register following the last bit of the 24th channel. As mentioned previously, the coder operates with a timing offset between the even and odd frames to provide for zero-setting. This offset is absorbed by the shift register where extra storage is provided. Signaling information is inserted at the proper frames in the eighth bit position of each word. This insertion takes place in the coder from which it enters one of the four coder output processors.

4.2 Decoder Input Processor

A single decoder is shared by four incoming digroups for decoding the PCM code words into analog samples. A simple queuing logic is used. The decoder has its own clock which operates at a rate higher than four times the incoming rate of each line. Each incoming digroup has a decoder input processor associated with it. These processors use two 8-bit registers. A complete 8-bit word is stored in one of the registers to wait for decoding. During this wait, the second register will be receiving the following serial 8-bit code word. The decoder timing circuit causes the decoder to poll each of the four decoder input processors. Whenever the decoder cycles to a particular processor and there is a complete 8-bit word ready for decoding, these eight bits are transferred to the input register of the decoder.

Since the decoder timing is faster with respect to any of the incoming lines, the decoder will at times find only partly filled serial-to-parallel registers. In this case, the decoder processor will not transfer any code words, and the decoder will rest for that period before moving on to the next processor. This variation in time between the arrival of a complete code word and the decoding time is one complete queuing cycle. Thus, at least two registers must be used in each processor to accommodate this delay (Fig. 5). It can be shown that even with a very sophisticated queuing logic, two registers are still needed for each digroup because the maximum delay cannot be less than three decoding times.

Removal of this delay variation is accomplished in the analog store of the select-and-hold circuit. Channel pulses derived from the incoming line sample the output of the select-and-hold circuit at the proper time position in the demultiplexing process, and thus provide uniformly spaced PAM pulses to the receiving lowpass filter.

Timing for the decoder is provided by a separate 820-kHz oscillator.

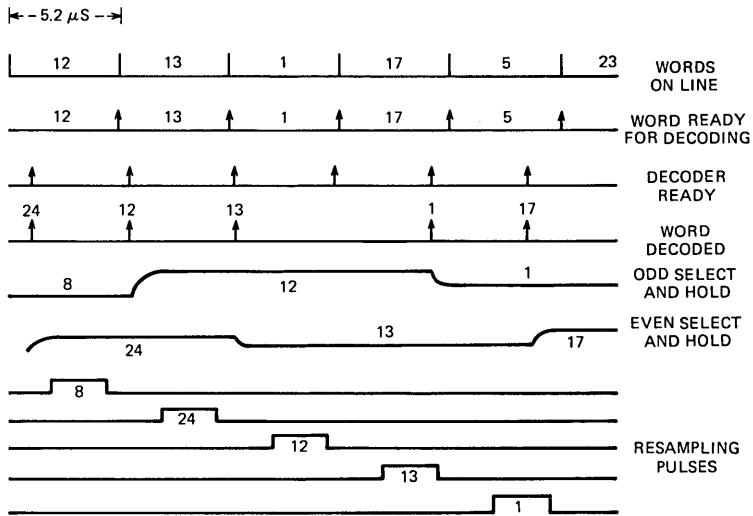


Fig. 5—Asynchronous decoding.

This is divided down to provide four phases of 205 kHz for use in polling the decoder input processors. This frequency must be high enough above the expected incoming word rate of 192 kHz to avoid queuing overflow during higher than normal incoming word rates caused by line jitter, and low enough to provide ample time for decoder settling. Another factor affecting the choice of the rate is the desire to minimize the effects of any beating frequencies between the incoming word rate and the decoding rate. The choice of 205 kHz meets these criteria.

Whenever the decoder becomes idle because an 8-bit word has not yet fully arrived at the input processor for digroup 4, the decoder will use this interval for zero setting.¹ With a line rate of 192 kHz, and a decoder cycling rate of 205 kHz, the rate of occurrence of no decoding is approximately 13 kHz for each digroup.

V. TRANSMITTING AND RECEIVING COMMON SIGNALING

Common signaling circuits in both the transmitting and receiving sections of the D2 Bank control the routing of signaling information from the channel units to the output bit stream and back. The signaling states of a particular channel are sampled by the appropriate channel pulse at the channel unit. Two signaling paths, allowing up to four

signaling states for each channel, result in two samples from each channel unit. These samples are multiplexed on two buses and sent to the transmitting common signaling circuit. The above action repeats every frame. During the sixth frame of a twelve-frame cycle, and only then, samples on one of the buses are sent to the coder where they take the place of digit 8, the least significant information pulse. During the twelfth frame, samples from the second bus are similarly sent to the coder.

To allow the receiving end to properly identify these signaling digits, signaling framing information must also be inserted in the 193rd position of every other frame. These bits follow the pattern 111000, completing a cycle every twelve frames. The frame following the signaling framing bit change from zero to one is defined as the frame containing signaling information in the eighth bit of every word for signaling path A. The frame following the one-to-zero change contains signaling information for signaling path B. In this manner, two signaling paths per channel are provided, each with a signaling capacity of 667 Hz.

At the receiving end, the frames containing the signaling bits are identified by testing the signaling framing bit. The locations of the signaling framing bit are indicated by the main framing pulse which shares the 193rd time slot in alternate frames. As with other receiving timing circuits, each digroup contains a separate receiving common signaling circuit. Each circuit extracts the eighth digit of each word during the appropriate frame, and places it on one of two buses, one for signaling channel A and the other for signaling channel B. The appropriate channel pulse gates the information on the buses to the individual channel-unit flip-flops which, in turn, drive the relays to reproduce the signaling information of the transmitting end.

The receiving common signaling circuit also controls the operation of the decoder during the signaling frame. Since only seven of the eight digits contain PCM information, the decoder is prevented from using the eighth digit for decoding. During these frames the decoder operates as a seven-digit decoder.

The transmitting and receiving common signaling circuits are designed so that they may be replaced by common-channel interoffice signaling (CCIS) plug-in units. When this is done, the CCIS bit stream will take the place of the signaling framing bits directly. This provides for a 4-kilobit per second channel. At the same time, the substitution of the eighth bit for signaling information will not take place and the full potential of eight-bit coding will be realized. At the receiving end, the decoder will no longer be asked to perform seven-digit decoding.

VI. SUMMARY

This article has discussed the system and circuit aspects of the digital functions in the D2 Channel Bank. The unique features of the D2 Channel Bank in this respect are:

- (i) Circuitry to accomplish fast framing by use of storage,
- (ii) Asynchronous operation of a single decoder shared by four incoming lines, and
- (iii) Signaling circuits that can later be changed for common-channel interoffice signaling systems.

REFERENCE

1. Dammann, C. L., McDaniel, L. D., and Maddox, C. L., "D2 Channel Bank: Multiplexing and Coding," B.S.T.J., this issue, pp. 1675-1699.

D2 Channel Bank:

Power Conversion

By S. D. BLOOM and G. F. SWANSON

(Manuscript received June 8, 1972)

The D2 Channel Bank requires a variety of voltage sources to operate. Some have lenient regulation requirements, while a few have extremely tight tolerance requirements. To derive these voltage sources from a 48-volt central office supply with a minimum of complexity and good conversion efficiency, a tailored design is chosen where, for the less critical voltages, a single switching regulator is used to supply a dc-to-dc converter with several dc voltage outputs and, for the critical voltage sources, series-type regulators are used.

I. INTRODUCTION

The D2 Channel Bank depends on a variety of power supply voltages for its operation. The Power Systems Converter Circuit (PSCC), driven by the minus-48-volt central office battery source provides the seven regulated dc voltages required by the D2 Channel Bank.

By using a high-powered switching-type regulator for primary regulation followed by a dc-to-dc converter for voltage transformation, high efficiency and compactness are achieved. Series-type regulators with high-gain feedback circuits are used to stabilize the voltages for three of the outputs where very stringent voltage requirements exist.

Overvoltage turn down, interlocking of critical voltages, voltage tracking and overcurrent protection are also provided. The switching regulator and dc-to-dc converter are designed to operate at commutating frequencies above 20 kHz to eliminate acoustical noise.

The power supply is packaged to facilitate manufacturing, testing, and maintenance. The unit is arranged to be mounted in the same type rack as the rest of the D2 Channel Bank.

1.1 Requirements

The power supply requirements are shown in Table I. It is seen that

TABLE I—D2 BANK ELECTRICAL REQUIREMENTS

Output Voltage (Volts dc)	Current Range (Amperes dc)	Output Voltage Tolerance (%)	Trouble Voltage (Volts dc)	Maximum Ripple	
				<20 kHz	>20 kHz
+24	2.1- 4.4	±5	+27.6	5 mV p-p	24 mV p-p
-24	0.9- 2.2	±5	-27.6	5 mV p-p	25 mV p-p
-12	1.8- 3.8	±5	-15.0	36 dBrnC	100 mV p-p
-5	2.2- 5.0	±10	*	*	*
+5	3.5- 10.0	±4	+ 5.30	50 mV p-p	50 mV p-p
+32	0.055-0.150	±0.75	+38.0	5 mV p-p	40 mV p-p
-32	0.035-0.100	†	-38.0	5 mV p-p	40 mV p-p

NOTE: The output voltage tolerance, trouble voltage and maximum ripple apply for an input voltage range of 42 to 53 volts dc and an ambient temperature range of 10 to 50 degrees Centigrade.

* No requirement specified.

† The minus-32-volt output must track the plus-32-volt output to within ±0.375 percent.

the first four of the voltages in Table I have relatively lenient requirements, whereas the last two, because they are used to generate reference currents in the coder, have extremely tight requirements. The plus-5-volt supply is also considered more critical than the first four because, although the allowed percentage variation is comparable, the absolute allowable voltage variation is quite small. Furthermore, because that voltage is used to power integrated circuits that are sensitive to overvoltages, greater protection is required.

The plus-24- and minus-24-volt supplies are used primarily for the channel counters, which are realized in blocking-oscillator form. The -5 and -12 are utility biasing voltages used for most DPD transistor circuits.

1.2 Powering Plan

The PSCC is divided into five sections, shown in Fig. 1. Four of these are plug-in units. The four plug-in units connect physically to the converter frame, the fifth section, which is hardwired into the D2 bay. The converter frame contains most of the filtering for the input, and for all of the outputs. It has no active elements. The switching regulator is the primary regulator in this circuit. It converts the unregulated minus-48-volt dc input into a regulated minus-24-volt dc voltage. The regulated minus-24-volt dc is one of the outputs, and is also used as the input for the converter. The converter inverts (converts to ac) the dc input and, by means of step-up and step-down transformers and rectifiers, furnishes six different dc voltages. Three of these voltages, minus 12 volts, minus 5 volts, and plus 24 volts, are

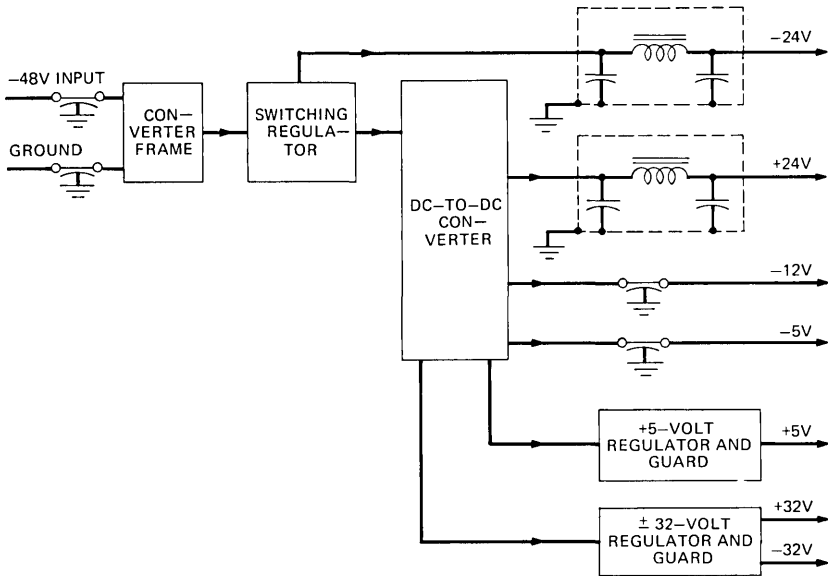


Fig. 1—Block diagram of PSCC.

delivered to the load through the necessary filtering in the converter frame. Thus a single switching regulator is shared by four output voltages. The other three voltages are fed to the inputs of three precision-series-type regulators. The precision regulators deliver plus 5 volts, plus and minus 32 volts to the load. Figure 1 illustrates the relationships of the circuitry among the various units. The filters shown for the input and output leads are physically in the converter frame.

II. REGULATORS AND CONVERTERS

2.1 *Switching Regulator*

The PSCC was designed to use a primary regulator in order to eliminate the need for separate regulators for each of the seven outputs. Only the three outputs with very tight tolerances have separate series regulators, and the power dissipation in these is held to a minimum because of the preregulation provided by the primary regulator. A switching-type regulator is used as the primary regulator in the PSCC because of the high efficiency that can be achieved with this type of circuit.^{1,2}

The switching regulator is illustrated in Fig. 2. The operation of

the switching regulator can be described briefly as follows. A simplified schematic is shown in Fig. 3. The unregulated minus-48-volt central office battery is connected to the regulated minus-24-volt load through a switch and filter inductor. By sensing the output voltage, comparing it against the voltage developed across a reference diode, and amplifying the resultant error signal, the off-time of three switching transistors connected in parallel are controlled to produce the regulated output. The on-time of the switching transistors is controlled primarily by the current in the secondary winding of the current transformer. This secondary current decreases exponentially with time due to the exponentially increasing exciting current in the transformer primary which subtracts from the constant primary current. This secondary current is used to keep the switching transistors in saturation and when the current falls below the level required to keep the switching transistors in saturation, switching occurs and the switching transistors turn off. The dc output voltage is derived by integration of the train of pulses from the switching transistors through the use of a large inductor with a fly-back diode and an output filter capacitor.

To increase the reliability of the switching regulator, the three high-power switching transistors are driven by another power transistor. These are arranged in a Darlington configuration to reduce the switching time and to provide current sharing among the paralleled units. The amplifier section of the switching regulator is located in

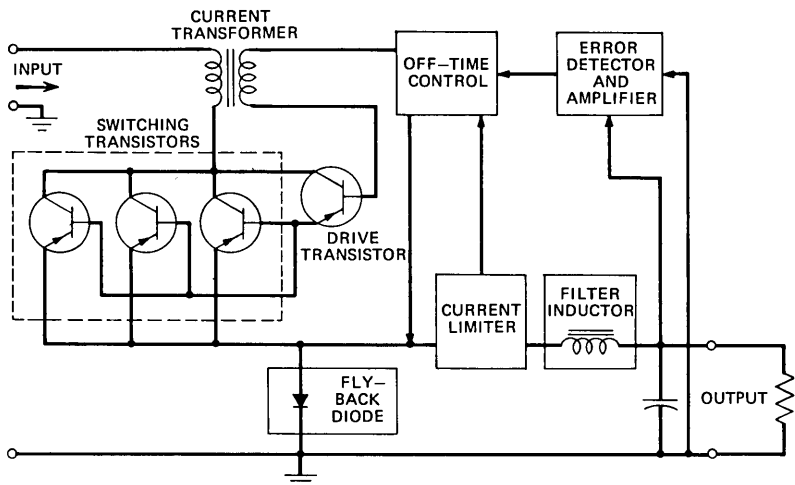


Fig. 2—Switching regulator.

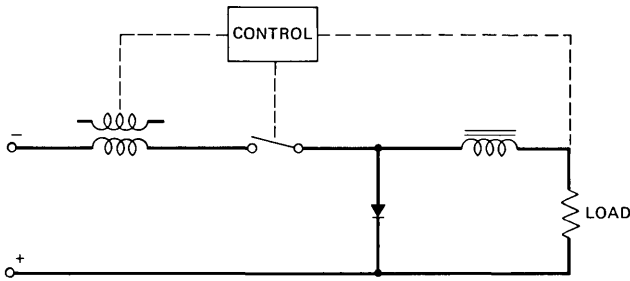


Fig. 3—Simplified schematic of the switching regulator.

an enclosed container within the switching regulator plug-in unit. This feature is used to protect the amplifier from the electromagnetic interference generated in the switching regulator and the dc-to-dc converter. A current limiter is used to protect the regulator against excessive current during turn-on and in the event of an output short.

2.2 DC-to-DC Converter

The dc-to-dc converter is shown schematically in Fig. 4. The converter, which operates from the regulated minus 24 volts dc supplied by the switching regulator, operates at approximately 24 kHz, which is beyond the audible range and so prevents objectionable disturbance to operating personnel. The converter consists of two main sections; the drive oscillator and the power converter.

The switching frequency of the drive oscillator is determined by the volt-seconds required to saturate the base drive transformer used in conjunction with the oscillator switching transistors. The base drive transformer alternately drives the two switching transistors into saturation and cut-off, thereby producing a square wave output across the primary of the oscillator transformer.

The power converter is driven by the square-wave base drive signal supplied by the output winding of the oscillator transformer. The high-power switching transistors in the power converter, therefore, also operate in a class "D" mode and generate a square wave. The switching duration is kept very short, thereby limiting power dissipation in the converter. The square wave generated by the high-power switching transistors is applied across the primary of four output transformers, the five secondary windings of which supply the voltages to the rectifiers and filters for the six outputs.

2.3 Plus and Minus 32-Volt Regulators

The plus and minus 32-volt regulators are high-gain series regulators

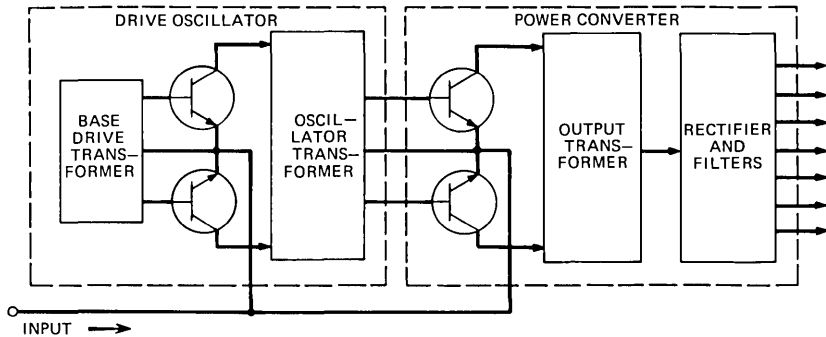


Fig. 4—DC-to-DC converter.

which derive their input from the plus and minus 45-volt outputs of the converter.

Both the plus and minus 32-volt regulators consist of a series-regulating transistor, a current amplifier, a differential error detector and a voltage reference. The series-regulating transistor acts as a variable series resistance and absorbs the difference between the partially regulated dc input voltage, and the highly regulated dc output voltage.

Since the minus-32-volt regulator must track the plus-32-volt output to within ± 0.375 percent, the error detector section of the minus-32-volt series regulator is connected across the plus and minus 32-volt outputs. Therefore, when the plus-32-volt output changes, the voltage change is also sensed by the minus-32-volt error detector. This arrangement causes the minus-32-volt regulator to track any change in the plus-32-volt output.

To protect the regulators in the event of an output short circuit, a zener diode is connected across the series transistor. It protects the transistor by limiting the maximum collector-to-emitter voltage of the transistor to the breakdown potential of the zener, until the 32-volt circuit breakers operate and shut down the regulators.

2.4 Plus-5-Volt Regulator

The plus-5-volt regulator is a high-gain series-type regulator which derives its input from the plus-7-volt output of the dc-to-dc converter through special filtering located in the converter frame. To reduce the power dissipation in the transistors used in the series element, the circuit uses a unique arrangement of paralleled transistors and resistors as shown in Fig. 5.

The operation of the series element can be described briefly as follows:

a portion of the input current flows through transistor Q2 and its series resistor R2, and the rest of the current flows through a parallel branch made up of transistor Q1 and series resistor R1. The voltage drop developed across resistor R1 increases the effective base-to-emitter voltage of transistor Q1 and forces more base current into transistor Q2.

Under normal conditions, most of the input current flows through Q2 and resistor R2. Resistor R2 dissipates most of the power developed across the series element, because Q2 operates near saturation. When the input voltage decreases due to changes in line or load, transistor Q2 saturates, and a greater portion of the input current flows through Q1. However, the voltage across Q1 also decreases so that the power dissipation in Q1 is maintained at a safe level.

A separate bias voltage is required by the plus-5-volt regulator to provide sufficient voltage to operate the differential error detector and voltage reference circuits.

III. PHYSICAL DESIGN

The PSCC is packaged as shown in Fig. 6. The overall width of 21 inches and depth of 12 inches was fixed by the bay requirements. The height of 16 inches was determined by the component density as illustrated later in this article.

Construction in the form of plug-in units was chosen for ease of manufacturing and field maintenance. The plug-in unit concept also contributed to lower cost, which was a design consideration. This type of unit also conforms to the D2 System philosophy of shipping separately all units that are not bay wired. The four plug-in units, i.e., dc-to-dc Converter (CONVERTER), Switching Regulator (REG-

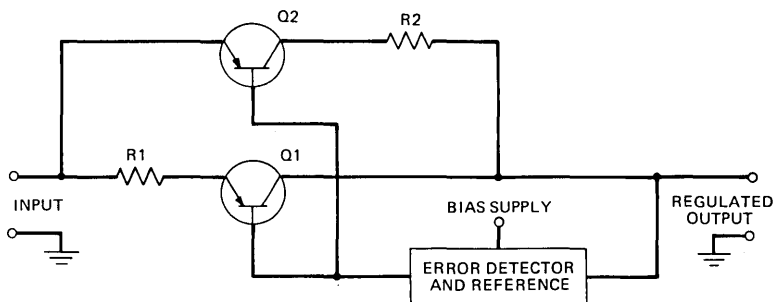


Fig. 5—Paralleled transistors and resistors used in the circuit of the plus-5-volt regulator.

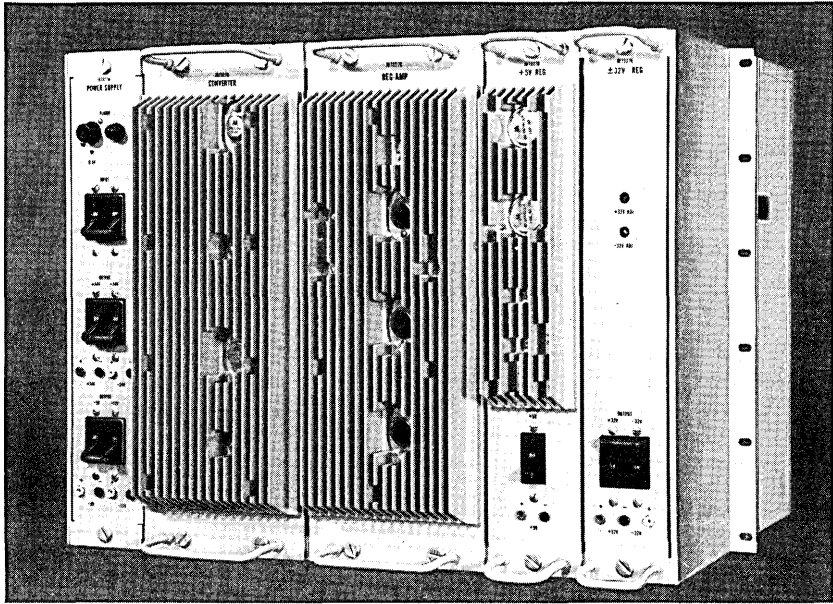


Fig. 6—D2 Channel Bank—J87327A power supply prototype, front/side view.

AMP), plus-5-volt Regulator (+5V REG) and the plus and minus 32-volt Regulator ($\pm 32V$ REG), are of different widths due to the amount of circuitry contained in each one. The largest unit weighs 20 pounds and the smallest unit weighs 4 pounds. Figure 7 shows the front/left/top view of the dc-to-dc converter illustrating the typical compactness of all the plug-in units. The plug-in units are held in place by $\frac{1}{4}$ -turn quick-release fasteners, and handles are provided for easy removal of the units for maintenance. The panel located on the left side of the converter frame contains the main control apparatus and may be released from the converter frame by means of two $\frac{1}{4}$ -turn quick release fasteners for ease of maintenance.

An electromagnetic interference shield is located between the switching regulator and the series regulators to prevent the radiated noise from affecting the circuits in the series regulators.

3.1 Shielding

The stringent noise requirements dictated by the D2 system required the shielding of the noise radiating components. Figure 8 displays the shielded compartment located at the rear of the converter frame with its cover removed. The compartment contains the interconnecting

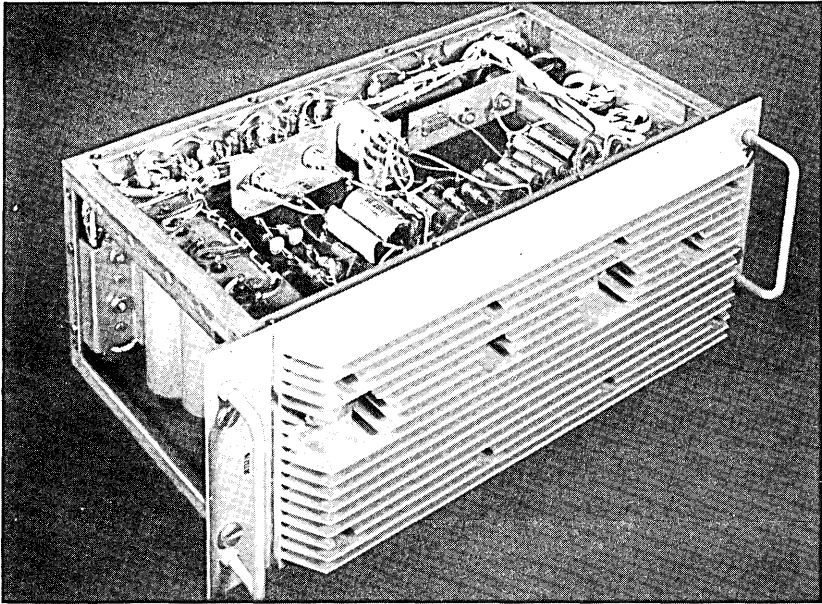


Fig. 7—D2 Channel Bank—J87327B power supply prototype (converter), front/left/top view.

wiring for the plug-in units and the input and output filtering. All leads that enter or leave this compartment are filtered by feed-through capacitors or feed-through filters. These capacitors and filters restrain the longitudinal conducted noise on the input and output leads. The cover when fastened in place will confine any EMI noise radiating from within the compartment.

3.2 *Plug-in Construction*

Figure 9 shows an exploded view of the dc-to-dc converter illustrating the method used to achieve low junction temperatures in the power semiconductors. The transistors are mounted on a massive extruded aluminum heat sink which projects 1-1/4 inches beyond the front panel to insure adequate air flow for convection cooling, thereby minimizing the temperature at the semiconductor junctions and within the unit. The angle framework construction is used to provide easy access to any of the five sides for assembly and maintenance. The open type of construction also provides for good air flow through the unit, thus contributing to the low equipment temperature. This method

of construction, coupled with the use of heat sinks, eliminates the need of a fan or blower for forced convection thus assuring additional reliability of the Power Supply.

IV. FEATURES AND ALARMS

4.1 *Special Electrical Features*

The D2 Bank transmission equipment requires special electrical features from the PSCC. The minus-32-volt output must track the plus-32-volt output very closely. A difference in the absolute value of these two voltages causes a shift in the voltage reference used for the coder and the decoder modules, and even a small reference voltage change can produce an error in the pcm code. Interlocking is required for the plus- and minus-24-volt outputs to provide simultaneous turn-on of the two outputs, and to remove the plus 24 volts when the minus-24-volt output falls below minus 18 volts in order to prevent destructive overheating of circuits biased by the plus and minus 24 volts.

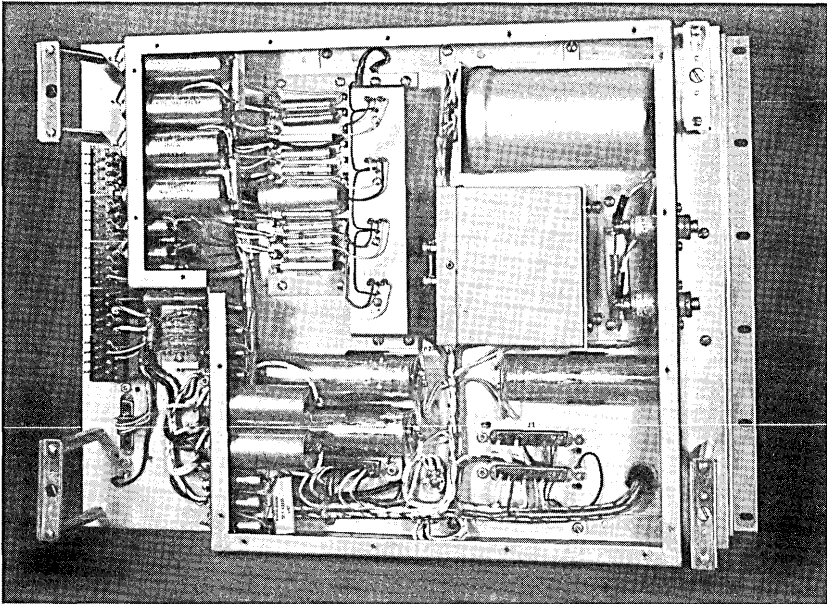


Fig. 8—D2 Channel Bank—J87327A power supply prototype (chassis), rear view, shield cover removed.

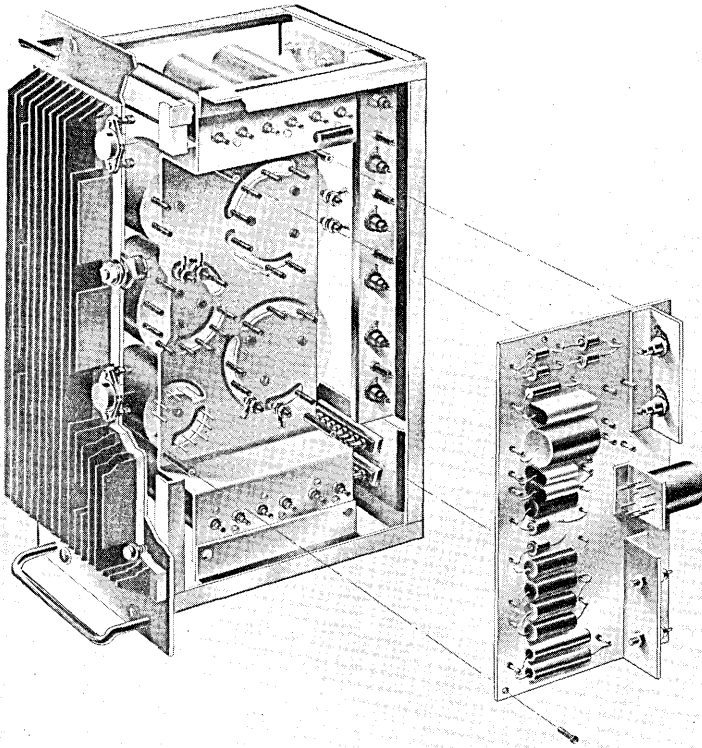


Fig. 9—Exploded view of dc-to-dc converter.

4.2 Alarm System

The D2 Bank bay must have all seven of the PSCC output voltages available in order to operate. Therefore, an extensive alarm system was designed so that, in the event of a PSCC failure, the affected system would be removed from service, the operating personnel alerted, and the location of the failure flagged to minimize down time.

To provide the alarm features, there is a circuit breaker in series with each of the outputs, and also one in series with the input. The circuit breakers will self-trip on input or output overloads. The circuit breakers associated with the three highly-regulated outputs (+5V, +32V, -32V) will also be tripped by their respective high-voltage shutdown circuits. The 24V circuit breaker may also be tripped by the low-voltage shutdown circuit on the minus-24-volt output.

Since all of the outputs would be affected by a high voltage at the output of the switching regulator, the INPUT circuit breaker is tripped by the minus-24-volt high-voltage shutdown circuit in the event of a switching regulator failure.

Whenever any of the PSCC circuit breakers trip, or when the output of the switching regulator is lost, a front panel alarm light is lit and an internal relay operates. One pair of contacts on the relay are used to busy the affected trunks out of service and another pair of contacts on the relay are used to light a bay alarm light and to activate an office alarm bell.

V. SUMMARY

The power-systems converter circuit described here features many design innovations. High-power switching regulator and converter operation above 20 kHz, resulting in silent operation with low power dissipation. A switching regulator used as the sole regulator for four outputs and as a preregulator for three other outputs. A plus and minus 32-volt regulator which provides precision tracking of the plus-32-volt output by the minus-32-volt output. A five-volt series regulator design producing low dissipation in the series transistors, and very reliable operation. The use of plug-in modules for ease of field maintenance, trouble shooting and repair. These features introduced several challenging problems in the area of electromagnetic noise suppression and mechanical design.

VI. ACKNOWLEDGMENTS

The authors wish to thank the many people who have contributed to the design and the development of this converter circuit. Particular credit is due to P. W. Ussery and F. C. LaPorta for their work in developing the power supply and F. T. Dickens for his design of the magnetic apparatus.

The circuit design was done under the supervision of G. W. Meszaros, the mechanical design was done under the supervision of S. Mottel and the magnet apparatus design was done under the supervision of T. G. Blanchard.

REFERENCES

1. Bomberger, D. C., Feldman, D., Trucksess, D. E., Brolin, S. J., and Ussery, P. W., "The Spacecraft Power Supply System," *B.S.T.J.*, 42, No. 4 (July 1963), pp. 943-972.
2. Schaeffer, J., and Lupoli, P. J., "Transistor Switching Regulator," *Electrical Design News*, April, 1965.

D2 Channel Bank:

Physical Design and Introductory Program

By D. J. VAN SLOOTEN and K. A. GLUCKOW

(Manuscript received June 22, 1972)

The D2 Channel Bank is designed to provide simplified engineering, installation, and maintenance. Integral voice-frequency alarm and access are provided in a packaged shop-wired frame with a centralized built-in test capability to facilitate initial line-up, testing, and trouble shooting.

Circuits are implemented with discrete components and with thin-film and silicon-integrated circuits. Low cost, reliable assembly and wiring techniques are employed. The frame organization and circuit partitioning provide a functional arrangement of circuits with good electrical isolation between critical multiplexing and coding functions.

An introductory program and an on-going reliability program have demonstrated the adequacy of both equipment and documentation. In the first 21 months of operation, approximately 3 percent of the circuit packs shipped have failed in initial line-up or in-service. This compares favorably with the performance of similar systems, and recent design modifications are expected to result in substantial improvements.

I. INTRODUCTION

The D2 Channel Bank multiplexes and codes the telephone traffic carried by 96 two-way toll-grade trunks into 4 two-way 1.544-megabit signals for transmission. Companion articles in this series discuss system aspects of the D2 Channel Bank, and the design of the circuits which embody it. This paper discusses the physical design of the D2 Channel Bank. This includes activities that have circuit information as input and, as output, the specification of a manufacturable design to Western Electric. The D2 introductory program and an on-going reliability program intended to ensure satisfactory service to the Operating Companies are also discussed.

II. PHYSICAL DESIGN OBJECTIVES

The primary goal of the physical design of the channel bank was

a manufacturable design which would satisfy the performance and cost objectives for toll service. An analysis of the cost of the D1 Channel Bank equipment as installed for toll-connecting service indicated that a considerable portion of the installed cost was attributable to engineering, installation, and maintenance. Objectives for the D2 Channel Bank, therefore, included simplified engineering, installation, and maintenance, as well as low manufacturing cost.

Operating Company engineering and installation would both be simplified if the channel bank could be furnished in a "packaged" frame containing all of the voice-frequency equipment associated with the channel bank. This contributes to simplification of installation because the only wiring required would be that which brings the trunks to the channel bank and takes the digital signals to the transmission facility and back. Additional simplification of installation would result if the wiring could be identical from the Intermediate Distributing Frame to the channel bank for all types of trunks, and if the wiring could terminate in the same way regardless of the particular option used in the channel banks.

In order to simplify maintenance, a number of objectives were established. First, all apparatus was to be pluggable to minimize outage time and circuit packs were to be compatible on an individual basis, so that circuit packs need not be replaced in sets. Centralized controls and indicators were to be provided for simplified line-up and fault diagnosis. Further simplification was planned by eliminating the need for setting carrier-group alarm options and the need to store voice-frequency attenuator pads.

III. SYSTEM CONFIGURATION AND FEATURES

The physical design of the D2 Channel Bank was customized to satisfy the system and electrical requirements of a toll digital channel bank. An early decision to provide 96-channel coding (four digroups of 24 channels each) in a single coder required modular physical arrangements based on multiples of four digroups to be used in 7-foot, 9-foot, and 11-foot 6-inch (2.1, 2.7, and 3.5 meters) frame heights. Figure 1 shows the layout for the 11-foot 6-inch frame which houses four digroups and Fig. 2 shows the 7-foot version consisting of a triple bay frame housing eight digroups. An expandable 9-foot version houses 12 digroups in a four-bay arrangement. A D2 Channel Bank standby bay is also available in 11-foot 6-inch, 9-foot, and 7-foot frame heights. It is designed to provide protection switching for the D2 Channel Bank

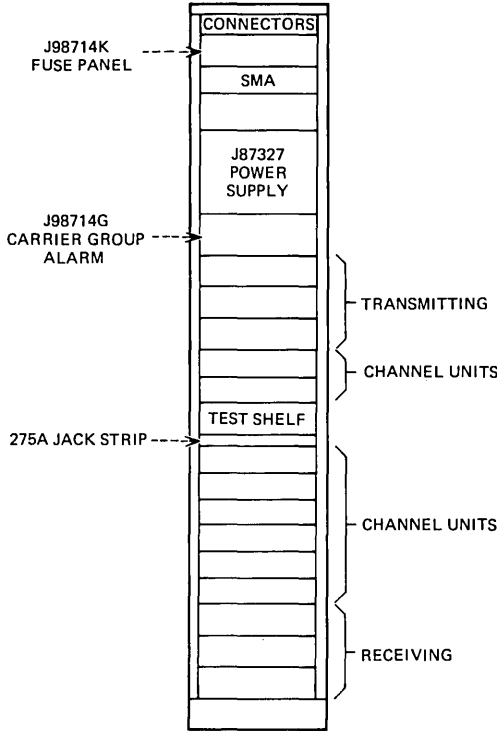


Fig. 1—D2 Channel Bank, 11'-6" bay.

on a digroup basis. Forty digroups may be protected by one D2 standby bay.

The interface between the voice-frequency plant and the common equipment of the channel bank is provided by the channel unit. The channel units mount in the center of the 11-foot 6-inch frame separating the transmitting and receiving common equipment. The 7-foot and 9-foot versions preserve the same arrangement of channel units and common equipment so that factory wiring is identical for all frame codes. Eight VF leads are brought to each channel unit from the intermediate distributing frame (IDF) which permits the various channel unit codes to be intermixed with identical wiring.

In addition to matching the particular types of trunk circuit to the channel bank, the channel unit provides standard signal level (transmission level of -16 dB, transmit, +7 dB, receive), and signaling jack access for maintenance as well as patching. Inter-bay patch jacks are also provided to facilitate restoration patching.

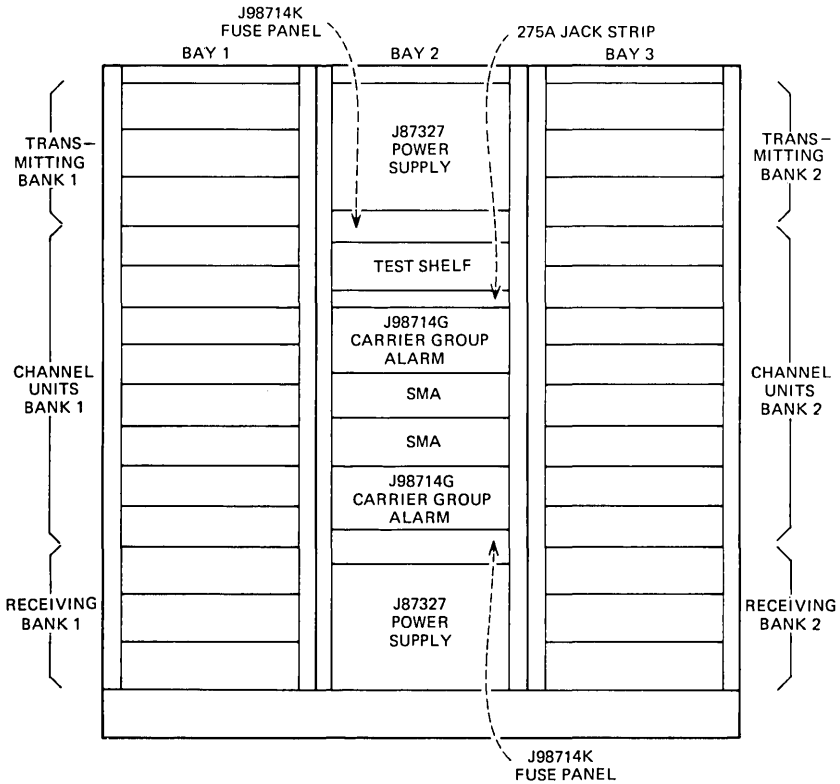


Fig. 2—D2 Channel Bank, 7' triple bay.

On certain channel unit codes, additional access is provided to allow for connection to external trunk equipment such as echo suppressors or delay equalizers, and for additional access to remote maintenance or to automatic protection switching.

Additional voice-frequency equipment included in the factory wired frame, as shown in Figs. 1 and 2, consists of the carrier group alarm (CGA) and the switched-maintenance access system (SMAS) connector. In addition to standard office alarm interfaces, an E2 status and control interface* is provided at the test shelf. Energy for the de-to-dc power converters in each frame are supplied from -48V office battery via a fuse panel.

* This allows remote monitoring and control of the channel bank at a distant centralized location.

The test shelf, mounted at a convenient height for craftsman operation, houses several circuit packs which form a built-in test facility. The test shelf is shown in Fig. 3. The facility includes a filter for measuring crosstalk resulting from a 1000-Hz test tone when applied at the 4-wire input of a channel bank. It is used in conjunction with access provided on the test shelf to the centralized transmission and noise-measuring system and to the milliwatt supply. A digital signal

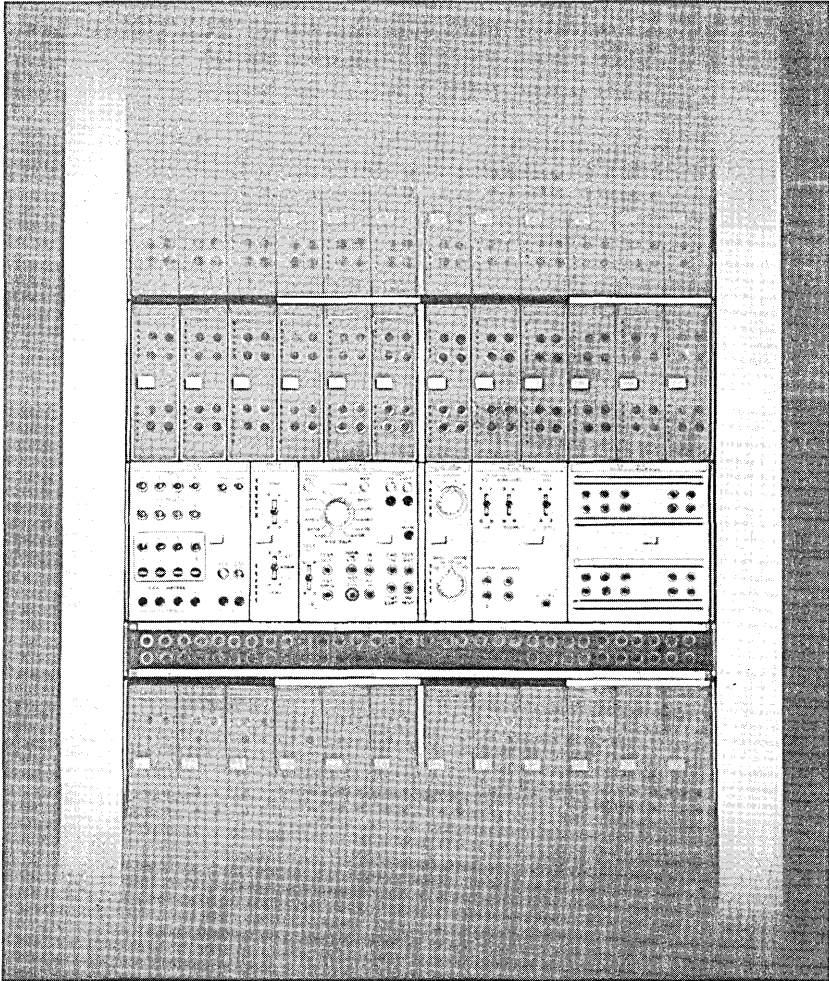


Fig. 3—Test shelf.

generator provides the digital equivalent of a 0-dBm0 level test tone as required for receiving gain adjustment of each voice channel. The talking test circuit is used to monitor transmission and noise on a high-impedance bridging basis and permits talking on the facility for maintenance and trouble shooting. In addition, alarm control circuitry on the shelf indicates various failure conditions, and provision is made for looping the digital signal for fault isolation. A miscellaneous jack unit is provided with jack, key, and lamp positions for order wires and jack-ended call numbers. The intent is that these features will be engineered by the operating companies to suit their individual needs.

In addition to the built-in test facility and the centralized transmission and noise measuring system (or their portable equivalent), a volt-ohm milliammeter and a portable trouble-locating test set are required for trouble-locating tests on the common equipment.

IV. PHYSICAL DESIGN CONSIDERATIONS

4.1 *Circuit Partitioning and Organization*

The constraints of circuit performance, operation, reliability, maintenance, and manufacturing cost are significant factors in circuit partitioning.

In the channel bank, functional partitioning was a primary consideration in order to provide adequate performance of critical circuit functions and to provide simplified factory testing and maintenance. A case in point was the decision to design the coder and decoder as single-circuit packs in spite of the large size of these functions.

In view of the multi-stage multiplexing scheme and 96-channel coding, considerable thought was given to minimizing the number of working trunks placed out of service by a circuit-pack maintenance removal. A notable exception was the packaging of eight per-channel gate and filters on a single circuit pack rather than placing the gate and filter on its channel unit. The need to place the gate and filter in close proximity to the multiplex and coding functions in order to minimize the length of the critical PAM bus was an overriding consideration.

An analysis of circuit function versus packaging volume indicated that a 5- by 10-inch (12.5×25 cm) printed wiring board was optimum for channel units and that a 6- by 10-inch (15×25 cm) printed wiring board was an appropriate choice for common units. This also provided a good pin match to the "908" series Western Electric 40-pin connector. Large functions could be accommodated within this scheme by the use of module boards and a second connector.

The arrangement of channel units and common-unit circuit packs is shown in Fig. 4. This arrangement was chosen to provide isolation of transmitting and receiving common equipment, and to provide for

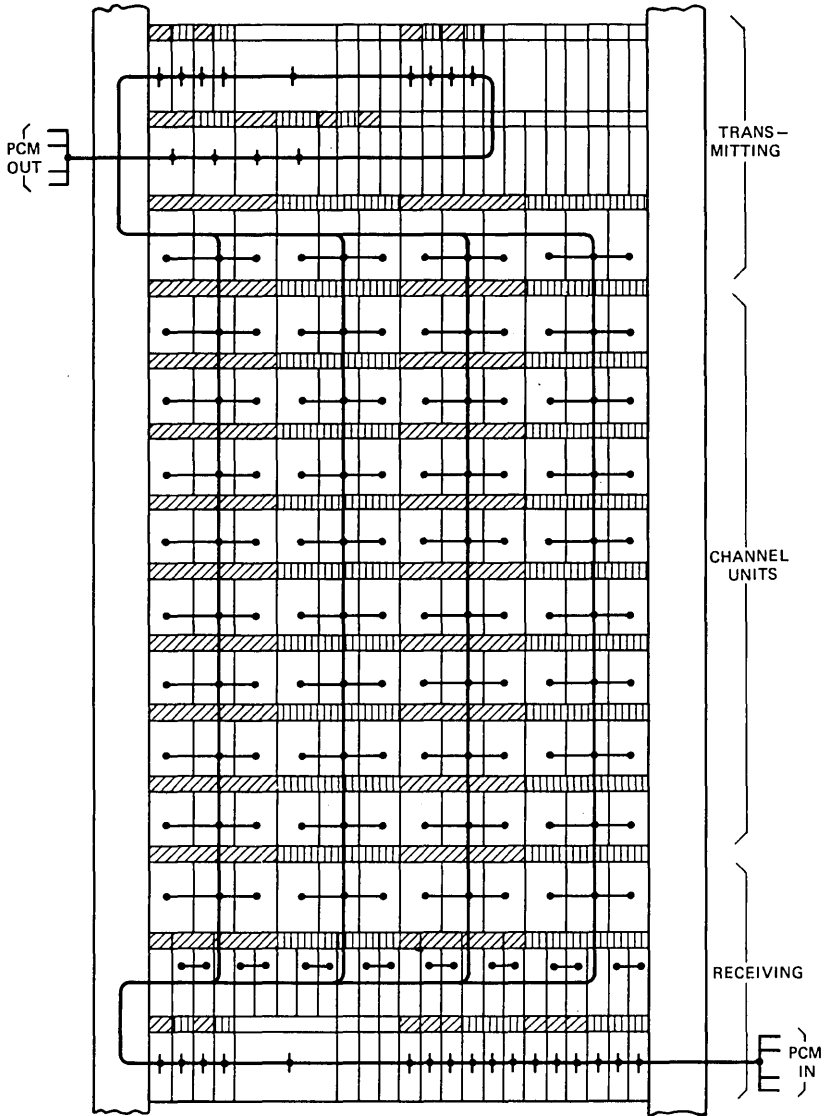


Fig. 4—Signal flow.

an orderly signal flow from the PCM input through the decoding and demultiplexing common functions to the channel units. The opposite direction of transmission is similar in physical arrangement. This gives good isolation of low-level analog signals and high-level digital signals and minimizes critical lead lengths.

The channel units for each of the four digroups are arranged in three adjacent vertical columns of eight units each. Each column of channel units is associated with the transmitting gate and filter circuit pack directly above it and the receiving gate and filter circuit pack directly below it. This facilitates level adjustment in the gate and filter circuit packs.

4.2 *Framework and Circuit Pack Mounting*

The D2 Channel Bank employs the transmission standard 23-inch (58 cm) unequal-flange cable-duct frame. The circuit packs are flush mounted with the wide flange of the frame for good appearance and to provide a full duct for cabling.

A flexible multiheight circuit pack mounting or shelf assembly was designed to accept the diversity of channel bank functions. The shelf is shown in Fig. 5 which is made from three identical die-cast aluminum parts that form the top, bottom, and rear of the assembly. The casting provides circuit pack guides and mounting slots for connectors with a 50 percent open area for ventilation. The use of a single part guarantees accurate alignment of card guides to connectors. The rear piece is trimmed to provide a shelf height from 5 to 10 inches in 1-inch increments. The side plates are fabricated with an integral

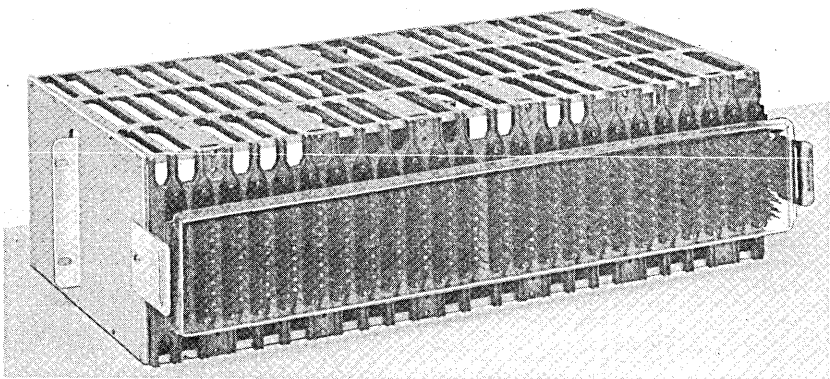


Fig. 5—Flexible multiheight circuit pack mounting or shelf assembly.

mounting bracket for assembly to the frame. The design provides a minimum circuit pack module width of seven-eighths of an inch with a maximum capacity of 24 circuit packs per shelf. Circuit packs of up to five module widths are used. Additional features include a notched locking bar for retaining circuit packs and for digroup designation. A rear plastic cover provides protection to the terminals.

The heat dissipation for the 11-foot 6-inch version of the channel bank is 720 watts. Thermal analysis and temperature measurement indicated satisfactory operation and reliability in the required range of office ambients from 2°C to 49°C. In some cases, relocation of components was necessary to eliminate hot spots in order to achieve the desired thin-film resistor aging characteristics.

The frame features a low-impedance power feed and ground system. The power supply voltages are distributed to the circuit packs by means of laminated distribution bars which have up to seven layers. The shelf assemblies for the transmitting and receiving sections of the frame contain an insulated ground plane at the rear within the connector field. Insulated ground straps separately connect the ground planes to the power supply to provide a radial single-point ground system within the frame.

Wiring volume within the frame is considerable as a result of the voice frequency interface requirements. For example, connections to the SMAS connector, to the carrier group alarm, to external trunk equipment require twelve, eleven, and eight leads per channel, respectively. The provision of protection switching necessitates an additional eight leads per channel. The voice-frequency input requires 768 (96 by 8) 24-gauge wires which are cabled from the IDF directly to the channel units. In retrospect, the elimination of a terminal strip at the top of the frame proved to be a poor choice since damage has occurred during installer wiring at the channel units. Future bays will have all voice frequency leads connectorized, which will correct the problem and further simplify installation.

4.3 *Circuit Pack Design*

The D2 circuit packs consist of an epoxy-glass printed wiring board, and a plastic faceplate and handle. Discrete components, thin-film and silicon integrated circuits mount through holes in the boards for connection by mass soldering. Printed wiring board modules mount on the master printed wiring board to realize large functions that could not be otherwise accommodated by a single planar board.

Epoxy glass was chosen as a board material because of its mechanical

strength and stability. Heavy apparatus could be accommodated without requiring the use of metal supporting frames. Epoxy glass permitted the use of plated-through via holes for codes where minimum circuit area and path length was desirable. Gold fingers are provided for contact to the "908"-type connector. The faceplates are constructed of fire-retardant PVC, and provide cavities for test points and a handle to facilitate removal. The circuit-pack code and title are placed on the faceplate for easy identification. In addition, color-coded handle labels are provided for identification of special service trunks.

Figure 6 illustrates a typical channel unit. All channel units are two modules (1-3/4 inches) wide. Jacks mounted on the faceplates provide standard level (-16 dB, +7 dB) access for maintenance and patching. Dark grey faceplates are used on channel units to distinguish them from common circuit packs. Thin-film loss-adjusting pads are provided for adjusting office losses over a range of 0 to 16.5 dB in 0.1 dB steps. Options on certain types of channel units are controlled by screws which are screwed down to insert the option and unscrewed to remove the option. The options provide for loop resistance compensation, network built out capacitance, and for matching trunk characteristics to the carrier group alarm. On certain codes of channel units, a second connector is provided to allow for access to external trunk equipment such as echo suppressors or delay equalizers, and for additional access to remote maintenance or to automatic protection switching.

Figure 7 shows a typical common-unit circuit pack. In Fig. 7, note

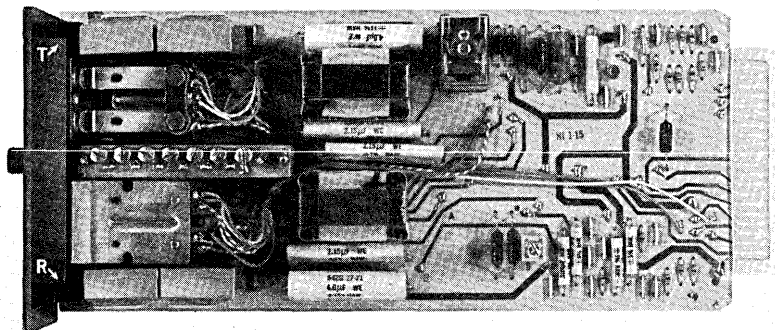


Fig. 6—Typical channel unit.

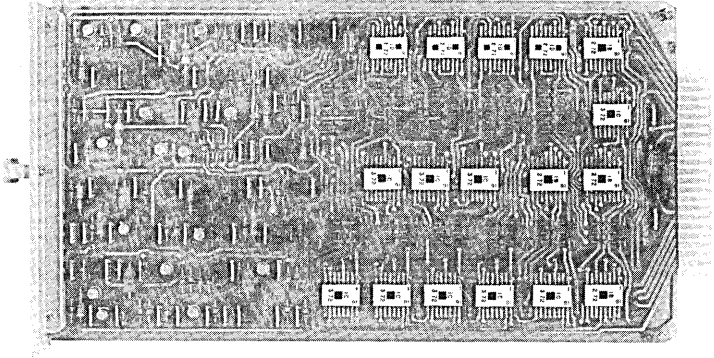


Fig. 7—Typical common-unit circuit pack.

the alignment of components on a single axis to facilitate machine insertion. Figure 8 shows the most complex circuit pack in the terminal, namely, the coder, which contains over 840 components. The module boards employ connectors to provide interconnection between boards which facilitate assembly and repair.

Table I summarizes the type and number of printed wiring boards used in the channel bank. Note that the packaging density is 5.5 com-

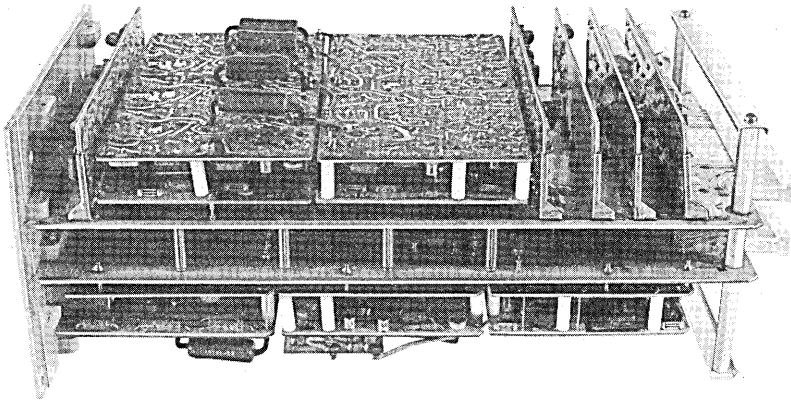


Fig. 8—Coder—the most complex circuit pack in the terminal.

TABLE I—PRINTED CIRCUIT BOARDS

Master boards.....	207
Module boards.....	496
	703
Number of different boards.....	90
Maximum components on single-space circuit pack with modules.....	388
Maximum components on single-space circuit pack without modules.....	292
Maximum components on any circuit pack.....	890 (coder)
53 square inches usable board area	
5.5 components per square inch	

ponents per square inch (one component per square centimeter).

Table II summarizes the type of circuit components used in the channel bank. The terminal is realized primarily with discrete components, although 19 codes of thin-film precision networks are used for multiplexing and coding, and one family of silicon-integrated logic circuit is used for digital processing. Figure 9 shows an assembly containing thin-film resistor networks. The resistor end-of-life tolerance is 0.04 percent absolute, and 0.02 percent in resistance ratio.

Selection criteria for components included minimization of code

TABLE II—COMPONENTS

Quantities shown are for fully-equipped bay with 96 dial-pulse-orig channel units, but does not include power supply components	
Resistors.....	9380
Diodes.....	5739
Capacitors.....	4833
Transistors.....	2006
Varistors.....	735
Transformers.....	574
Attenuators.....	384
Jacks.....	309
Potentiometers.....	213
Connectors.....	207
Filters.....	193
Relays.....	106
Inductors.....	130
Integrated circuits.....	68
Networks.....	51
Special thin-film resistors.....	23
Miscellaneous switches, keys, etc.....	55
Total.....	25,946

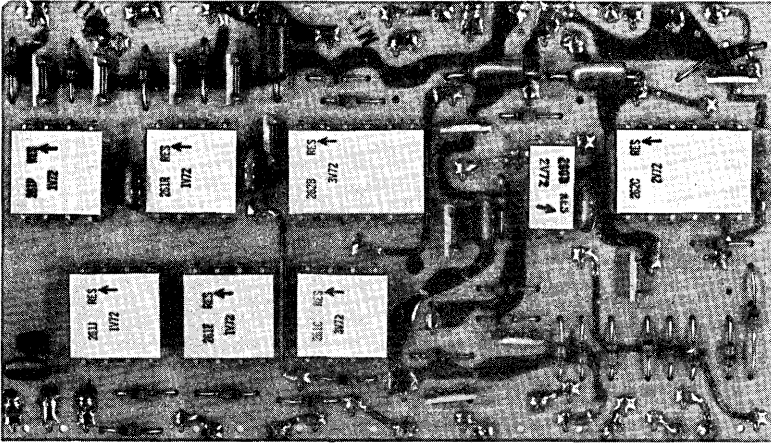


Fig. 9—An assembly containing thin-film resistor networks.

types consistent with performance and cost objectives. Good reliability was a key concern, and components were derated to assure satisfactory life. The predicted total channel-bank failure rate is 114,000 FU (Failure Units) corresponding to a mean time to failure of 8760 hours. This compares favorably with similar equipment. (The actual failure rate during the life of the D2 Performance Study was 6300 hours mean time, or 1.4 the predicted rate. This represents the first two years of service.)

V. INTRODUCTORY PROGRAM

5.1 *AT&T Introductory Program Planning*

As enunciated by the American Company, the Introductory Planning Program for New Transmission Systems involves four activities:

- (i) Organization of AT&T, Bell Laboratories, and Western Electric managing teams to coordinate the introduction of each new product.
- (ii) Development and continuous monitoring of schedules for all necessary activities needed to assure smooth introduction of the new product.
- (iii) Establishment of a close working relationship with the Operating Companies in the early stages of development to assure that field needs are met.

- (iv) Organization of follow-up programs for initial production to verify proposed operation of field conditions and adequacies of written Bell System Practices (BSP's).

The fourth point applied to the activity of the Laboratories performed during the initial manufacture and during the initial installation.

5.2 *Initial Installation*

Following a number of AT&T, Bell Laboratories, and Western Electric tri-company meetings, the Western Electric Quality Service Management (QSM) Organization chose a Los Angeles central office of the Pacific Telephone and Telegraph Co. for the New Product Survey of the D2 Channel Bank. As the survey is organized, participation of personnel from all three companies is required as a means of implementing and evaluating the New Product Survey.

The QSM Organization provided observers at the chosen site to monitor the adequacy of packaging, bay and circuit pack installation, including related documentation (Western Electric drawings, installation handbooks, and Bell Laboratories BSP's).

An expedited repair procedure to enable Bell Laboratories and Western Electric personnel to quickly repair and evaluate circuit pack failures was implemented with the cooperation of the Western Electric Merrimack Valley Works merchandising organization.

A summary of problems found during the New Product Survey was submitted to members of the team within four months after the survey began. The group reviewed these problems, and referred them to the responsible organizations for corrective action. The New Product Survey served its purpose most effectively in the identification of problems in the areas of packaging and bay wiring.

VI. RELIABILITY PROGRAM

6.1 *Objectives*

In the early part of 1964, it became increasingly clear to the Bell Laboratories and Western Electric personnel involved in the repair of D1-T1 circuit packs that a thorough study was required to examine the reliability of these plug-ins. Thus began the first effort within the transmission area to record a history of failures for circuit packs of a given transmission system.

Beginning in 1964, the D1-T1 Repair Study was able to clearly identify problem areas associated with individual circuit pack codes, and served as a means of providing the development organization with clearly defined current engineering experiences through these reports.

In 1967 it became evident that a more detailed accountability of D1-T1 failures was necessary to evaluate improvements made in the various circuit packs, design changes incorporated into the D1 bays, as well as trouble-shooting procedures used by Operating Company craftsmen. A further problem uncovered was the excessively large number of plug-ins returned for repair with no-trouble-found (NTF). As a means of providing more valid data concerning the overall reliability of D1-T1, a study was begun in the Operating Company central office to evaluate initial lineup and in-service performance of this system.

These two studies were the forerunners for the D2 Repair Study and the D2 Performance Study. The objective of these studies is to determine:

- (i) the nature and scope of D2 circuit packs troubles
- (ii) the adequacy of BSP's and craftsmen usage
- (iii) the nature of channel bank outages
- (iv) the cause of no-trouble-found returns
- (v) the results of circuit pack improvements
- (vi) the requirements for telephone company circuit pack spare inventory.

6.2 *D2 Repair Study*

The Repair Study is intended to continue throughout the manufacturing interval. Data for this program are collected by Western Electric personnel at the repair locations which are presently the Merrimack Valley Works and the Los Angeles Service Center. The Repair Study examines and analyzes information on all D2 apparatus returned for repair.

During the first year and a half of repair, the information received from Western Electric repair locations was manually analyzed by physical design personnel. Quarterly summaries of the results were provided for general information to outside organizations, while constant feedback was provided to development organization personnel for action as required.

During the second quarter of 1972, the outputs of the incoming data were summarized and analyzed through a computer program written and maintained by a programming support group.

6.3 *Performance Study*

The purpose of the D2 Performance Study, as stated earlier, is to provide a means of evaluating the nature and scope of D2 troubles

beyond the simple analysis of individual circuit pack failures. Data for this program are provided by central office craftsmen at three locations. These are the Los Angeles and San Diego regions of Pacific Telephone, the Chicago area of Illinois Bell Telephone, and the Dallas area of Southwestern Bell. The basic objective is to obtain information on craftsmen-BSP-equipment interactions and thus determine where there are defects or weaknesses in the BSP or the equipment. When difficulties are experienced by craftsmen in either equipment, line-up routines, or in connection with in-service failures, they fill out a form reporting the nature of the failure. If the trouble is one involving the return of apparatus for repair, the craftsman identifies the circuit pack with a sticker serialized to agree with the number on the form he is returning. In this way, it is possible to correlate the craftsman's report with repair information received from the repair center. In this connection, the incidence of no-trouble-found is of particular interest. The Performance Study, as originally planned, was completed in December, 1971, after a life of 1-3/4 years.

6.4 The Results of the D2 Repair and Performance Studies

The D2 Performance Study and its interaction with the D2 Repair Study has provided a means of analyzing D2 outages much more rapidly than for other transmission systems in the past.

Approximately 1000 digroups of D2 apparatus were installed in the three performance-study locations during the life of the study. This quantity of installed digroups represents a significant percentage of D2 digroups installed through the entire country during this period.

During the life of the study, 782 circuit packs were replaced during initial line-up representing a 2.8-percent replacement of all circuit packs shipped. One hundred seventy-five circuit packs were replaced on an in-service basis, and reflected an actual replacement rate of 1.4 times that which was predicted on the basis of component FU rates for a fully equipped N2 Channel Bank. The actual equipage of D2 Channel Banks during the life of this study was 2.8 digroups per bank. Although some individual circuit pack codes reflect high replacement rates, the overall replacement rate was shown to be far better than previous transmission systems had experienced during such early stages of production. The actual initial lineup replacement rates of common plug-ins have shown significant reductions which reflect craftsmen experience and product improvement. This drop-off with time is shown in Fig. 10.

A further area of concern has been the high initial-lineup replace-

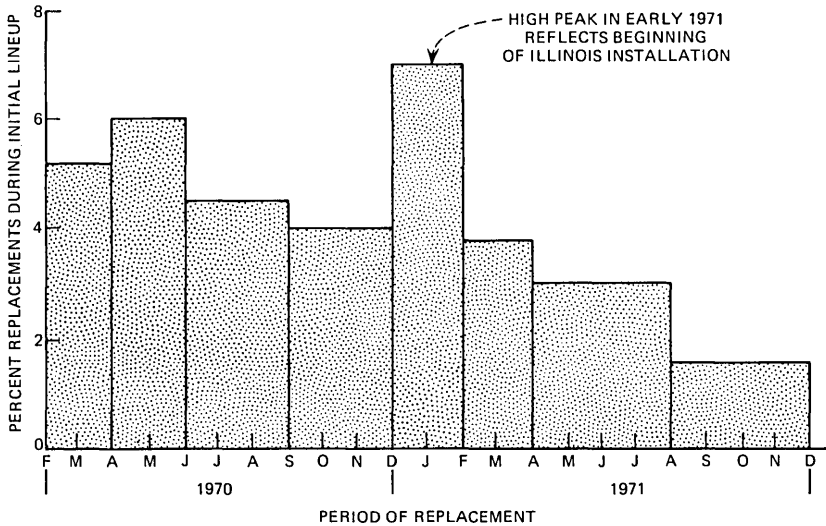


Fig. 10—Replacement rates for D2 circuit packs.

ment rate of D2 power supplies, and this continues to present a problem at the close of the Performance Study. In-service replacement rates for the power supply plug-ins appear reasonable at 1.5 the predicted rate.

The Performance Study is the quickest means of uncovering the broad patterns of failure which develop during initial production. The information received through the D2 Repair Study is, however, significant in that it reflects all plug-ins returned for repair from the Operating Companies. It further provides us with a means of examining changes incorporated into existing products by relating failure to dates of manufacture and series numbers appearing on the faceplate of each plug-in.

One means of evaluating craftsmen's performance is to examine the relationship between plug-ins removed and no-trouble-found with the equipment when it is repaired at the Service Center. In the case of D2, the no-trouble-found (NTF) rate for repaired plug-ins is approximately one-half that of its predecessor. A further examination of NTF removals shows that, in the case of 86 multiple plug-ins removed (that is, for a single failure, more than one plug-in is removed), only four cases were reported of no-trouble-found in any of the removals. This supports other evidence of the adequacies of the D2 BSP routines

and of the ease in which craftsmen were able to go through line-up and trouble-shooting procedures with the equipment.

The overall replacement rate of D2 plug-ins received from all Operating Companies is 2.8 percent. When no-trouble-found rates are subtracted from replacements or all plug-ins returned for repair, the failure rate for returned plug-ins drops to 2.4 percent.

As noted, the failure rate of D2 common circuit pack in-service failures is presently 1.4 times higher than predicted. With the improvements being incorporated into the various plug-ins which have shown excessive replacement and failure rates, it is expected that this ratio will drop to approximately 1.1 by the end of 1972.

Further work is presently under way to establish relationships between failure rates and circuit outages. A standard spare ratio based on reliability considerations is being compiled which will permit the Operating Company to reduce their present inventory of D2 spares which currently amounts to approximately 10 percent of the total installed circuit packs. It is expected that this number will drop to between 4 and 6 percent and will reflect a substantial cost reduction to the Operating Companies.

D2 Channel Bank:

Manufacturing and Testing

By J. E. D. BATSON, JR., and J. W. GORMAN

(Manuscript received June 28, 1972)

Previous articles in this issue have covered the design concepts of the D2 Channel Bank. This article is about the contribution by Western Electric, the manufacturing and supply organization for the Bell System. The discussions will focus on the effort at Western Electric, Merrimack Valley Works which includes bay and circuit pack manufacturing. Although certain portions of these topics are rather routine to those in manufacturing, they will be briefly mentioned in order to give completeness in explaining how a new product line is put into production. This article thus serves to give the reader a better understanding of how a product is introduced by telling, in narrative form, the D2 Channel Bank story.

I. INTRODUCTION

Manufacturing of the D2 Channel Bank takes place at Merrimack Valley and Kearny Works of the Western Electric Company. The introduction of any new system by Western Electric involves considerable effort by many locations within the system. This article will concentrate on those activities at Merrimack Valley which are directly associated with the manufacture of circuit packs and bays for the D2 Channel Bank.

Manufacturing development for the D2 Channel Bank began in 1965, and ended with the first shipments late in 1969. Major activities during this period were (i) early development (ii) construction of field trial units (iii) production planning (iv) test planning and (v) initial production. The project is presently on a continuing production basis.

1.1 *Early Development*

For those concerned with development to manufacture the D2 Channel Bank, the first stage was to become familiar with the details of the product and to make preliminary production plans. One of the first steps in this operation was to establish an estimated cost. This

was a difficult task since the design work was not complete at the time. However, the known factors were calculated and the rest estimated using judgment gained from experience. Next, a schedule was established for completion of the design and development work in order to determine an availability date. This date and the pricing information were used with technical information from Bell Laboratories to transmit Engineering Letters to telephone companies that described the system. As development progressed, the data was updated to take into account new information.

Product engineers began studying the system for areas of new technology. One example was the precision thin-film resistors used in the coder and decoder. Steps were taken to ensure that these resistors could be manufactured in time for initial production. Models were made and submitted to the Laboratory for approval in this area of new technology.

At the same time, test engineers began studying D2 to learn how to test it efficiently and at a reasonable cost. One engineer from Merrimack Valley was assigned to Bell Laboratories, Holmdel to work directly with the design engineers during their final design stages. His specific duties were:

- (i) Write preliminary test specifications.
- (ii) Estimate test set requirements.
- (iii) Coordinate test requirements with test set capabilities.
- (iv) Learn the details of the D2 system.
- (v) Establish personal contacts between the two engineering groups.

1.2 *Field Trial*

As the design of the D2 Channel Bank approached completion, a field trial was planned. Several bays and their circuit packs were constructed and installed in the Philadelphia area for evaluation. All equipment for the field trial was manufactured in the model shop at Merrimack Valley. This was the first direct contact with the actual product for most engineers in the group.

During field trial production, the product engineers had their first real opportunity to learn what kind of problems to expect in manufacturing this new channel bank. These problems were concerned with assembly techniques, module board handling, integrated circuit mounting and many others.

The magnitude of this project also became apparent during this period, as the nearly 160 different codes of circuit packs or module boards were assembled.

Test engineers used this time to good advantage by setting up "mucket" facilities and actually testing units. If available, preliminary test specifications were used to prove their adequacy. This procedure established a small group of test engineers who began to have a detailed understanding of the intricacies of the D2 Channel Bank's operation. Time spent during this period was repaid many times over during the difficult days of early production.

Another beneficial aspect of this period was the beginning of close relationships between individual Bell Laboratories and Western Electric engineers. The two parties discussed the problems as they came up, and arrived at solutions. This small operation repeated many times over a few months became a good foundation for mutual understanding. The authors feel strongly that these relationships were one prime factor in the successful introduction of this product.

II. PRODUCTION PLANNING

With the field trial complete, preparations were made to begin manufacturing the circuit packs and bays of the D2 Channel Bank. When a new product is introduced, many different organizations play an active role. In keeping with the scope of this article, only those engineering activities directly associated with the product will be discussed. The reader should be aware that this is the tip of the iceberg representing only a fraction of all the necessary activities.

2.1 *Drawings*

Before any product can be built, appropriate drawings must be made and distributed. This was one of the major tasks due to the many codes of module boards and circuit packs.

Initially, Laboratories' Design Information (LDI) is transmitted to Western Electric for a preliminary analysis of each code. Both Product and Test Engineers at Western Electric review the LDI, which usually consists of a schematic, stocklist, and a proposed layout, so that any particular assembly problem or test problem which might arise can be anticipated. Previous agreements between Bell Laboratories and Western engineers had established ground rules for the proposed layouts that included considerations for ease of manufacturing. This was done to reduce the huge load on the Western drafting organization.

Next, the Western drafting organization began preparation of the official manufacturing drawings. The draftsmen familiarized themselves with the schematic and proposed layout.

Variations from the proposed layout were initiated by the draftsman

or engineer when alternates were seen that would substantially reduce the cost. An example of these changes was the elimination of a module board from certain codes of channel units when a method was found to incorporate the components on the main board. These proposed design variations were checked with Bell Laboratories, particularly when critical circuits such as the coder or decoder were involved.

The draftsman's layout was checked for accuracy by another member of the drafting organization, and the art master was made, usually by the original draftsman.

Prints of the formal drawings were sent to the product engineer for final analysis and approval. By working closely with the draftsman on the physical layout of the circuit, the product engineer's main task was to check the artwork against the schematic. With this accomplished, signatures were affixed making the drawings official.

2.2 *Production Facilities*

In planning for production facilities, the following phases of manufacturing were given attention: (i) initial low-volume production when the processes were unfamiliar to the shop personnel. (ii) intermediate volume production when the processes were familiar to some personnel and yet unfamiliar to the new personnel required for the increased volume, and (iii) normal expected future volume which will have a normal amount of experienced shop personnel with mechanized and automated production machinery to supplement the hand assembly requirements.

2.2.1 *Manufacturing Layouts*

Manufacturing layouts are the detailed instructions for the assembly of the various circuit packs. The initial set was written using the intermediate volume of production as a guide. This was done because a shop goes through the low-volume stage rather quickly. It would have been a waste of time to put too much effort on methods that would have only a short life.

Industrial engineers played an important role in writing layouts. They supplied base hour rate information for each assembly operation in the layout. Where possible, standard rates were used for normal operations such as inserting components, soldering, etc. Any new processes were estimated and later refined with actual time-motion studies. Accurate rate information is essential to a successful product. If the rates are too loose, the price of the product will be adversely affected; if too tight, the shop personnel will be frustrated in their efforts to meet these rates

and could cause a severe morale problem. Industrial engineers also contribute by designing the individual work positions and visual aids.

2.2.2 *Manual Assembly*

Good manual assembly bench layouts are essential to efficient manufacturing in early low-volume production. They are equally valuable on larger volume production in the area of impractical automatic assembly and low-volume miscellaneous codes.

The typical D2 manual assembly position consists of a five-foot-wide standard bench position and chair. On this bench position is a semicircular rack which holds from 16 to 28 individual removable stacking bins (Fig. 1). Precut and formed components are held in these numbered bins and are available for insertion into a printed wiring board as called for on the associated visual aid. If more than 28 different components are required for this position, an additional rack can be used.

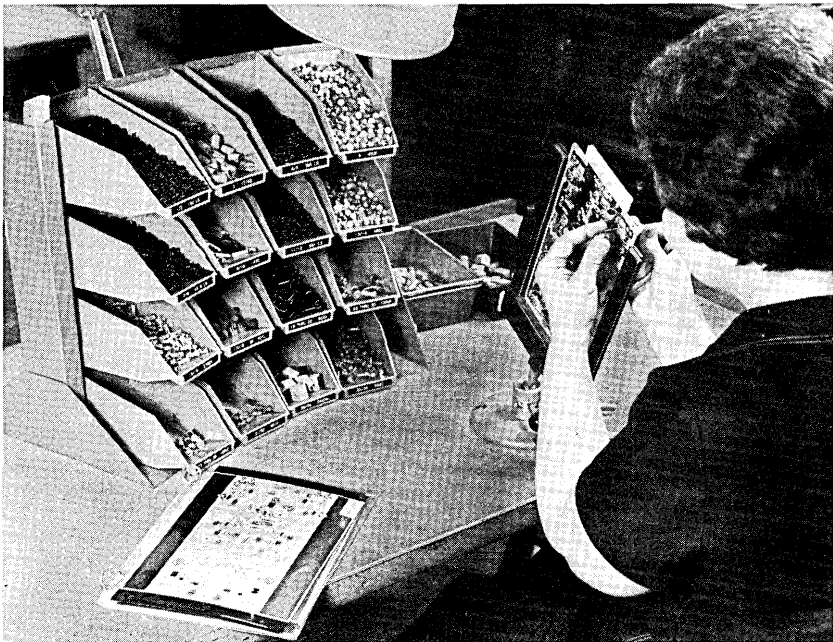


Fig. 1—Hand insertion position showing components in the rack and the board holder with PWB.

In the center of this semicircular arrangement is a spring-loaded device for holding the printed wiring board so that both hands of the operator are free for assembly work. The usual hand tools for this operation are also available. New and completed work is stored behind the operator on a hand truck.

2.2.3 *Visual Aids*

Visual aids are essential in any assembly shop if efficient use of time and a high degree of quality is to be expected. Two visual aids used in the D2 Channel Bank shop proved to be very helpful to the hand assembly personnel, semiautomatic machine insertion operators, and to the inspectors.

The inspection of interfacial "C" strap connections, which are the first items put on the printed wiring boards, is made faster and more accurate with the aid of a sheet of black phenolic cut to the same shape as the printed wiring board to be checked. Holes were drilled in the same location as each "C" strap connection with the remainder of the black phenolic being left blank. This configuration allows the inspector to see only the "C" strap connection area and makes it quite easy to determine if a connection has been omitted or not soldered. See Fig. 2.

The other type of visual aid is used where some complexity is involved in hand assembly of components. This visual aid consists of a color photograph of the printed wiring board with a series of numbers showing where components are to be inserted. These numbers indicate the correspondingly numbered bin where the operator will find the proper component. Those components actually inserted by the operator are shown on the visual aid. The industrial engineer specifies these components which will result in the most efficient learning pattern for the operator. See Fig. 3.

2.2.4 *Component Insertion Machines*

High-volume production requires the use of automatic machinery whenever possible. In manufacturing the D2 circuit packs, many of the components are inserted using automatic insertion equipment. Components with either axial or radial leads are inserted.

Axial components such as diodes, resistors and capacitors are inserted using Dual Center Distance (DCD) and Variable Center Distance (VCD) insertion machines (Fig. 4). The VCD is capable of inserting components on centers varying from 1 to 3 cm. The variability feature permits a maximum number of components on a circuit pack to be inserted. The DCD machine is used primarily for inserting "bulk"

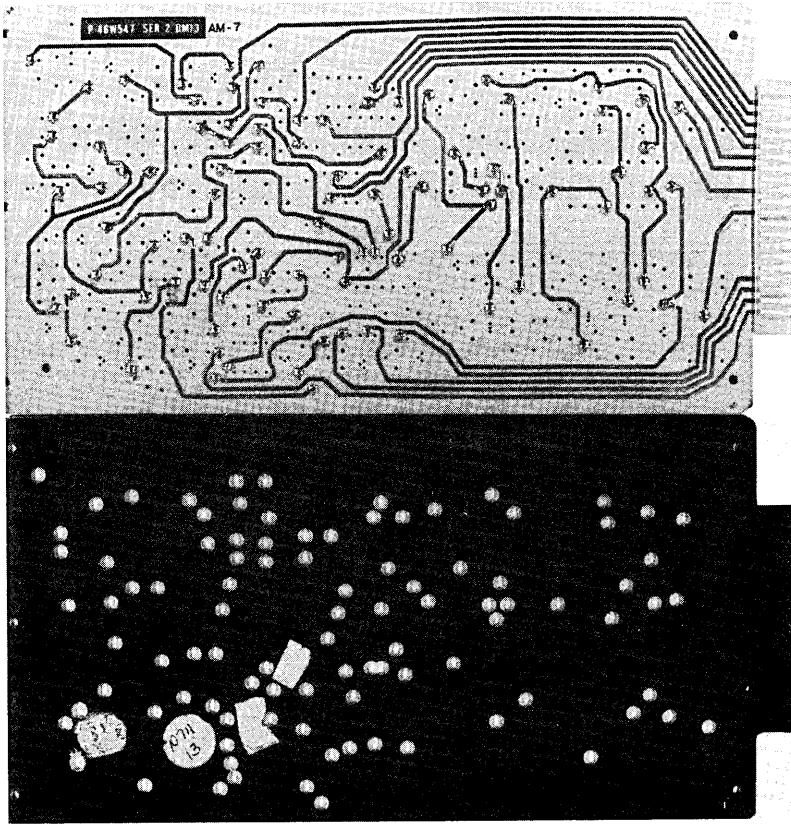


Fig. 2—C-strap inspection template.

components. The insertion machines are controlled by either a punched tape or by a small computer.

The normal work table has two positions: one rotated 180 degrees from the other. This permits polarized components such as diodes and tantalum capacitors to be inserted without stocking both polarities on the sequencer. A modified work table has recently been introduced with four positions at 90-degree intervals permitting even more components to be inserted. Radial leaded components are inserted using similar commercial equipment.

2.2.5 Component Sequencing

A normal part of any automatic component insertion operation is component sequencing. In this operation, the components are taken

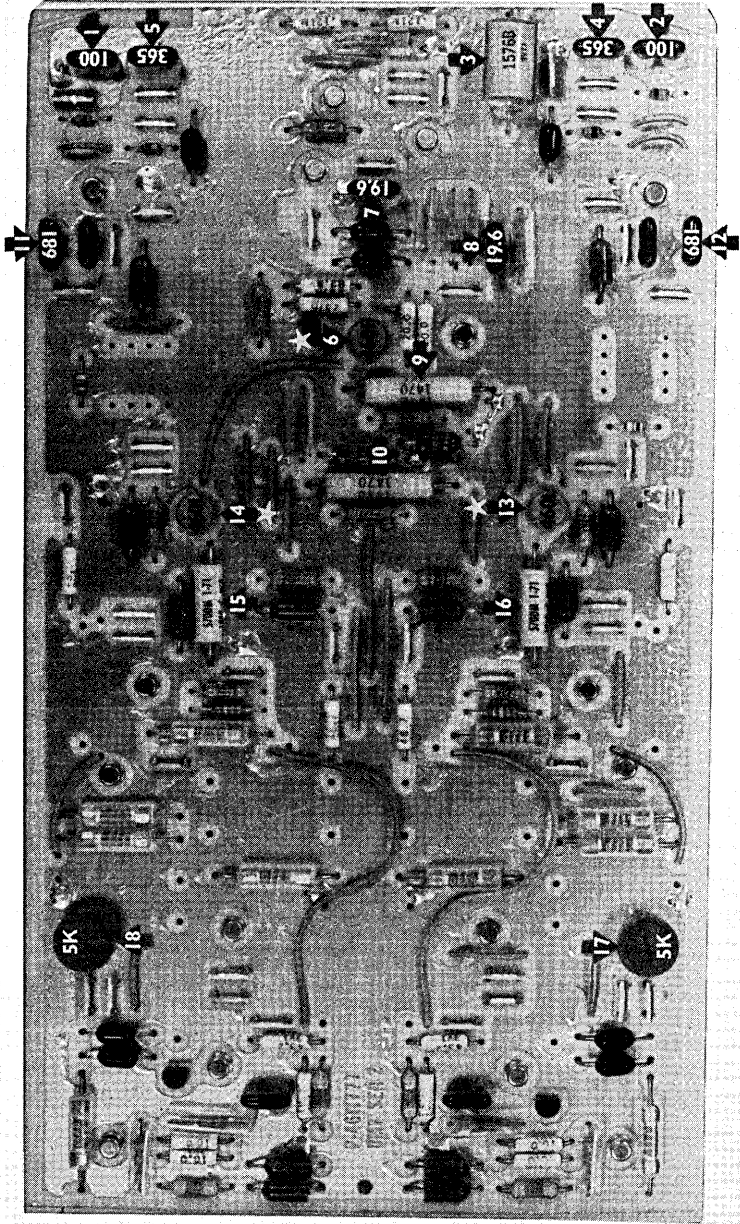


Fig. 3—Hand-assembly visual aid.

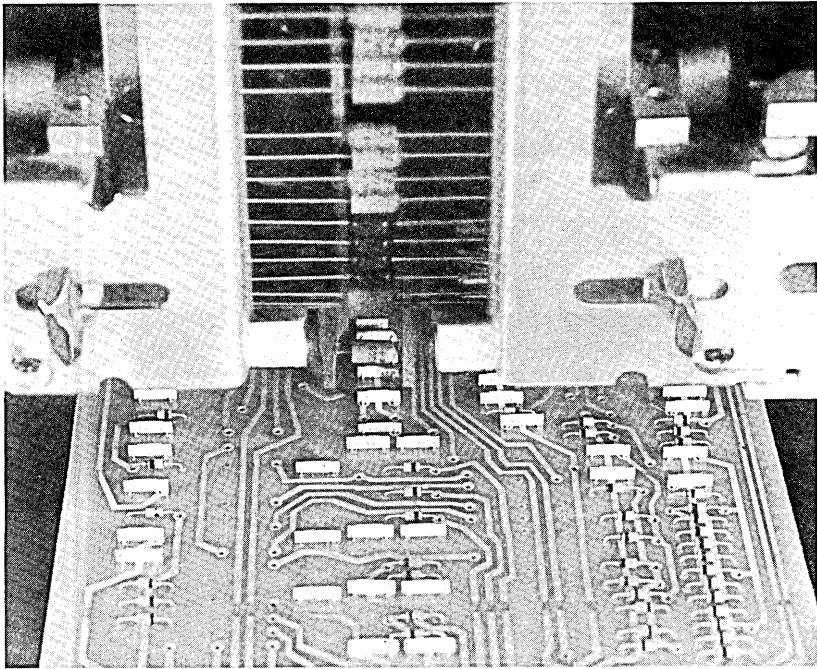


Fig. 4—Axial lead component insertion machine.

from reels of identical components and put onto reels of mixed components programmed for a specific circuit pack. The only exception to this general rule is in the insertion of diodes (458 series). Since they are used extensively throughout many of our circuit packs, we found it practical to use one machine to insert only these diodes. Proper polarity is achieved by stopping the machine after all diodes of one polarity are inserted, reversing the board and inserting the remainder.

An axial leaded sequencer was purchased from a commercial supplier. This machine has the capability of selecting from 39 types of components, sequencing in the reverse order to which they are to be used and then taping and winding the components on a reel. Only 39 types of components can be used on one code. A very thorough study was made to determine the maximum utilization of the machine and the maximum coverage of components and circuit packs.

The radial leaded 257-type resistors are used numerous throughout the D2 circuit packs. A commercial machine takes the container in which the resistors are shipped and vibrates the resistors down a track

where they are sequenced automatically first onto a conveyor, and then to a taping station. This sequencing machine has positions available for 30 different values of resistors. These positions can be changed for different circuit packs. As in the axial leaded component sequencer, maximum utilization was determined for the machine to cover as many codes as much as possible with a minimum of change over. See Fig. 5.

The possibility of selecting the wrong value of 257 resistor for a particular machine head appeared to be quite high. Also, the 257 type does not lend itself to rapid visual identification of component values. Therefore, a test station was designed to measure the value of each resistor on its way to being taped. See Fig. 6. The only resistor that is sequenced here is a 257J which has a 3-percent tolerance. Test limits were set at 5 percent in order to detect gross errors but not to interfere with normal product variation. After introducing this operation, a further benefit was soon discovered. The sequencer would occasionally remove small chips from the resistors which changed their values slightly (5 to 10 percent). These chips were nearly invisible and might not have been detected by any other tests in our shop. This caused us to closely examine the various mechanical features of the sequencer to find and eliminate the source of the chipping.

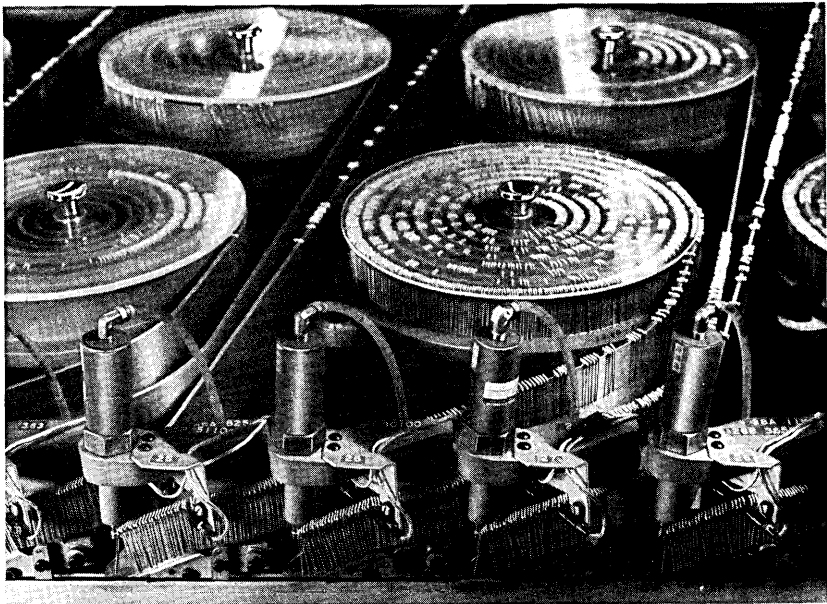


Fig. 5—Resistor containers on radial sequencer machine.

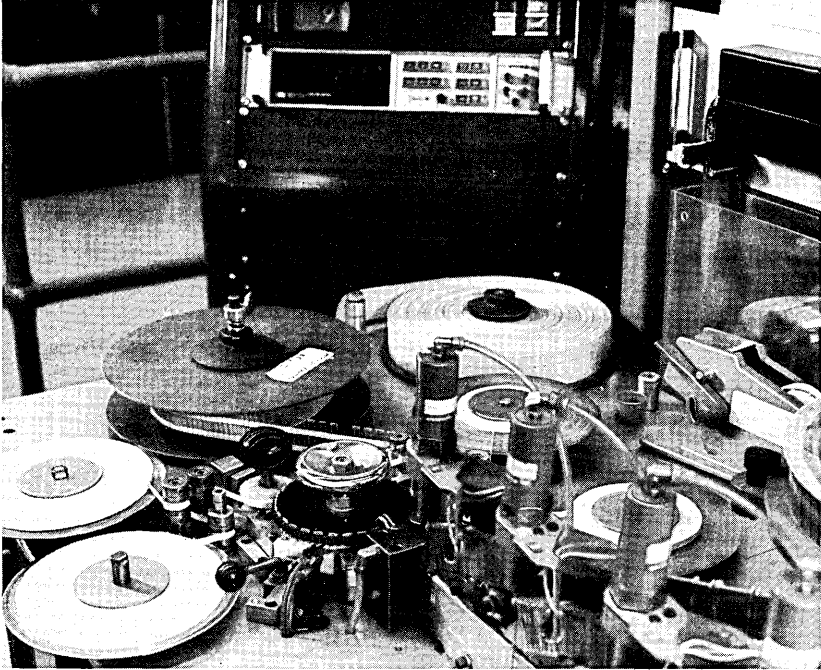


Fig. 6—Resistor taping station on radial sequencer machine with test station in background.

2.2.6 *Mass Soldering*

Much material has been written about the benefits of mass soldering and its acceptance in the printed wiring board field. The major considerations of the mass soldering machine for the D2 shop were quality, reliability, and ease of maintenance. Several units were investigated before one commercial machine was accepted.

2.3 *Precision Thin Film Resistors*

The D2 film circuits include 19 resistor codes with two to six resistors per code. The resistors range in value from 150 to 77,000 ohms. Several codes have initial resistance tolerance requirements of 0.04 percent absolute, and 0.02 percent matching when specified at 38°C operating ambient temperature. These exceptionally tight tolerances were required for circuits in the coder and decoder in order to attain the high level of precision to ensure toll grade performance for the D2 Channel Bank.

The D2 Channel Bank precision resistor networks were the first film-integrated circuits processed in the Merrimack Valley Process Capability Laboratory (PCL). The reasons for the introduction of this project through the PCL were to evaluate lead frame bonding, precision anodizing to extremely tight tolerances and to prove in the PCL facilities.

The initial work in the PCL brought the following recommendations:

- (i) Reduce the numbers of circuit sizes to aid future production inventory control of prescored substrates. This was before laser scribing was introduced as a production tool.
- (ii) Change from the Ta-Ni-Cr-Cu-Pd metal system to Ta-Ti-Au, to allow the use of stable gold-gold thermo-compression bonds for the lead frames. Minor changes in contact pad locations permitted a common lead frame to be used on all circuit sizes. This lead frame allowed simultaneous bonding instead of the individually bonded nail-headed leads.
- (iii) Reposition resistors to permit the use of a single substrate anodizing head for all codes to minimize the anodizing tool cost.¹

These recommendations were carried out and reasonable yields were realized in production through the use of a computer-directed thin-film trim anodization processor and tester, which would optimize total adjustment time with respect to a minimum chance of overshoot.

The capability of trim anodization to precise resistor values depends not only on the basic absolute accuracy and repeatability of measurement, but also on the method of controlling the percent resistance change per anodization cycle. Using constant current as the anodization power sources, one method of controlling the percent resistance change is by reducing the current. Another method is to reduce the time per anodization cycle. These considerations led to a binary step approach which would allow for a total variation of 50 percent in the anodization constants and still not overshoot by more than one binary step. A moderately-fast precision-measuring system coupled with a small process control computer was used to obtain an efficient shop set up and use.²

III. TESTING AND SHIPMENT

Concurrent with the production planning, test planning took place. It included the development of an overall test philosophy and the design, prove-in, and introduction of 14 new test sets to the production floor.

3.1 *Test Philosophy*

The development of a test philosophy involved considerations of the product to be tested (both physical and electrical), and the methods to use for both an inexpensive and thorough test. The philosophy evolved gradually with inputs from the test engineers, Bell Laboratories engineers, and shop supervisors. The following test philosophy was adopted for the circuit packs:

- (i) Component verification and circuit pack integrity test,
- (ii) Module test,
- (iii) Circuit pack test,
- (iv) Common equipment terminal test.

For the bays, a two-step procedure was adopted:

- (i) Back plane wire test,
- (ii) Bay wire test.

These elements will be discussed in greater detail. Supporting the test philosophy are the test specifications (called X-Specs). These are Bell Laboratories controlled documents. For the D2 project, however, most X-Specs were initiated by the Western Electric test engineer who had the responsibility for designing the test facility for the code. Many discussions concerning testing techniques and parameters to be tested were held as the X-Specs began to take form. These discussions had the dual purpose of educating the Western engineer to design considerations in the product and the Bell Laboratories engineer to testing methods in the manufacturing environment. The net result was a set of test specifications that both parties were happy with and that formed a well planned foundation for the design of test facilities. We feel, particularly as products become more complex, that test considerations must be included during the design of the product. Based on our experience in D2, the test and design engineers should begin discussions very early in the design of the product.

3.2 *Component Verification and Circuit Pack Integrity Test*

The circuit packs in the D2 Channel Bank are characterized by discrete components attached to double-sided printed wiring boards. Some of these boards are very densely packed as shown in Fig. 7. Testing such a circuit pack would be easier if there were a way to check the overall integrity of the board and to measure the components to specified tolerances. To accomplish these ends, a commercial test set was purchased. The test set uses punched tape control to set the

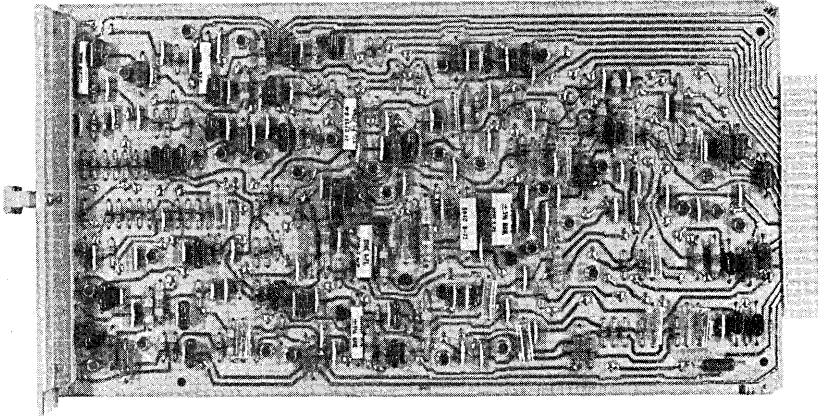


Fig. 7—DM2 Circuit pack.

input parameters for each test. A vacuum fixture and spring loaded plungers give random access to the circuit pack under test. See Fig. 8. The tests are made automatically (about 1000 per minute) with any defects being printed on paper tape for analysis.

The effect of parallel components is removed by using the guarding technique. In Fig. 9, let R_x be the unknown resistor with R_o and R_i as parallel components. A specified voltage is supplied at point A by the input amplifier. A ground is put at point Z by the test set. The other side of R_x (point B) is connected to the input of an operational amplifier. Since this point is a virtual ground, no current flows through R_i . The current through R_o does not affect the measurement. The current through R_x is sent into the op amp and is easily measured.

Each circuit pack is tested twice; the first pass is for overall card integrity; the second pass is an actual component check. Some of the typical tests used on the first pass are as follows:

- (i) Continuities—most “long” paths on the card are checked for continuity, especially those paths that have C-straps or plated-through holes.
- (ii) Shorts—tests for solder crosses on both sides of the circuit pack are made.
- (iii) Diodes—forward voltage drop is measured. For example, the 458 series is measured to $0.7V \pm 10$ percent. This has proven sufficient to detect most diode failures.

- (iv) Zener diodes—measured to specified voltage ± 10 percent.
- (v) Transistors—junctions CB and BE are tested to $0.7V \pm 20$ percent.
- (vi) Transformers and Relays—resistance of windings measured to ± 20 percent.
- (vii) Jacks—continuity through jack contacts.

After the first pass, the defects are removed and the unit is given a second test. Resistors are measured to slightly more than their rated tolerance. For example, a 257J-type resistor with a 3-percent rating is tested to 3-1/2 percent. Capacitors could be measured but are not in our application. Most have small values and would be very difficult to check. Capacitors are given an implied test during the functional test of the circuit pack.

3.3 Module Test

Some of the circuit packs contain modules which are attached to the master board. The modules are given a functional test prior to

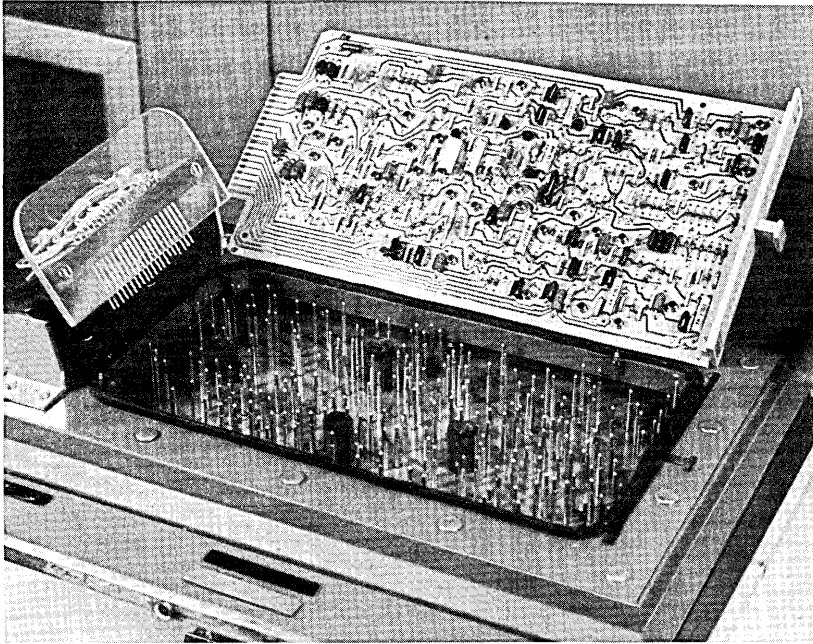


Fig. 8—Component verification test fixture.

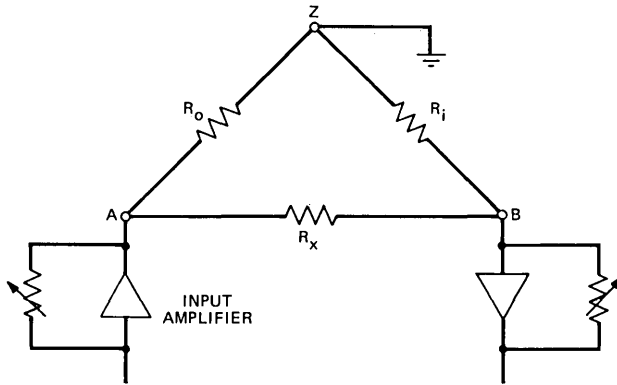


Fig. 9—Circuit diagram of "guarding" technique.

mounting on the master board. This test simulates the actual operating conditions of the module in the circuit pack. Input and output impedances are carefully matched. Driving voltages and currents are chosen to duplicate actual conditions. The purpose of this test is to eliminate the work of removing a module if it causes a test failure and in some cases to test a specific function that would be difficult to test in the complete assembly. For example, the zero code suppression of the coder is checked during a module test.

3.4 Circuit Pack Test

Each circuit pack is given a complete functional test during the testing procedure. In general, the unit-under-test is placed in a circuit that simulates conditions in the bay. Actual D2 circuit packs are used to provide the input signals and output loads. This technique provides the same operating conditions found in the bay. Test measurements are divided into two categories: digital and analog.

The digital measurements are made with an oscilloscope. Factors such as pulse height, width, rise time and phase relationship are measured for all pulses as they leave the circuit pack and for specific internal wave forms. The present test sets are manual or semiautomatic in operation.

However, increased volume permits the introduction of automatic pulse-measurement techniques. When the automatic sets are completed, the original manual sets will be used for analysis and repair of defects. This is the normal evolution of test sets.

Analog circuit packs are tested in a different manner. Basically, they are tested within a complete operating system.

A test set (Coder-Decoder Test Set) has been designed which is a two-digroup system with the transmitter looped into the receiver. The unit-under-test is plugged in as a replacement for one of the units in the test set. The unit then is tested using transmission measurements such as noise and signal-to-distortion. On some units, internal wave forms are examined for specific attributes.

All circuit packs are "margin" tested to reduce the chances of incompatibility. In margin test, the power voltages to the unit-under-test are raised or lowered by 10 percent, while the voltages to the load circuits are held fixed. This is an attempt to find the marginal failure that will cause trouble later in the field. The 10-percent figure is based on twice the allowable deviation of the bay power supplies in order to allow for differences in bays and other circuit packs. In most cases, the same performance limits apply to a unit under margin test as to a unit powered with the normal voltages.

The testing method described is referred to as "testing product with product". It is not a perfect test method since the selection of "standard" units within the test set is not well defined. An alternative method would be to actually measure the output of individual circuit packs after stimulating them with artificial circuits. This method has its hazards also. All inputs must be well defined and exactly simulated. The detectors must be sensitive enough to detect small traces of noise and other imperfections on waveforms such as a PAM sample. We must also be able to define exactly how much signal degradation is permissible in each unit of the entire channel bank. Not that the alternative technique is impossible, the authors merely feel that this is impractical to implement for a new system that is still going through its growing pains.

3.5 *Terminal Test*

All of the common equipment circuit packs are given a terminal test. In this test, groups of units are plugged into an actual D2 bay. Measurement of idle circuit noise, signal-to-distortion, gain tracking and crosstalk are made. Additional tests of signaling and alarm circuits are made. Most of these tests are duplicates of tests that will be made by the operating companies. The rationale for making these tests is to provide the customer with an assurance that his product is free from incompatibilities. Many design improvements have been initiated as a result of difficulties in this test area. This is expected when the complexity of the completed system is considered. The terminal test also gives engineers the opportunity to continuously monitor system performance and to detect unfavorable trends.

The Terminal Test Set was designed to provide random access to each of the 96 channels in the D2 bay. To accomplish this, two Western Electric crossbar switches were mounted to each bay. See Fig. 10. These switches (each a 10 by 10 matrix) can be individually controlled to allow for random access to one channel in the transmitting and receiving sections. The test set itself contains a minicomputer which controls the crossbars and the measuring instruments. Operation is automatic after an initial alignment of the system. This set could be used (with minor modifications) to test any voice frequency channel bank. All that would be necessary to adapt to another system is a different program for the computer and a scheme to attach the crossbar switches to the bay.

3.6 Data Collection

A data collection and reduction scheme is needed for any complex test system. To be reliable, the data should be collected automatically.

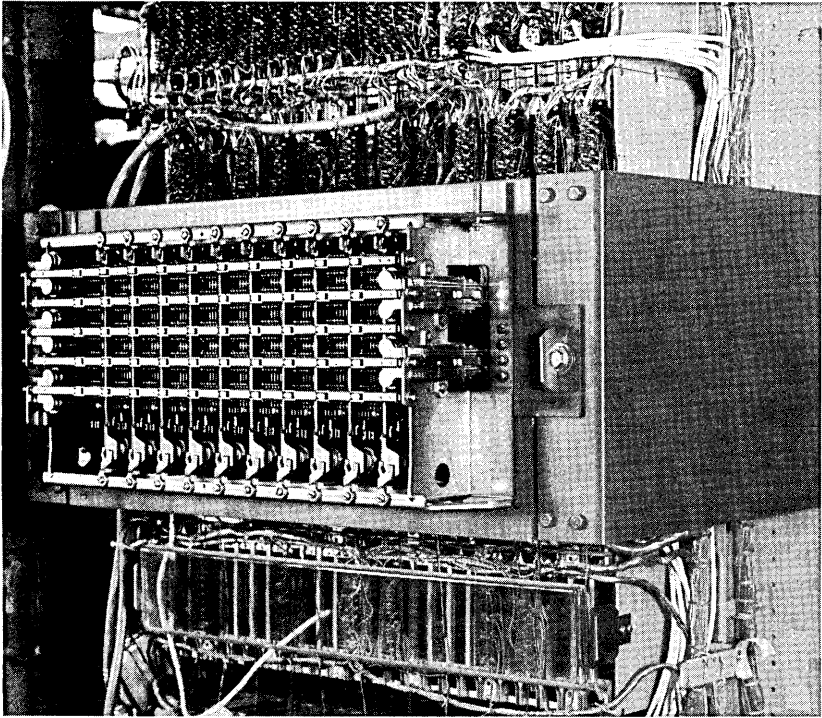


Fig. 10—Crossbar switch attached to D2 bay.

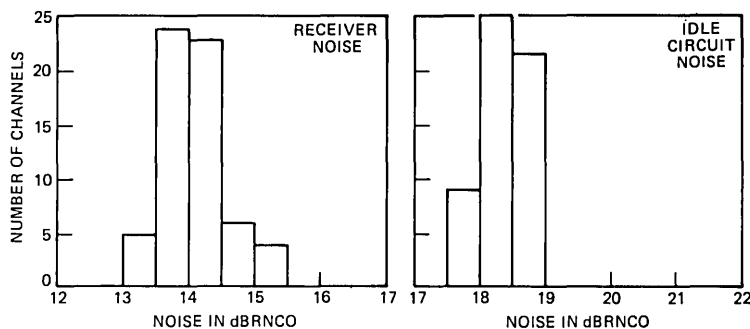


Fig. 11—Histogram plot of typical data.

Such an automatic system was designed into an early version of the Terminal Test Set.³ Each test reading was sent via a Data-Phone link to a process control computer. The computer stored the readings of a complete system test (304 readings) and at the end of the test gave the following outputs:

- (i) A UNIT SUMMARY of averages by test number.
- (ii) A set of punched cards with all 304 readings.
- (iii) Update of a history file of readings.

Each week or on demand the history file was printed out. The cards were available for off-line analysis such as histogram plots, trend studies and others. See Fig. 11. Although the data collection feature was not included in the newest Terminal Test Set, these data have given us a very complete picture of how the D2 Channel Bank met its requirements and which tests should be worked on to improve the performance. It also gave a data base to evaluate the effects of any major change to the system.

3.7 Bay Test

The D2 bays are given a wire verification test and shipped directly to the customer. Two steps of testing are used: a process test and a final bay test.

In the process test, each shelf backplane is given a wire verification. A tape-controlled test set checks for continuity (< 5 ohms) and insulation (> 1 megohm at 500 volts) for all wires.

For final test, a larger tape-controlled test set is used to test all wires for opens and shorts as in the process stage. Additional tests are

made on the Carrier Group Alarm by operating the various relays in these circuits. Another test is made to see if all the empty terminals in the bay are really empty.

The test sets used for bay test are generally purchased from outside suppliers. They are functionally similar, and use a punched tape to supply the input parameters. In the D2 bays, the test program has about 8000 tests and takes about ten minutes to run. Defects are printed for analysis and repair.

Programming these test sets requires supplying a list of wires to be tested and specifying the type of test to be made. Creating a test program by manual methods would be very laborious. Therefore, a computer program that generates the test tape has been developed.⁴ The program, which runs on a general purpose computer, is general in nature and is not restricted to any specific type of bay. Test tapes for other product lines have been produced using this program.

There are two inputs to the program:

- (i) A Fixture List telling how the test set is connected to the product.
- (ii) A Wire List telling how the product is wired and the type of test desired.

The output is a punched tape which is capable of controlling the test set. The program has the following features:

- (i) Validation of fixture and wire listings.
- (ii) Addition of empty terminal tests.
- (iii) Automatic sorting and printing of multiple wire runs.
- (iv) Listing of Machine to Product Points.
- (v) Punching of cards of the test program.

3.8 *Shipment*

D2 circuit packs are packed for shipment in two ways: in systems according to various list structures and singly for channel units and spares.

In systems packaging, the circuit packs are packaged by list number. Each circuit pack is fitted into a corrugated shipping container with grooved polystyrene inner details that guide and hold the printed wiring board.

The spare circuit packs and channel units are packaged in individual containers made of two pieces of expanded polystyrene to support the circuit pack properly. The packaged unit is then inserted into a form-fitting corrugated sleeve and secured at the ends with gummed tape.

Similar-shaped circuit packs use similar containers to keep the number of different containers to a minimum.

The advantages of using individual containers are (*i*) the existence of suitable reusable containers in which to return defective units from the operating companies to the repair centers and (*ii*) the freedom in stocking circuit packs by the Merchandising organization.⁵

IV. INITIAL PRODUCTION

The initial production phase began in the summer of 1969 and continued for approximately six months until the first systems were shipped in December 1969. This was an exciting time as all the planning and organizing began to show results. It was also a period characterized by close cooperation between Bell Laboratories and Western Electric as all the final details were examined and reconciled.

As finished drawings of each code and its detailed piece-parts became available from drafting, models were constructed in the shop. Both the supervisory and assembly personnel became familiar with the details of each code. They also contributed valuable advice on potential trouble areas. When the shop and engineering found trouble areas, solutions to correct them were proposed to Bell Laboratories. By this time, close personal relationships between the individuals had formed. Thus, both sides respected the advice and judgment of the other. Problems were discussed; solutions were synthesized; and final corrective action was formulated in a highly efficient manner.

After models were built, they were turned over to a Western Electric test planning engineer. He examined the electrical performance of the unit, using the experience he gained during the field trial evaluation period. Again, problem areas were discussed with his Bell Laboratories counterpart and solutions were formulated. Once the test engineer was satisfied with a unit, he shipped it to the Bell Laboratories engineer for final evaluation.

By October of 1969, most of the circuit packs had had their initial evaluation. An actual production bay was erected in our shop and the final evaluation phase, prove-in of the system, began. This was one of the busiest periods of the entire project. As the system gradually took shape in the bay, each circuit was closely examined for proper operation. The results of this examination were fed back to the individual circuit pack test in an effort to duplicate actual operating conditions. The entire group, Bell Laboratories and Western, equipment and circuit engineers, were collectively working on the single project of final prove-in of the

D2 Channel Bank. Traditional lines between designer and manufacturer were set aside as each individual contributed according to his ability. It was a proud moment when the first system was packed and sent to a customer.

V. CONTINUING PRODUCTION

With the first systems shipped to the customer, Western is now in a continuing production phase. Some aspects of this phase are customer support, cost reduction, and design improvement.

5.1 *Customer Support*

Western Electric takes the position that a customer should be satisfied with a system he purchases. A new system is likely to cause some difficulties in the earliest installations. Therefore, we have traveled to a few field installations in order to offer assistance. The primary purpose of these trips is to get the equipment on-line and make money for the telephone company.

We at Western benefit from these trips by finding inadequacies in our own production methods. In the first year, Western Electric engineers made four trips to field sites due to requests from the telephone companies. In each case, the engineers (including the authors) returned with a greater appreciation of how the units arrive and how they are used in the field. This valuable information was fed back to our associates for corrective action. In addition to field trips, we consulted with operating company personnel by telephone to offer assistance in solving installation and line-up problems. The customer support effort has provided the required installation assistance and also brought information back to Western to improve the product.

5.2 *Cost Reduction*

Cost reduction is a very important part of all Western Electric product lines. The Bell System strives to continue to offer service at low prices in spite of rising costs of labor and materials. Cost reduction is one tool used to achieve this end.

In the case of a project such as the D2 Channel Bank, cost reduction can come from sources such as circuit redesign, physical redesign, component substitution, and automation. One particular case, which included these elements, is discussed to show the technique of cost reduction.

The case involved the filter and gates in the transmitting section.

This unit had a large number of module boards and LC filters mounted in metal cans. It had some redundant circuitry and it was produced in high volume.

One aspect of the case was to reduce the number of module boards from 19 to 1. This involved a layout of the circuit and repositioning circuit elements from the many small module boards to one larger module board.

The reason for a savings at this point may not be obvious and so will be expanded. Previously, the unit had 19 modules of six different codes. Each one was a different size and shape. Some were used once; others were used eight times per circuit pack.

All these odd sizes and shapes discouraged the use of automatic component insertion, mass soldering, etc. In addition, each individual PWB had a cost that was significant. The use of only two PWB's, one main and one module board, saves in material costs and also allows the use of automatic equipment.

Another portion of the case was to remove the metal cans from the filters. This may seem elemental and even a reflection on the designer's ability when he specified the cans. Actually, this is far from obvious. Removal of the cans included a detailed study of the effects of cross-modulation and other circuit effects. This investigation showed improvement in crosstalk performance and no degradation in any other parameters. Removal of the cans would have been extremely risky if attempted before the system was well characterized and being routinely produced.

Another part of the project is to utilize Cap-Pak in place of discrete capacitors. Cap-Pak is an assembly of groups of capacitors mounted in a common case similar in shape to a transformer. Savings result from reduced labor effort in the manufacture of the capacitors and reduced labor effort in inserting the capacitors.

The example above is a typical large cost reduction case. Numerous small cases are constantly in progress as the manufacturing engineer constantly searches for ways to improve his product.

REFERENCES

1. Mushial, R. G., unpublished work, 1970.
2. Raymond, D. H., "Computer-Directed Anodization and Testing of Precision Thin Film Resistor Circuits," *The Western Electric Engineer* (April 1971), pp. 2-9.
3. Batson, J. E. D., and Fiore, A. R., unpublished work, 1969.
4. Dickerson, N. O., unpublished work, 1970.
5. Salvage, C., unpublished work, 1970.

Optical Power Flow in Multimode Fibers

By D. GLOGE

(Manuscript received May 8, 1972)

Loss, coupling, and delay differences among the modes of multimode fibers influence their transmission characteristic in a complicated way. An approximation of the modes by a continuum leads to a comprehensive description of these interrelations. We relate the mode power distribution to the far-field output and calculate these distributions as functions of the fiber length and the input. We report measurements of the far-field distributions at various lengths of a cladded low-loss multimode fiber. A comparison of theory and experiment yields a quantitative estimate of the mode coupling involved. We associate this coupling with random irregularities of the fiber configuration and straightness, and construct a quantitative model of such irregularities.

I. INTRODUCTION

Some sources considered for use in optical communication systems have a spatially incoherent or multimode output and require overmoded fibers for efficient transmission. The fibers consist of a highly transparent core surrounded by a cladding of lower refractive index. Liquid core prototypes with losses as low as 20 dB/km have been built.¹ Solid multimode fibers have slightly higher losses.² A recent study of their propagation and dispersion characteristics³ showed a rather intricate behavior complicated by the fact that hundreds of modes could propagate simultaneously. These modes underwent a perpetual mixing process. The attenuation coefficient appeared to vary from mode to mode causing a relatively fast loss of the high-order modes.⁴ An increase of delay with mode number (and fiber length) was observed as expected, but mixing and attenuation seemed to influence this relationship in a complex way.

An exact knowledge of the processes involved is of considerable interest not only to understand the sources of loss in the fiber, but in order to determine the signal distortion in long fibers. It has been predicted⁵ that under certain circumstances increased mixing reduces the signal distortion (ultimately forcing all energy to propagate at an

average velocity). But it remains to be determined what actual signal improvements can be gained in practice from this effect. Previous investigations of these problems,^{5,6} although suited to show the concepts involved, were limited to model studies involving relatively few modes.

In this work, we replace the modes by a continuum. This results in a relatively simple differential equation which describes the power distribution as a function of time, fiber length, and the continuous mode parameter. The differential equation can be solved rigorously for certain conditions which satisfactorily match experimental results. Explicit relations result which describe the propagation characteristics as a function of the modal coupling, attenuation, and delay coefficients. The coupling is then related to specific imperfections of the configuration or straightness of the fiber.

This paper is primarily devoted to the time-independent solution of the problem. Signal distortion and, specifically, the (baseband) impulse response of long fibers can be derived from a slight modification of the above equations and this will be done in a subsequent paper. The concept underlying our results developed from experiments with solid-core fibers³ but, in the meantime, measurements of long liquid-core fibers⁷ have proven that these fibers follow the same concept.

II. TRANSITION TO MODAL CONTINUUM

For large mode numbers, the characteristic mode parameters change so little between neighboring modes that their discrete values can be replaced by one continuous variable. Consider the two-dimensional dielectric guide—a thin film, for example—sketched in Fig. 1. We assume that the relative index difference

$$\Delta = 1 - \frac{n_c}{n} \quad (1)$$

between core index n and cladding index n_c is small compared to unity. In that case, the critical angle* for total internal reflection

$$\theta_c = \sqrt{1 - \left(\frac{n_c}{n}\right)^2} \approx \sqrt{2\Delta} \quad (2)$$

is small as well, and we can use small-angle approximations in the following relations.

Within the high-index material, the field distribution of the m th mode is essentially sinusoidal (see Fig. 1) with transverse wave number

* Defined here as the angle measured from the reflecting surface (see Fig. 1).

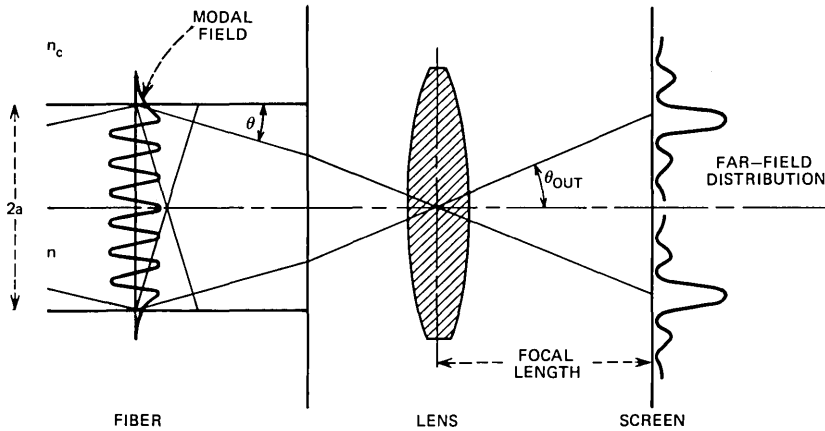


Fig. 1—Sketch to illustrate the wave nature of the modes in the dielectric slab and in the fiber.

$$u = \frac{\pi m}{2a} \tag{3}$$

where $2a$ is the guide width. If $k = 2\pi/\lambda$ is the free-space wave number, the propagation direction of a mode (i.e., its representative plane wave) follows from

$$\theta = \frac{u}{nk} = \frac{m\lambda}{4an} \tag{4}$$

Because of Snell's law, this angle becomes

$$\theta_{out} = \frac{m\lambda}{4a} \tag{5}$$

outside of the guide. In the far field (or the focal plane of a lens), the plane waves are concentrated about the directions $+\theta_{out}$ and $-\theta_{out}$. The aperture of the guide determines the angular concentration of the two far-field "spots." If the guide width is $2a$, the spot width is of the order of λ/a .

As we learn from (4), the propagation directions of neighboring modes differ by

$$\Delta\theta = \frac{\lambda}{4an} \tag{6}$$

and hence by $\lambda/4a$ outside the guide. The modes thus form a partly

overlapping sequence of spots in the far field, ordered according to mode number. Consequently the far-field distribution represents a direct image of the modal power distribution.

The transition to the modal continuum uses a continuous angle θ instead of the discrete values (4). In this way, we arrive at a continuum of plane waves which, in the following, will be represented by rays. The power distribution $P(\theta)$, in this continuum, is obtained by replacing θ_{out} by θ in the (average) far-field power distribution.

The cylindrical configuration lacks part of the conceptual clarity associated with the plane-wave representation, but a formal similarity permits us to arrive at an equivalent ray model which is satisfactory for almost all problems related to multimode fibers. We refer again to Fig. 1, considering now a cylindrical core of radius a imbedded in cladding material. The modal field distributions are given by Bessel functions. In the case of a small index difference, there are degenerate mode pairs ($HE_{l+1,a}$ and $EH_{l-1,a}$) whose transverse mode number u is determined by the q th root of the Bessel function⁸

$$J_l(ua) = 0. \quad (7)$$

Here l is the azimuthal order number. Fig. 2, which lists a few low-order

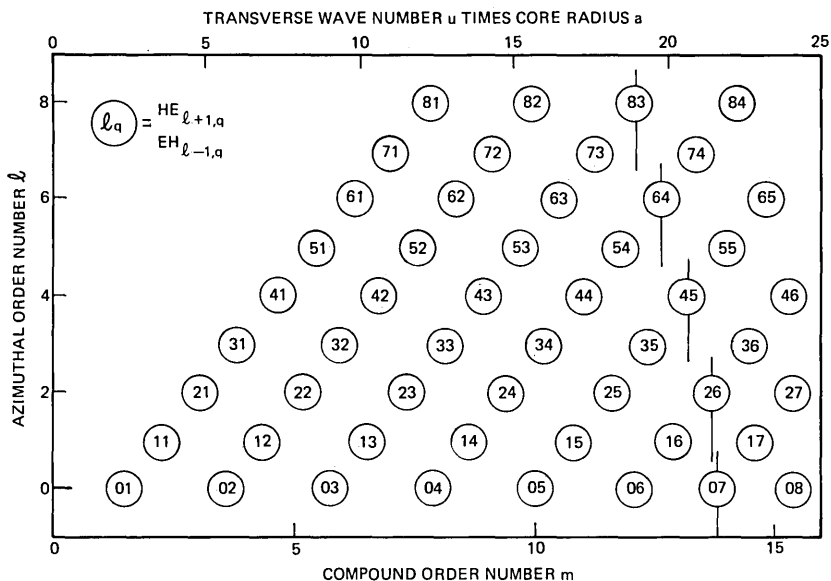


Fig. 2—The order numbers of degenerate fiber modes plotted versus ua and the effective group number $m = 2au/\pi$.

roots, has the purpose of indicating how much the exact roots deviate from the approximation

$$u = \frac{\pi}{a} \left(q + \frac{l}{2} \right), \quad (8)$$

which is to be used in the following. The group with $q + l/2 = 7$, for example, which has $ua = 22$ according to (8), is marked by vertical lines in Fig. 2.

Most problems of interest in multimode fibers (coupling, scattering loss, delay) require only the transverse wave number u for a satisfactory description of each mode. Furthermore, u can be related to a propagation angle θ and a far-field angle θ_{out} —in exact formal agreement with (4) and (5)—through an Hankel transformation of the mode field at the fiber end. This transformation shows that a mode of azimuthal order l produces l far-field spots located on a circle which is defined by the angle θ_{out} of (5). Figure 1 illustrates the situation if viewed as a meridional cross section through a cylindrical configuration.

These facts suggest a description of the cylindrical modes by a single mode number

$$m = 2q + l. \quad (9)$$

Equations (1) through (6) then obtain for the cylindrical guide as they do for the slab. The important difference is hidden in the fact that m of (9) comprises a group of modes with different q and l . As can be seen from Fig. 2, the number of possible combinations for a given m is the nearest integer below $m/2$. As mentioned earlier, every combination of q and l represents two (degenerate) modes. Consequently, each m describes a group of (approximately) m modes. In the far-field pattern, this mode group covers an annular area of "radius" θ_{out} and approximate "width" λ/a .

The transition to the continuum again converts θ to a continuous parameter. But θ is now considered as a radial variable which covers the solid angle $\pi\theta_c^2$. We have a conceptual model which consists of a continuum of rays within the cone $\pi\theta_c^2$, whereby the modal power distribution $P(\theta)$ is obtained by replacing θ_{out} in the (average) far-field power distribution by θ of (4).

To compute the total number of modes, we determine the highest possible group number $m_c = 4an\theta_c/\lambda$ by inserting θ_c into (4). If we consider also that each group has m modes; and each mode has two possible states of polarization, we have for the mode volume

$$\sum_{m=1}^{m_c} 2m = \left(\frac{4an}{\lambda} \right)^2 2\Delta. \quad (10)$$

A comparison with the more accurate number $(2\pi an/\lambda)^2 \Delta$ from Ref. 8 gives an indication of the quality of the approximations used here.

III. POWER FLOW EQUATION

For the sake of simplicity, the following derivation is based on a model which seems to have limited validity at first glance. We assume that mode coupling takes place only between next neighbors. It will become apparent later that the error involved in this approximation is small if other modes couple also, but the coupling strength decreases sufficiently fast with the mode spacing. There is experimental evidence^{3,7} that a mechanism of this kind is indeed present in real multimode fibers. Fig. 3, for example, shows a measurement performed with the solid-core multimode fiber mentioned previously.^{2,3} The core diameter was $55 \mu\text{m}$, the relative index difference $\Delta = 0.0046$, and the nominal loss 33 dB/km . About 700 modes could propagate. By injecting a very narrow cone of light (through an index-matching cell), we excited about 150 of the low-order modes. This was measured by scanning the far-field of the output after 30 cm of fiber. Similar measurements with longer fibers revealed a slow but steady increase in the number of excited modes (the angular far-field width) with fiber length. This slow increase is considered as strong evidence of a power exchange which favors near neighbors and decreases rapidly with mode spacing.

To simplify matters, let us again consider the two-dimensional case first. As long as the coupling mechanism is a statistical process, we can

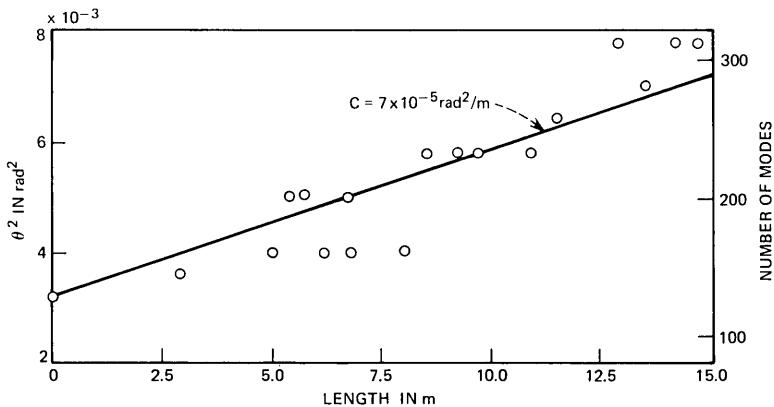


Fig. 3—The increase of mode volume with fiber length measured in a multimode fiber which propagated 700 modes.

ignore the individual mode fields and obtain the power distribution directly from some form of power rate equations.⁹ These consider the variation dP_m in the power P_m of the m th mode along a guide increment dz . In the time-invariant case, the variation dP_m has two causes: (i) dissipation and loss (scattering) to the outside, which we comprise in a term $-\alpha_m P_m dz$; (ii) coupling to other modes. Our simplified model assumes coupling between neighboring modes only. Thus, if d_m is the coupling coefficient between the modes of order $m + 1$ and m , we have

$$\frac{dP_m}{dz} = -\alpha_m P_m + d_m(P_{m+1} - P_m) + d_{m-1}(P_{m-1} - P_m). \quad (11)$$

The transition to the continuum requires power differences to be replaced by differentials. Especially, we set

$$\frac{P_{m+1} - P_m}{\theta_{m+1} - \theta_m} = \frac{dP_m}{d\theta}. \quad (12)$$

With $\theta_m - \theta_{m-1} = \Delta\theta$ from (6), we can rewrite (11) in the form

$$\frac{dP_m}{dz} = -\alpha_m P_m + \Delta\theta \left(d_m \frac{dP_m}{d\theta} - d_{m-1} \frac{dP_{m-1}}{d\theta} \right). \quad (13)$$

The remaining difference requires an analogous transition

$$d_m \frac{dP_m}{d\theta} - d_{m-1} \frac{dP_{m-1}}{d\theta} = \Delta\theta \frac{d}{d\theta} \left(d_m \frac{dP_m}{d\theta} \right). \quad (14)$$

After replacing the index m everywhere by a functional dependence of θ , we finally obtain the power flow equation

$$\frac{\partial P}{\partial z} = -\alpha(\theta)P + (\Delta\theta)^2 \frac{\partial}{\partial \theta} \left[d(\theta) \frac{\partial P}{\partial \theta} \right]. \quad (15)$$

In the cylindrical case, the index m stands for a group of m modes. To obtain the power equation for the m th mode group, we must therefore sum (11) over all m members. The coefficients α_m and d_m depend only on m , and hence are the same for all group members.¹⁰ However, the coupling to the lower group ($m - 1$) can occur only between $m - 1$ members. Thus

$$m \frac{dP_m}{dz} = -m\alpha_m P_m + m d_m (P_{m+1} - P_m) + (m - 1) d_{m-1} (P_{m-1} - P_m). \quad (16)$$

Using (12) and a transition analogous to (14), we obtain

$$\frac{\partial P}{\partial z} = -\alpha_m P_m + (\Delta\theta)^2 \frac{1}{m} \frac{\partial}{\partial \theta} \left(m d_m \frac{\partial P_m}{\partial \theta} \right). \quad (17)$$

With the help of (4) this finally leads to

$$\frac{\partial P}{\partial z} = -\alpha(\theta)P + (\Delta\theta)^2 \frac{1}{\theta} \frac{\partial}{\partial \theta} \left[\theta d(\theta) \frac{\partial P}{\partial \theta} \right]. \quad (18)$$

Because of the symmetry involved, we can expand α in the form

$$\alpha(\theta) = \alpha_o + A\theta^2 + \dots,$$

where α_o comprises loss common to all modes. A loss of this kind can later be accounted for by multiplying the final solution by a term $\exp(-\alpha_o z)$. For the moment, we ignore this part of the loss. Among the higher orders, the term $A\theta^2$ is the most important one, because it essentially comprises the loss caused at the core-cladding interface. This is so because the power density at the interface increases quadratically with the transverse wave number u of a certain mode⁸ and hence quadratically with θ . In the following, we retain only this important term.

The coupling coefficient $d(\theta)$ requires the same expansion. Its zero-order term is essential and cannot be accounted for later on. Although no estimates exist on the magnitude of other terms, the following derivation merely retains this first term. Its physical significance will become clearer as we proceed. Thus with

$$d(\theta) = d_o \quad (19)$$

or

$$D = (\Delta\theta)^2 d_o = \left(\frac{\lambda}{4an} \right)^2 d_o, \quad (20)$$

we can write (15) and (18) in the form

$$\frac{\partial P}{\partial z} = -A\theta^2 P + D \frac{\partial^2 P}{\partial \theta^2} \quad \text{for the slab,} \quad (21)$$

and

$$\frac{\partial P}{\partial z} = -A\theta^2 P + \frac{D}{\theta} \frac{\partial}{\partial \theta} \left(\theta \frac{\partial P}{\partial \theta} \right) \quad \text{for the fiber.} \quad (22)$$

The form of the last terms in (21) and (22) identifies next-neighbor mode coupling as a diffusion process in the continuum.

Solutions that are independent of z can be obtained from the substitution

$$P = Qe^{-\gamma z}, \quad (23)$$

where γ denotes a (power) attenuation constant related to the "steady-state" solution Q . Equations (21) and (22) take the form

$$D \frac{\partial^2 Q}{\partial \theta^2} = (A\theta^2 - \gamma)Q \quad \text{for the slab,} \quad (24)$$

and

$$\frac{D}{\theta} \frac{\partial}{\partial \theta} \left(\theta \frac{\partial Q}{\partial \theta} \right) = (A\theta^2 - \gamma)Q \quad \text{for the fiber.} \quad (25)$$

The first of these equations is satisfied by the Hermite-Gaussian, and the second by the Laguerre-Gaussian polynomials; both are well known from the theory of the open resonator. The attenuation parameters γ associated with each of these solutions increase with the order of the polynomial.

Both for the slab and the fiber, the solutions of least loss have the form

$$\exp(-\theta^2/\Theta_\infty^2) \quad (26)$$

with

$$\Theta_\infty = (4D/A)^{\frac{1}{2}}. \quad (27)$$

The power loss associated with this distribution is

$$\gamma_\infty = (AD)^{\frac{1}{2}} \quad \text{for the slab,} \quad (28)$$

and

$$\gamma_\infty = 2(AD)^{\frac{1}{2}} \quad \text{for the fiber.} \quad (29)$$

The distribution (26) constitutes an optimum balance between the loss in high-order modes and the steady outflow of power into those modes through coupling. It is assumed, of course, that the critical angle θ_c is so large compared to Θ_∞ , that the steady-state distribution (26) is not significantly influenced by the boundary relations at $\theta = \theta_c$. If this is not the case, the solutions of (24) and (25) have to take these boundary relations into account.

IV. BUILD-UP FROM GAUSSIAN INPUT

Any z -dependent solution of (21) or (22) can of course be constructed from the infinite set of solutions of (24) and (25). But in the case of the fiber, there is a certain interest in special solutions which have an

arbitrary Gaussian input

$$P_{i_0} = P_0 \exp [-\theta^2/\Theta_0^2] \quad (30)$$

as initial condition. This is because it is convenient to study multimode fibers by using a Gaussian laser beam for excitation. A high-power lens converts this beam into the angular Gaussian distribution (30). By observing the change in $P(\theta)$ with fiber length, and the loss as a function of various (Gaussian) input distributions, one obtains valuable information on the power flow in the fiber.

Since both the input and the steady-state are Gaussian, it is reasonable to try the solution

$$P = f(z) \exp [-\theta^2/\Theta^2(z)]. \quad (31)$$

Although this approach is useful both for the two- and the three-dimensional configuration, we shall concentrate in the following on the fiber only. Introducing (31) into (22) yields the two differential equations

$$d\Theta/dz = -\frac{A}{2} \Theta^3 + 2D/\Theta \quad (32)$$

and

$$df/dz = -4Df/\Theta^2. \quad (33)$$

We can solve the first of these equations for Θ and obtain

$$\Theta^2 = \Theta_\infty^2 \left\{ \begin{array}{l} \tanh \gamma_\infty(z + z_0) \\ \coth \gamma_\infty(z + z_0) \end{array} \right\} \quad (34)$$

with the steady-state parameters Θ_∞ and γ_∞ from (27) and (29). The choice of \tanh or \coth and the coefficient z_0 are determined by the initial conditions. Eq. (33) can be solved with the help of (34) and yields

$$f = f_0 \left\{ \begin{array}{l} \sinh \gamma_\infty(z + z_0) \\ \cosh \gamma_\infty(z + z_0) \end{array} \right\}. \quad (35)$$

A Gaussian input stays indeed Gaussian, its width approaching monotonically that of the steady state. The transition function is the hyperbolic tangent if the input width is smaller than the steady-state width and the hyperbolic cotangent in the opposite case. For the initial conditions (30), the solutions (34) and (35) can be written in the form

$$\Theta^2(z) = \Theta_\infty^2 \frac{\Theta_0^2 + \Theta_\infty^2 \tanh \gamma_\infty z}{\Theta_\infty^2 + \Theta_0^2 \tanh \gamma_\infty z} \quad (36)$$

and

$$f(z) = \frac{P_o \theta_o^2}{\Theta_\infty^2 \sinh \gamma_\infty z + \Theta_o^2 \cosh \gamma_\infty z}. \quad (37)$$

To obtain the total power in the guide, $P(\theta)$ must be integrated over all angles up to θ_c . If we assume, as previously, that $P(\theta)$ is sufficiently small at the critical angle θ_c and beyond, we can extend the integration to infinity. With (31), the total power is

$$2\pi \int_0^\infty P(\theta) \theta d\theta = \pi f \Theta^2. \quad (38)$$

The power loss per unit length is consequently

$$\gamma(z) = -\frac{1}{f \Theta^2} \frac{d}{dz} (f \Theta^2). \quad (39)$$

By using the differentials $d\Theta/dz$ and df/dz from (32) and (33), we obtain

$$\gamma(z) = A \Theta^2(z). \quad (40)$$

With (27) and (29) this can also be written in the form

$$\frac{\gamma(z)}{\gamma_\infty} = \frac{\Theta^2(z)}{\Theta_\infty^2}. \quad (41)$$

The ratio (41) is plotted in Fig. 4 versus the fiber length for some specific input conditions. The plot illustrates the loss, the solid angle covered by the fiber output and, since this is proportional to the mode volume, also the number of modes propagating in the fiber.

If the measured width $\Theta(z)$ is small compared to the steady-state width Θ_∞ , we can approximate (36) by

$$\Theta^2 = \Theta_\infty^2 \gamma_\infty z + \Theta_o^2. \quad (42)$$

In this case, because of (27) and (29), $\Theta^2(z)$ is a straight line with the slope

$$\Theta_\infty^2 \gamma_\infty = 4D. \quad (43)$$

Thus measuring $\Theta^2(z)$ under these conditions yields directly the coupling parameter D . The approximate linear increase of the data in Fig. 3 is an indication of the validity of (42). A straight-line approximation of the measured data yields

$$D = 7 \cdot 10^{-5} \text{ rad}^2/\text{m}. \quad (44)$$

This value represents a first approximation for the zero-order coupling

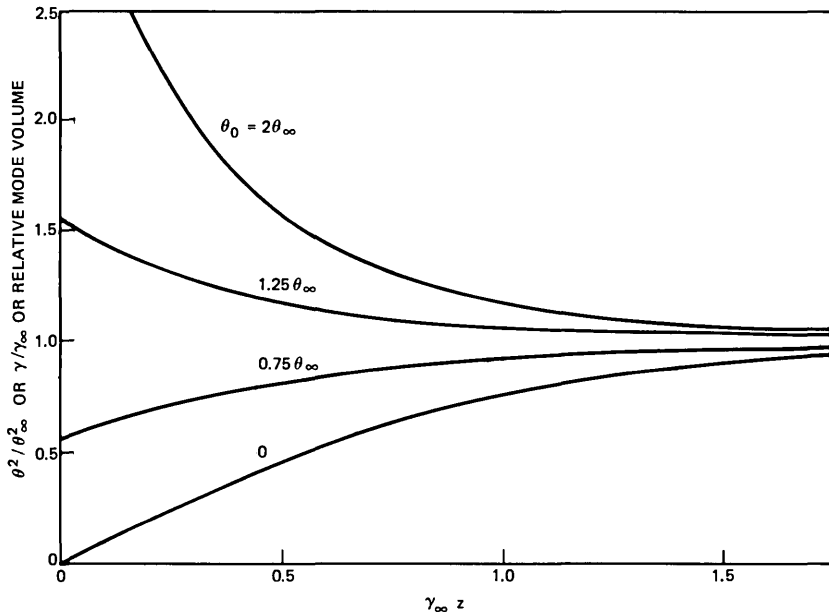


Fig. 4—Calculated increase of mode volume and loss with fiber length. The data are related to the steady-state values.

coefficient d_0 which, because of (20), becomes $d_0 = 16 \text{ m}^{-1}$ at the measuring conditions ($a = 50 \text{ } \mu\text{m}$, $\lambda = 0.63 \text{ } \mu\text{m}$). The range of the measured data is not sufficient to draw any conclusions on the size of higher-order coefficients in $d(\theta)$.

V. TWO SOURCES OF COUPLING

Marcuse¹⁰ has studied a dielectric slab guide with slightly distorted interfaces. He finds that two modes are coupled if the surface imperfections comprise a component of wavelength Λ that coincides with the "beat wavelength" between the two modes. The beat wavelength is the distance in which the phase difference between two modes increases to 2π . It can be calculated from the wave number

$$\beta_m = (k^2 n^2 - u^2)^{\frac{1}{2}} \approx kn - \frac{1}{2} \frac{u^2}{kn} \quad (45)$$

of the m th mode. With u from (3), we obtain for two neighboring modes of order m and $m + 1$,

$$\Lambda = \frac{2\pi}{\beta_m - \beta_{m+1}} = \frac{4a}{\theta}. \quad (46)$$

This is also the distance that a ray of angle θ requires to complete a zigzag period in a slab of width $2a$.

To describe the distorted slab walls, let us define a "power" spectrum $\phi(1/\Lambda)$ of the (random) deviations $\delta(z)$ from ideally straight interfaces. If the irregularities of both walls are uncorrelated, the coupling coefficient $d(\theta)$ between adjacent modes is¹⁰

$$d = \frac{1}{2} \left(\frac{nk}{a} \right)^2 \theta^4 \phi \left(\frac{\theta}{4a} \right). \quad (47)$$

Using this relation and (20), we can calculate the spectrum associated with the measured coupling parameter D of (44): The result

$$\phi \left(\frac{1}{\Lambda} \right) = \frac{D\Lambda^4}{32\pi^2} \quad (48)$$

suggests a decrease of ϕ with the fourth power of the (spatial) frequency $1/\Lambda$.

This result permits us to estimate the coupling among nonadjacent modes. As is evident from (45) and (46), modes which differ by a small number r have the beat wavelength $4ra/\theta$. Because of (48), coupling among such modes decreases with the fourth power of their order difference. It is this rapid decrease which permitted us to neglect all coupling except that between next-neighbors in (11). The error involved in this approximation can be estimated for the case that the power distribution $P(\theta)$ changes slowly within $r\Delta\theta$. In that case,

$$\frac{P_{m+r} - P_m}{\theta_{m+r} - \theta_m} \approx \frac{dP}{d\theta} \quad \text{for } r \ll m. \quad (49)$$

The transition from (11) to (15) then allows us to lump the coupling of all modes from m to $m+r$ in the coefficient $d_m = d(\theta)$ which assumes the form

$$d \sum_{r=1}^{\infty} \frac{1}{r^4} = \frac{\pi^4}{90} d \approx 1.08d. \quad (50)$$

This suggests that the relative error in our approximation is as small as 8 percent.

Random bends in the guide are another source of coupling. This problem has been studied by tracing rays through a randomly curved slab guide.¹¹ The result relates the statistics of the ray angle θ to the

“power spectrum” $C(1/\Lambda)$ of the curvature components. Reference 11 demonstrates that components of wavelength Λ predominantly influence rays with the same zigzag wavelength. For a short guide length z , the probability distribution of the ray angle θ is found to have a variance which increases as

$$\sigma^2(z) = \frac{8}{\pi^2} z C\left(\frac{\theta}{4a}\right) + \sigma_0^2 \quad (51)$$

where σ_0^2 is the variance at the input.¹¹ These results presuppose Gaussian statistics, for which the probability to find a ray at θ has the form $\exp(-\theta^2/2\sigma^2)$.

Let us compare this distribution to the distribution

$$P(\theta) = \exp[-\theta^2/\Theta^2]$$

of (31). For negligible mode attenuation, the width of this Gaussian is given by the simple relation (42). Like the variance of (51), it grows linearly with length. A comparison of the growth factors involved must take the factor 2 into account which enters because of the definition of the variance. This leads to the relation

$$C\left(\frac{1}{\Lambda}\right) = \frac{\pi^2}{4} D. \quad (52)$$

To compare this result with ϕ of (48), let us consider correlated deviations of the form $\delta \sin(2\pi z/\Lambda)$ at both walls. Twofold differentiation with respect to z transforms this into a curvature component of the form $(2\pi/\Lambda)^2 \delta \sin(2\pi z/\Lambda)$. Accordingly, we can relate the curvature spectrum to a spectrum of (correlated) irregularities

$$\phi_c\left(\frac{1}{\Lambda}\right) = (\Lambda/2\pi)^4 C. \quad (53)$$

This result transforms (52) into

$$\phi_c\left(\frac{1}{\Lambda}\right) = \frac{D\Lambda^4}{64\pi^2}. \quad (54)$$

The factor 2 which distinguishes this result from (48) results from the correlation of the wall deviations assumed in (53) contrary to (48);¹⁰ to represent a curved slab, the deviations must be equal and in phase.

Irregularities of this kind couple only modes which differ by an odd order number $r = 1, 3, \dots$. Coupling across even numbers results from irregularities in anti-phase or—in the case of the fiber—from irregularities of at least twofold cross-sectional symmetry. Specific parts of the spectrum $\phi(1/\Lambda)$ are likely to be dominated by certain kinds of irregularities.

The region of interest is determined by $\Lambda = 4a/\theta$ and was in our case between $\Lambda = 1$ and 10 mm. Since these lengths are significantly larger than the core diameter (150 μm), we believe that bends were the dominant source of coupling. The following example therefore uses the relation (52) for a quantitative estimate of the irregularities involved. Although (52) applies specifically to the slab model, we shall combine it with the fiber data (44), confident that this will illustrate at least the orders of magnitude involved.

To obtain a more tangible description of the random curvature, we assume it to be composed of randomly distributed singular deviations of the kind illustrated in Fig. 5a. We model these deviations by single-period sinewaves of the form $\delta \sin(2\pi z/\Lambda)$. Essentially, only those with a width larger than $\Lambda/2$ contribute to the curvature spectrum at Λ . If there are η of those per unit length, the curvature spectrum has approximately the value¹¹

$$C = 8\pi^4 \eta \delta^2 / \Lambda^2 \quad (55)$$

in the vicinity of Λ . Because of (52)

$$\eta \delta^2 = \frac{\Lambda^2 D}{32\pi^2}. \quad (56)$$

Fig. 5b evaluates this relation for the case of the fiber measured. Plotted is the density η versus the amplitude δ for ray periods (and angles) of interest. For example, an average of 1000 singular irregularities per meter would account for the coupling measured, if their magnitudes obeyed the relation $\delta = 15 \cdot 10^{-6} \Lambda$. That would mean that the magnitude of the irregularities increases linearly with their length reaching a value of 15 nm at $\Lambda = 1$ mm.

VI. CONCLUSIONS

A comprehensive description of a multimode fiber by one differential equation is possible, if the modes are approximated by a continuum. Under certain realistic conditions, this equation has a rigorous solution. We calculate here the far-field output distribution as a function of fiber length and compare this to experiments performed with a low-loss multimode fiber. We find neighboring modes to be coupled by an average 1.6 percent per mm. Coupling among other modes seems to be at least an order of magnitude less.

Among the possible sources of coupling taken into consideration, we believe random bends to be the most likely. In this case, the curvature

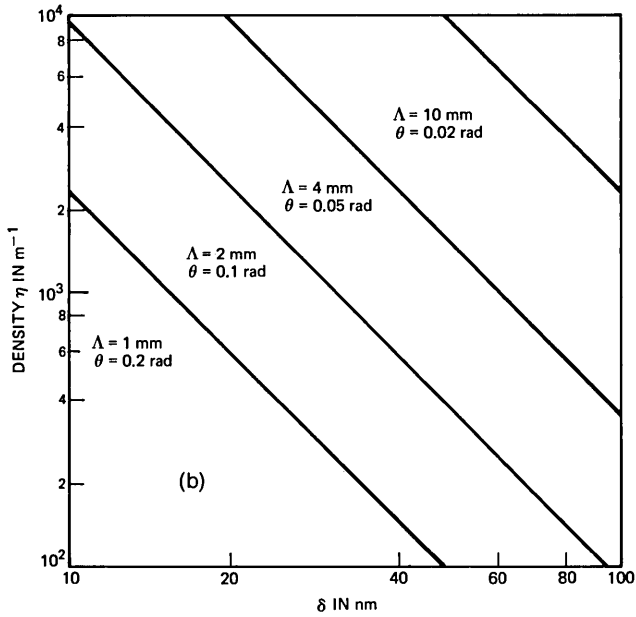
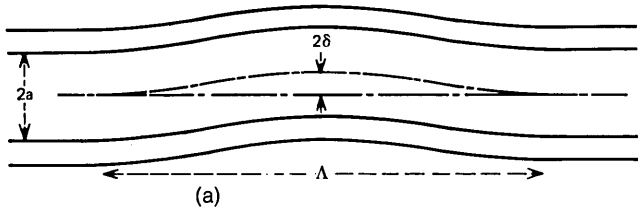


Fig. 5—Deviations from straightness as indicated in (a) account for the coupling measured, if their density η , and magnitude δ are related to their length Λ as plotted in (b).

spectrum has a value of $2 \cdot 10^{-7} \text{ mm}^{-1}$ in the region of spatial frequencies between 0.1 and 1 mm^{-1} . Another equivalent description of this result is by small hump-shaped deviations from straightness—on the average 1000 of them per meter—the magnitude of which increases proportionally to their length and is about 15 nm for a length of 1 mm.

The theory derived here can be modified to include the velocity differences among the modes and, in this way, to describe the impulse response in the presence of coupling and loss. This will be the subject of another paper.

VII. ACKNOWLEDGMENT

The experiment was performed together with E. L. Chinnock whose skillful cooperation is gratefully acknowledged.

REFERENCES

1. Stone, J., "Optical Transmission Loss in Liquid-Core Hollow Fibers," Topical Meeting on Integrated Optics—Guided Waves, Materials, and Devices, Las Vegas, Nevada, February 7–10, 1972.
2. Manufactured by Corning Glass Works, Corning, New York.
3. Gloge, D., et al., "Dispersion in a Low-Loss Multimode Fiber Measured at Three Wavelengths," Topical Meeting on Integrated Optics—Guided Waves, Materials and Devices, Las Vegas, Nevada, February 7–10, 1972.
4. Gloge, D., et al., "Picosecond Pulse Distortion in Optical Fibers," IEEE J. Quantum Elec., *QE-8* (1972), pp. 217–221.
5. Personick, S. D., "Time Dispersion in Dielectric Waveguides," B.S.T.J., *50*, No. 3 (March 1971), pp. 843–859.
6. Marcuse, D., "Pulse Propagation in Multimode Dielectric Waveguides," B.S.T.J., *51*, No. 6 (July–August 1972), pp. 1199–1232.
7. Gloge, D., Chinnock, E. L., and Stone, J., "Dispersion Measurements in Liquid-Core Fibers," unpublished work.
8. Gloge, D., "Weakly Guiding Fibers," Appl. Opt., *10*, 1971, pp. 2252–2258.
9. Marcuse, D., "Derivation of Coupled Power Equations," B.S.T.J., *51*, No. 1 (January 1972), pp. 229–237.
10. Marcuse, D., "Mode Conversion Caused by Surface Imperfections of a Dielectric Slab Waveguide," B.S.T.J., *48*, No. 10 (December 1969), pp. 3187–3215.
11. Gloge, D., "Bending Loss in Multimode Fibers With Graded and Ungraded Core Index," to appear in Applied Optics.

Pulse Propagation in a Two-Mode Waveguide

By D. MARCUSE

(Manuscript received May 2, 1972)

The results of an earlier paper, describing pulse propagation in multimode dielectric waveguides with random coupling, are specialized to the two-mode case. Because of their greater simplicity, the results for this special case provide more insight into the mechanism of pulse shortening due to mode coupling. The two-mode theory yields a formula for the width of a pulse carried by coupled guided modes that is found to hold also for four modes, so that it may be true for an arbitrary number of modes. This formula [eq. (23)] contains only the measurable distance required to establish the steady-state power distribution and the length of uncoupled pulses. The pulse length formula is identical with Personick's important result. Our treatment suggests that the characteristic length appearing in this formula may be accessible to measurement.

I. INTRODUCTION

In a series of earlier papers, the theory of multimode propagation in dielectric waveguides was analyzed with the help of stochastic coupled power equations.¹⁻³ The propagation of Gaussian-shaped pulses in a waveguide with randomly coupled modes was treated quite generally for N modes in Ref. 3. The theory was based on the following form of the stochastic coupled power equations:

$$\frac{\partial P_\nu}{\partial z} + \frac{1}{v_\nu} \frac{\partial P_\nu}{\partial t} = -\alpha_\nu P_\nu + \sum_{\mu=1}^N h_{\nu\mu} (P_\mu - P_\nu). \quad (1)$$

P_ν is the average power in mode ν , v_ν the group velocity, α_ν is the attenuation coefficient of mode ν in the absence of coupling to other guided modes, and $h_{\nu\mu}$ is the power coupling coefficient. To second order of perturbation theory, and assuming a Gaussian shape (in time) of the input pulse, the solution of (1) can be expressed as:³

$$P_\nu(z, t) = \sum_{i=1}^N \frac{2\tau}{\Delta t_i} k_i B_\nu^{(i)} e^{-\alpha_\nu^{(i)} z} \exp \left\{ -\left(\frac{t - z/v}{\Delta t_i/2} \right)^2 \right\}. \quad (2)$$

Δt_i , the width of the i th Gaussian function in (2), is given by

$$\Delta t_i = 2(\tau^2 + 4\alpha_2^{(i)}z)^{\frac{1}{2}}. \quad (3)$$

The input pulse with half width τ and amplitude G_v is assumed to be

$$P_v = G_v \exp\left(-\frac{t^2}{\tau^2}\right). \quad (4)$$

The coefficient k_i appearing in (2) is determined by the input pulse

$$k_i = \sum_{\nu=1}^N G_\nu B_\nu^{(i)}. \quad (5)$$

The vectors with components $B_\nu^{(i)}$ and the parameters $\alpha_o^{(i)}$ are the i th eigenvectors and eigenvalues of an algebraic eigenvalue problem defined in Refs. 2 and 3. The parameter $\alpha_2^{(i)}$ is the second-order perturbation of the eigenvalue $\alpha_o^{(i)}$ and is defined as follows:

$$\alpha_2^{(i)} = \sum_{\substack{j=1 \\ j \neq i}}^N \frac{\left\{ \sum_{\nu=1}^N \left(\frac{1}{v_\nu} - \frac{1}{v} \right) B_\nu^{(i)} B_\nu^{(j)} \right\}^2}{\alpha_o^j - \alpha_o^i}. \quad (6)$$

v is the average group velocity.

This approximate theory of pulse propagation in multimode waveguides holds for random coupling between the guided modes under the assumption that the correlation length of the coupling function is short compared to the distance over which the mode power P_v changes appreciably.

II. APPLICATION TO THE TWO-MODE CASE

In its full generality, the theory of multimode pulse operation is hard to evaluate. Computer solutions are being provided in Refs. 2 and 3. Here we want to derive expressions for the special case of two modes, that allows us to gain more insight into the meaning of the theory. The two-mode case has been treated previously by several authors.^{4,5} Our results are thus not all new.

For two modes, we write $h_{12} = h_{21} = h$. We now have to solve the eigenvalue problem^{2,3}

$$\left. \begin{aligned} (\alpha_o - \alpha_1 - h)B_1 + hB_2 &= 0 \\ hB_1 + (\alpha_o - \alpha_2 - h)B_2 &= 0 \end{aligned} \right\}. \quad (7)$$

The equation system (7) has the solution

$$\alpha_o^{(1)} = h + \frac{\alpha_1 + \alpha_2}{2} - \left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}} \tag{8}$$

$$\alpha_o^{(2)} = h + \frac{\alpha_1 + \alpha_2}{2} + \left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}} \tag{9}$$

for the first and second eigenvalue. The two components of the first eigenvector are

$$B_1^{(1)} = \frac{h}{\left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}} \left\{ \alpha_1 - \alpha_2 + 2 \left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}} \right\}^{\frac{1}{2}}} \tag{10}$$

$$B_2^{(1)} = \frac{\left\{ \alpha_1 - \alpha_2 + 2 \left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}} \right\}^{\frac{1}{2}}}{2 \left[\frac{(\alpha_2 - \alpha_1)^2}{4} + h^2 \right]^{\frac{1}{2}}} \tag{11}$$

The components of the second eigenvector can be expressed in terms of the components of the first eigenvector.

$$B_1^{(2)} = B_2^{(1)} \quad B_2^{(2)} = -B_1^{(1)}. \tag{12}$$

III. DISCUSSION OF THE TWO-MODE CASE

In the special case $\alpha_1 = \alpha_2 = 0$ we have the eigenvalues

$$\alpha_o^{(1)} = 0 \tag{13}$$

and

$$\alpha_o^{(2)} = 2h, \tag{14}$$

while the eigenvectors are

$$B_1^{(1)} = \frac{1}{\sqrt{2}} \quad B_2^{(1)} = \frac{1}{\sqrt{2}} \tag{15}$$

and

$$B_1^{(2)} = \frac{1}{\sqrt{2}} \quad B_2^{(2)} = -\frac{1}{\sqrt{2}}. \tag{16}$$

There are several interesting features apparent in this special solution. In the absence of loss, both modes carry equal power. We see immediately from (15) and (16) that the sum of the squares of the components of each eigenvector adds up to unity, while the inner product of the two vectors vanishes. This is a general property that is also shared by the solutions (10) through (12).

The "loss coefficient" $\alpha_o^{(1)}$ of the term with $i = 1$ in (2) vanishes, so that this term does not decrease in amplitude as the pulse moves along the z -axis. The "loss coefficient" $\alpha_o^{(2)}$ of the second term of (2) (with $i = 2$) is equal to twice the coupling coefficient h , so that this term becomes vanishingly small for large values of z . This too is a general feature of (2). The lowest order eigenvalue $\alpha_o^{(1)}$ is smaller than all other eigenvalues, so that only the first term of the series (2) remains for large values of z while all the other terms have become vanishingly small. Even though $\alpha_o^{(1)}$ is not zero in the general (lossy) case, the multimode waveguide always reaches a steady state which is described by the first term of the series expansion in (2), provided that the modes are coupled. Only in the lossless case do we find $\alpha_o^{(1)} = 0$. It is noteworthy that the "loss terms" $\exp(-\alpha_o^{(i)}z)$ ($i > 1$) all become vanishingly small even in the absence of losses. A steady-state distribution of mode power versus mode number is thus established that is independent of the initial excitation of the waveguide. The decay of the higher-order terms in (2) does not necessarily indicate power loss. We see indeed, from the solution (16), that the sum of the components of the second eigenvector adds up to zero indicating that no power is carried by the second term ($i = 2$) of (2) in the absence of loss. The individual terms of (2) must not be confused with waveguide modes. They have no independent physical meaning except for the first term with $i = 1$, which is the steady-state power distribution. It is apparent that it is sufficient to study the behavior of the first term in (2) alone, since all other terms (the second term is the only other term in the two-mode case) become negligible for large values of z .

We can easily define a characteristic distance that is required for the steady state to establish itself. Once the exponential factor $\exp(\alpha_o^{(1)} - \alpha_o^{(2)})z$ has become small, the steady state is reached. We thus define the characteristic length as follows

$$L_s = \frac{\kappa}{\alpha_o^{(2)} - \alpha_o^{(1)}}. \quad (17)$$

The parameter κ is a number of order unity. For $\kappa = 1$ we have $\exp(\alpha_o^{(1)} - \alpha_o^{(2)})L_s = 1/e$. Thus $\kappa = 1$ is too small to consider the steady state as reached. However, we can still define L_s by (17) with $\kappa = 1$. If we use $\kappa = 4.6$, we have $\exp(\alpha_o^{(1)} - \alpha_o^{(2)})L_s = 0.01$. This number is small enough to consider the second term in (2) as negligibly small. For the two-mode case we obtain from (14) and (17) for the case of low losses

$$L_s = \frac{\kappa}{2h}. \quad (18)$$

Finally, we study the steady-state pulse width which follows from (3), with $i = 1$. We also neglect τ in this equation, assuming that the pulse has spread to a size much longer than the input pulse. From (6), (8), (9), (10), and (11) we find

$$\Delta t_1 = 4\Delta T \frac{h}{[(\alpha_2 - \alpha_1)^2 + 4h^2]^{\frac{3}{4}}} \frac{1}{\sqrt{L}}. \quad (19)$$

We used $z = L$, with L designating the length of the waveguide. The factor ΔT , the width of the pulse in the absence of coupling, is defined as

$$\Delta T = \left(\frac{1}{v_2} - \frac{1}{v_1} \right) L. \quad (20)$$

Compared to the width ΔT of the uncoupled modes, the pulse length in case of coupled modes is improving with length. The pulse width formula can again be considered in the two limiting cases. If $\alpha_2 - \alpha_1 \ll h$ we have

$$\Delta t_1 = \sqrt{2} \frac{\Delta T}{(hL)^{\frac{3}{4}}}. \quad (21)$$

This formula shows clearly that the pulse length shortens with increased coupling strength. In the other extreme, $h \ll |\alpha_2 - \alpha_1|$, we have

$$\Delta t_1 = 4\Delta T \frac{h}{[(\alpha_2 - \alpha_1)^3 L]^{\frac{3}{4}}}. \quad (22)$$

It appears strange at first that the pulse length now increases as the coupling strength is increased. However, in this mode of operation, the pulse length is primarily determined by the differential loss of the two modes. If both modes travel uncoupled, one will die out while the other carries the pulse all by itself. In this case, the pulse width is determined only by the dispersion of the surviving mode, which is not included in our theory. For $h = 0$, we thus obtain a vanishing pulse length. As the coupling is increased, power is flowing from the lower loss mode to the high loss mode so that the pulse width is increased by the different delay time of each mode. It is thus clear that a small amount of coupling causes the pulse to lengthen.

Finally, we combine the formula (21) for the low loss case with the formula (18) defining the characteristic length that indicates where the steady state is reached. We thus obtain the interesting result

$$\Delta t_1 = \frac{2}{\sqrt{\kappa}} \Delta T \left(\frac{L_s}{L} \right)^{\frac{1}{2}}. \quad (23)$$

The factor in front of this equation becomes unity if we use $\kappa = 4$. We have seen that this value is large enough to ensure that steady state is essentially reached. Equation (23) has been derived by Personick.⁴ We see from our derivation that the characteristic length l_c in Personick's formula can be interpreted as the length that is required to reach the steady-state power distribution.

Equation (23) was derived for the case of only two modes. It is tempting to use this equation also for the multimode case. In order to test the mode dependence of this formula, I solved the four-mode problem under the assumption that all off-diagonal elements of $h_{\nu\mu}$ vanish with the exception of the elements directly adjacent to the main diagonal. All non-vanishing elements of $h_{\nu\mu}$ were set equal to the same value h . For the lossless case, and assuming that the inverses of the group velocities are evenly spaced, the following formula was obtained for the four-mode case.

$$\Delta t_1 = 0.79 \frac{2}{\sqrt{\kappa}} \Delta T \left(\frac{L_s}{L} \right)^{\frac{1}{2}}. \quad (24)$$

L_s is again defined as the length required to achieve the steady state. ΔT is the length of the uncoupled signal for the four-mode case. Since (23) and (24) are essentially identical [(24) is even slightly more favorable], it might be assumed that (23) may hold independently of mode number.

Equation (23) was derived for the case where the coupling coefficient h is larger than the loss coefficients α_1 or α_2 . In the opposite case, where the losses determine the rate at which the power distribution approaches the steady state, no simple relationship exists between L_s and Δt . Formula (23) is quite useful for estimating the length of the Gaussian pulse if the distance L_s , at which steady state is reached, can be observed. If it is known that radiation losses are small compared to the coupling coefficient h , the conditions exist for which (23) was derived. However, it may well be that the region of dominance of radiation losses over coupling strength is different for different modes. Experiments have shown that the coupling mechanism in optical fibers consists of two parts.⁶ A Rayleigh-type background with a wide mechanical Fourier spectrum is responsible for most of the radiation losses, while a very sharp peak at zero frequencies of the mechanical power spectrum is responsible for most of the coupling between guided modes. The

broad spectrum causes more radiation loss for higher-order modes, provided that the coupling is caused by core-cladding interface irregularities. The narrow peak at zero mechanical frequencies couples lower-order modes much more strongly than high-order modes, because of the closer spacing (in β -space) of the low-order modes.³ The combined effect of these two spectral regions causes a steady-state distribution that favors the lower-order modes. In this situation, it may still be possible to estimate the pulse performance of the multimode waveguide by ignoring those modes that do not carry power in the steady-state distribution and interpret ΔT as the pulse length that would be obtained by the remaining modes in the absence of coupling. The steady-state distance L_s is best observed by launching only the lowest-order modes and measuring the distance that is required until the power versus mode number distribution ceases to change its shape.

REFERENCES

1. Marcuse, D., "Derivation of Coupled Power Equations," B.S.T.J., 51, No. 1 (January 1972), pp. 229-237.
2. Marcuse, D., "Power Distribution and Radiation Losses in Multimode Dielectric Waveguides," B.S.T.J., 51, No. 2 (February 1972), pp. 429-454.
3. Marcuse, D., "Pulse Propagation in Multimode Dielectric Waveguides," B.S.T.J., 51, No. 6 (July-August 1972), pp. 1199-1232.
4. Personick, S. D., "Time Dispersion in Dielectric Waveguides," B.S.T.J., 50, No. 3 (March 1971), pp. 843-859.
5. Rowe, H. E., and Young, D. T., "Transmission Distortion in Multimode Random Waveguides," IEEE Trans. MTT, MTT-20, No. 6 (June 1972), pp. 349-365.
6. Rawson, E. G., "Measurements of the Angular Distribution of Light Scattered from a Glass Fiber Optical Waveguide," to be published in Applied Optics, 11, No. 11 (November 1972).

Fluctuations of the Power of Coupled Modes

By D. MARCUSE

(Manuscript received May 2, 1972)

Using perturbation theory, an expression for the variance of the power of each mode of a multimode waveguide with randomly-coupled modes is derived. The variance builds up from zero to a constant value as a function of z (length along the waveguide). For most cases of interest, the variance is equal to the square of the average power. This means that the power of each mode of a system of randomly-coupled modes of a multimode waveguide fluctuates like the short-term time averaged power of a narrowband electrical signal the voltage of which is a random variable with Gaussian probability distribution.

I. INTRODUCTION

The behaviour of waves propagating in multimode waveguides can be described by coupled equations for the amplitudes of each mode.¹ This description is rigorous, but has the disadvantage that the coupled wave equations usually cannot be solved. It has been shown that a much simpler description is possible if we limit our interest to knowledge about the average power carried by each mode.²⁻⁴ Coupled equations for the average mode power have been derived and applied to the problem of wave propagation in multimode dielectric waveguides.^{4,5} However, the description of multimode waveguides in terms of average power is incomplete unless some information is available about the fluctuations of the actual power about the average value. With the help of the same perturbation approach that was used to derive the coupled power equations,⁴ we derive in this paper a differential equation for the variance of the power.

The result of our perturbation theory is expressed in terms of the cross-correlation and the average power of the modes. In order to evaluate this expression, we need to make several assumptions. It has been shown in an earlier paper⁵ that the average power settles down to a steady-state distribution of power versus mode number that is

independent of the initial excitation of the waveguide. Using this concept of the steady-state distribution and the further assumption that the cross-correlation between the modes is small, we can solve the differential equation for the variance. We find that the variance builds up from zero values at $z = 0$, to a constant value which is equal to the square of the average of the mode power. The relative fluctuation of the power of each mode is thus 100 percent. This means that the power in each of the randomly-coupled modes behaves like the short-term time-averaged power of a narrowband electrical signal the voltage of which is a Gaussian random variable.

II. DERIVATION OF THE DIFFERENTIAL EQUATION FOR THE VARIANCE

Our starting point is the set of coupled wave equations for the slowly varying wave amplitudes (envelops) A which are defined by

$$a_\nu = A_\nu e^{-i\beta_\nu z} \quad (1)$$

with a_ν being the rapidly oscillating mode amplitude. The coupled wave equations can be expressed in the form⁴

$$\frac{dA_\nu}{dz} = \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^N c_{\nu\mu} A_\mu e^{i(\beta_\nu - \beta_\mu)(z - z')} \quad (2)$$

z' is used as a convenient reference point. The parameters β_ν are the propagation constants of the modes. The coupling coefficient can be expressed as a product of a constant term times a function of z .

$$c_{\nu\mu} = K_{\nu\mu} f(z). \quad (3)$$

If we define

$$K_{\nu\nu} = 0, \quad (4)$$

we can drop the restriction $\mu \neq \nu$ in (2). Conservation of power leads to the relation⁴ (the asterisk indicates complex conjugation)

$$K_{\mu\nu} = -K_{\nu\mu}^*. \quad (5)$$

The perturbation theory uses the approximate solution of (2)

$$A_\nu(z) = A_\nu(z') + \sum_{\mu=1}^N K_{\nu\mu} A_\mu(z') \int_{z'}^z f(x) e^{i(\beta_\nu - \beta_\mu)x} dx. \quad (6)$$

The power of mode ν is

$$P_\nu(z) = |a_\nu|^2 = e^{-\alpha_\nu(z - z')} |A_\nu|^2. \quad (7)$$

$\alpha_\nu = -2\text{Im}\beta_\nu$ is the power attenuation coefficient of mode ν in the absence of coupling. Throughout our derivation we assume that the losses are so slight that we can approximate $\exp(-\alpha_\nu(z - z'))$ by unity.

We use the loss term in (7) only to modify our equations for the lossy case.

The variance of the power of mode ν is defined as

$$(\Delta P_\nu)^2 = \langle P_\nu^2 \rangle - \hat{P}_\nu^2, \tag{8}$$

with the simplified notation

$$\hat{P}_\nu = \langle P_\nu \rangle. \tag{9}$$

The derivative of the variance can be written with the help of (7) (replacing the exponential term by unity)

$$\frac{1}{2} \frac{d}{dz} (\Delta P_\nu)^2 = \left\{ \left\langle A_\nu A_\nu^* \frac{dA_\nu}{dz} A_\nu^* \right\rangle + \text{c.c.} \right\} - \hat{P}_\nu \frac{d\hat{P}_\nu}{dz}. \tag{10}$$

The expression c.c. indicates that the complex conjugate of the first term in the bracket must be added. The derivative of the average power has already been evaluated so that we do not need to express it in terms of the wave amplitudes. With the help of (2), (10) can be written as follows:

$$\frac{1}{2} \frac{d}{dz} (\Delta P_\nu)^2 = \left\{ \sum_\mu \langle A_\nu A_\nu^* A_\mu A_\mu^* f(z) \rangle K_{\nu\mu} e^{i(\beta_\nu - \beta_\mu)z} + \text{c.c.} \right\} - \hat{P}_\nu \frac{d\hat{P}_\nu}{dz}. \tag{11}$$

We now follow the technique that was developed in Ref. 4. We replace all the amplitudes in (11) with the approximate solution (6), but keep only terms up to second order in $K_{\nu\mu}$. The first-order terms vanish if we assume that $f(z)$ is statistically independent of $A(z')$. This assumption is justified if we let $z - z'$ be much larger than the correlation length of $f(z)$. For the same reason, we write the ensemble average of products of the field amplitudes with terms containing $f(z)$ as a product of an ensemble average containing only amplitude terms, times an ensemble average of a term that contains only $f(z)$. We thus obtain

$$\begin{aligned} & \frac{1}{2} \frac{d}{dz} (\Delta P_\nu)^2 \\ &= \left\{ \sum_{\mu, \delta} \left[K_{\nu\mu} K_{\nu\delta} \langle A_\delta A_\nu^* A_\mu A_\nu^* \rangle e^{i(2\beta_\nu - \beta_\mu - \beta_\delta)z} \right. \right. \\ & \quad \cdot \int_{z'}^z \langle f(z)f(x) \rangle e^{i(\beta_\nu - \beta_\delta)(x-z)} dx + 2K_{\nu\mu} K_{\nu\delta}^* \langle A_\nu A_\delta^* A_\mu A_\nu^* \rangle e^{i(\beta_\delta - \beta_\mu)z} \\ & \quad \cdot \int_{z'}^z \langle f(z)f(x) \rangle e^{-i(\beta_\nu - \beta_\delta)(x-z)} dx + K_{\nu\mu} K_{\mu\delta} \langle A_\nu A_\nu^* A_\delta A_\nu^* \rangle e^{i(\beta_\nu - \beta_\delta)z} \\ & \quad \left. \cdot \int_{z'}^z \langle f(z)f(x) \rangle e^{i(\beta_\mu - \beta_\delta)(x-z)} dx \right] + \text{c.c.} \left. \right\} - \hat{P}_\nu \sum_\mu h_{\nu\mu} (\hat{P}_\mu - \hat{P}_\nu). \tag{12} \end{aligned}$$

The amplitudes A_i are understood to have the argument z' . The last term was obtained by using the lossless coupled equations for the average power derived in Ref. 4. The power coupling coefficient is defined as⁴

$$h_{\nu\mu} = |K_{\nu\mu}|^2 F(\beta_\nu - \beta_\mu). \quad (13)$$

The power spectrum $F(\beta_\nu - \beta_\mu)$ is the ensemble average of the absolute square value of the Fourier coefficient of $f(z)$.

The next step in the derivation is based on the realization that only nonoscillatory terms contribute appreciably to the growth of the variance as a function of z . We thus neglect all but the nonoscillating terms in (12). This procedure is reinforced by the fact that the ensemble averages of cross terms of amplitudes are likely to be smaller than the ensemble averages of absolute squares of the amplitudes. The first term in (12) causes some concern since it appears that there may be several combinations of μ and δ in addition to $\mu = \delta = \nu$ that contribute nonoscillatory terms. However, the uneven spacing of the modes along the β axis makes it appear unlikely that combinations of modes can be found for which the exponent of the exponential function in front of the integral vanishes. However, even if a few such combinations could be found, we could still consider the term belonging to such combinations as small because of the lack of correlation between the amplitude coefficients belonging to different modes. Since $K_{\mu\mu} = 0$, we find that the first term in (12) does not contribute appreciably to the derivative of the variance and can be neglected. The integrals can be expressed in terms of the power spectrum of the function $f(z)$ as was shown in Ref. 4. We thus obtain, with the help of (5) and (13), and dropping all oscillatory terms

$$\frac{d}{dz} (\Delta P_\nu)^2 = 2 \sum_{\mu=1}^N h_{\nu\mu} [2\langle P_\nu P_\mu \rangle - \hat{P}_\nu \hat{P}_\mu - (\Delta P_\nu)^2]. \quad (14)$$

Assuming that $P_\nu(z') \approx P_\nu(z)$, we use z as the argument of P_ν . Finally, we introduce losses into the theory. The example of the coupled power equations serves well to illustrate the procedure. Neglecting losses we obtain⁴

$$\frac{d\hat{P}_\nu}{dz} = \sum_{\mu=1}^N h_{\nu\mu} (\hat{P}_\mu - \hat{P}_\nu). \quad (15)$$

According to our derivation, we have used the approximation $\hat{P}_\nu = \langle |A_\nu|^2 \rangle$. Using (7), we obtain

$$\frac{d\langle |A_\nu|^2 \rangle}{dz} = \left(\alpha_\nu \hat{P}_\nu + \frac{d\hat{P}_\nu}{dz} \right) e^{\alpha_\nu(z-z')}. \quad (16)$$

The left-hand sides of (15) and (16) are identical according to our derivation. By substituting (16) and assuming again that $\alpha(z - z') \ll 1$, we obtain

$$\alpha_\nu \hat{P}_\nu + \frac{d\hat{P}_\nu}{dz} = \sum_{\mu=1}^N h_{\nu\mu} (\hat{P}_\mu - \hat{P}_\nu). \quad (17)$$

Equation (17) is identical with equation (29) of Ref. 4. There, we introduced the loss simply as a phenomenological parameter. Our present treatment shows how the loss term can be obtained directly from the derivation based on perturbation theory. By applying the same reasoning to (14), we obtain

$$\frac{d}{dz} (\Delta P_\nu)^2 = -\kappa_\nu (\Delta P_\nu)^2 + 2 \sum_{\mu=1}^N h_{\nu\mu} [2\langle (P_\nu - \hat{P}_\nu)(P_\mu - \hat{P}_\mu) \rangle + \hat{P}_\nu \hat{P}_\mu]. \quad (18)$$

The parameter κ_ν is defined as

$$\kappa_\nu = 2\alpha_\nu + 2 \sum_{\mu=1}^N h_{\nu\mu}. \quad (19)$$

We regrouped the terms under the summation sign in (18) in order to express $\langle P_\nu P_\mu \rangle$ in terms of the cross correlation $\langle (P_\nu - \hat{P}_\nu)(P_\mu - \hat{P}_\mu) \rangle$. Integration of (18) yields, finally, the desired expression for the variance of the mode power

$$\begin{aligned} (\Delta P_\nu)^2 &= (\Delta P_\nu)_{z=0}^2 + 2e^{-\kappa_\nu z} \\ &\quad \cdot \int_0^z e^{\kappa_\nu x} \sum_{\mu=1}^N h_{\nu\mu} [2\langle (P_\nu - \hat{P}_\nu)(P_\mu - \hat{P}_\mu) \rangle + \hat{P}_\nu \hat{P}_\mu] dx. \end{aligned} \quad (20)$$

Equation (20) is the solution of the variance problem. The expression in brackets under the summation sign can be positive or negative, so that the variance can increase or decrease with increasing z .

III. EVALUATION OF THE VARIANCE FOR SPECIAL CASES

In order to be able to evaluate the general expression (20) for the variance, we would need to know the cross correlation and the average power as functions of z . The average power can be obtained by solving the coupled power equations. However, the cross correlation is not known. It appears reasonable to assume that the cross correlation may be small in many cases of practical interest. One would not expect to obtain small values of the cross correlation for only two modes because as one mode gains power the other must lose an equal amount of power.

However, for large numbers of modes, it appears reasonable to expect that the cross correlation between different modes may be small.

It is known from the theory of coupled power equations that the distribution of power versus mode number settles down to a steady state.⁵ Once the steady state is reached, each mode decays with the same attenuation coefficient. The shape of the distribution of average power versus mode number remains unchanged, but its level decreases exponentially with a power attenuation constant α_s . If we launch a power distribution at $z = 0$ that corresponds to the steady-state distribution, we obtain power averages that do not change with z except for a common exponential decay term. Assuming, therefore, that the cross correlation is negligible and that the steady-state power distribution is launched into the guide, allows us to solve (20) immediately. Using $\hat{P}_\nu = \hat{P}_{\nu 0} e^{-\alpha_s z}$ we obtain for $(\Delta P_\nu)_{z=0}^2 = 0$

$$\frac{\Delta P_\nu}{\hat{P}_\nu(z)} = \left\{ 2 \frac{1 - e^{-(\kappa_\nu - 2\alpha_s)z}}{\kappa_\nu - 2\alpha_s} \sum_{\mu=1}^N h_{\nu\mu} \frac{\hat{P}_{\mu 0}}{\hat{P}_{\nu 0}} \right\}^{\frac{1}{2}} \quad (21)$$

with κ_ν given by (19). Equation (21) represents the relative fluctuation of the power of mode ν . It shows clearly that the relative fluctuations build up from zero to a constant value which is reached when the z -dependent exponential function in (21) becomes negligibly small. The shape of the steady-state power distribution depends on the interplay between the coupling between the guided modes and the loss of power to radiation. The loss coefficient α_s that appears in the coupled power equations (17), depends on the mode number. Usually higher-order modes suffer more losses than lower-order modes. Modes with a large loss coefficient carry only little power once the steady-state power distribution is reached. Modes with small average power are of little interest. Concentrating on those modes that carry appreciable amounts of power allows us to neglect the attenuation coefficient α_s that appears implicitly in (21) through relation (19). If all the guided modes couple strongly to each other, they are also strongly coupled to the radiation field, thus losing a large amount of power by radiation. Since reasonably low loss operation is of most interest in practical applications, we can limit our discussion to the situation where only neighboring guided modes are coupled to each other. This means that $h_{\nu\mu}$ is small for large values of $|\mu - \nu|$. The sum over $h_{\nu\mu}$ thus extends only over those values of μ which are close to ν . The steady-state power distribution is continuous in the sense that neighboring modes carry nearly equal amounts of power. Neglecting the small steady-state loss coefficient α_s compared to the sum over the coupling coefficients $h_{\nu\mu}$, and using the fact that

the power of neighboring modes is nearly equal, allows us to obtain for $z \rightarrow \infty$ from (21) the important relation

$$\frac{\Delta P_\nu}{\hat{P}_\nu} \approx 1. \quad (22)$$

IV. DISCUSSION OF THE RESULT

The relative fluctuation of the power of those modes that carry appreciable amounts of power is approximately 100 percent. Such fluctuations are not unusual, however. The short-term time-averaged power carried by a narrowband electrical signal, the voltage of which is a Gaussian random variable, is known to fluctuate in the same way. The probability distribution for P_ν can, in analogy to the electrical case, be assumed to be

$$W(P_\nu) = \frac{1}{\hat{P}_\nu} \exp\left(-\frac{P_\nu}{\hat{P}_\nu}\right). \quad (23)$$

From this analogy we can immediately state that the relative fluctuations of the power of M modes (assumed to be uncorrelated) is equal to $M^{-1/2}$.

The fluctuations that we are considering do not occur in time at the output of any given waveguide. They are fluctuations of random variables in an ensemble sense. If we were to measure the power in a given mode for each of a large number of similar waveguides, we would expect to obtain results that fluctuate according to (22). Equation (22) thus tells us the accuracy of predicting the value of the power in a given mode on the basis of the coupled power equations. Since it is very hard to measure the power carried by one individual mode of a multimode waveguide, we are more likely to observe the power P_M in a fairly large number of M modes simultaneously. In this case, we expect to obtain fluctuations according to the law

$$\frac{\Delta P_M}{\hat{P}_M} = \frac{1}{\sqrt{M}}. \quad (24)$$

It is helpful to remember that the power of all N modes does not fluctuate at all.

V. CONCLUSIONS

We have discussed the problem of the relative fluctuations of the power in individual modes of a multimode waveguide in the case that

the modes are coupled by a random coupling function. Our discussion was limited to the c.w. case and does not directly apply to pulsed operation. We derived the general expression (20) for the variance of the power in terms of the cross correlation and the average power carried by the modes. Under the assumption that the modes are approximately uncorrelated among each other, and assuming further that only neighboring modes are coupled, we found that the relative fluctuations are nearly 100 percent. This result is reminiscent of the fluctuations of the short term time averaged power of a narrowband electrical signal the noise voltage of which is a Gaussian random variable.

Cross correlation between the modes can either increase or decrease the fluctuations depending on the sign of the cross correlation term in (20). It is reasonable to assume that the sign would tend to be negative. As stated earlier, coupling between only neighboring modes is necessary for low loss operation. If mode ν should, at a given point on the z axis, carry more than the average amount of power, we conclude that this power has been transferred from the neighboring mode (or modes) μ so that this mode is expected to have less than the average amount of power. The sign of the two factors in the cross correlation term must thus be different so that the term assumes a negative sign. This qualitative discussion indicates that correlations among the modes would tend to reduce the variance $(\Delta P_\nu)^2$. I have observed fluctuations of the mode power as large as those predicted by this theory in numerical solutions of coupled line equations with random, band-limited coupling function. This "experimental" result confirms the assumption that cross correlation between modes does not appreciably reduce the variance of the power fluctuations.

REFERENCES

1. Miller, S. E., "Coupled Wave Theory and Waveguide Applications," B.S.T.J., 33, No. 3 (May 1954), pp. 661-719.
2. Rowe, H. E., and Young, D. T., "Transmission Distortion in Multimode Random Waveguides," IEEE Trans. MTT, MTT-20, No. 6 (June 1972), pp. 349-365.
3. Morrison, J. A., and McKenna, J., "Coupled Line Equations with Random Coupling," B.S.T.J., 51, No. 1 (January 1972), pp. 209-228.
4. Marcuse, D., "Derivation of Coupled Power Equations," B.S.T.J., 51, No. 1 (January 1972), pp. 229-237.
5. Marcuse, D., "Power Distribution and Radiation Losses in Multimode Dielectric Slab Waveguides," B.S.T.J., 51, No. 2 (February 1972), pp. 429-454.

Higher-Order Scattering Losses in Dielectric Waveguides

By D. MARCUSE

(Manuscript received April 17, 1972)

This paper discusses the scattering losses of dielectric slab waveguides that are caused by higher-order grating lobes of the sinusoidally distorted core-cladding interface. The results of this paper are used in a companion paper to evaluate the radiation losses of multimode guides with intentional mode coupling. An exact system of equations is derived for the amplitudes of all grating orders. This system is used to derive first- and second-order approximations that hold for small amplitudes of the sinusoidal interface distortion. The theory is used to derive formulas for the average power loss coefficient for first- and second-order scattering processes.

I. INTRODUCTION

Scattering losses in dielectric waveguide caused by core-cladding interface irregularities have been studied extensively by means of first-order perturbation theory. The principle result of this theory can be stated as follows:¹ Two modes with propagation constants β_v and β_μ are coupled only if a Fourier component of the core-cladding interface function exists the mechanical frequency of which, ϕ , satisfies the relation

$$\phi = |\beta_v - \beta_\mu|. \quad (1)$$

The propagation constants β_v and β_μ may both belong to guided modes or one may belong to a guided mode while the other belongs to the continuum of radiation modes. Coupling of a guided mode to radiation modes results in power loss of the guided mode. This first-order coupling process is very strong and leads to high radiation losses if suitable Fourier components of the mechanical core-cladding interface irregularity function exist.

The result of first-order perturbation theory can be understood by viewing the core-cladding interface as a diffraction grating.² Since it

modulates the phase of the incident wave passing through it, the dielectric interface acts as a phase grating. The guided modes of dielectric waveguides can be decomposed into plane waves.³ This decomposition is particularly simple in the case of the modes of a slab waveguide. The following discussion is thus applied to this structure. Two plane waves are superimposed to form a traveling wave in the z direction—the direction of the waveguide axis—and a standing wave in the direction transverse to the z axis. The coupling coefficients for guided mode coupling and the radiation loss coefficients can be calculated by solving the plane wave scattering problem at the dielectric interface.⁴ The geometry of the problem is shown in Fig. 1. For clarity of discussion, it was assumed that the incident plane wave approaches the interface at right angles. The actual plane waves making up the guided mode of the slab waveguide approach the interfaces at grazing angles. Figure 1 is drawn with a sinusoidally distorted core-cladding interface. In this case, the incident plane wave decomposes into a wave that continues to travel in the original direction after passing the interface and into a reflected wave plus a series of side lobes that are labeled by positive

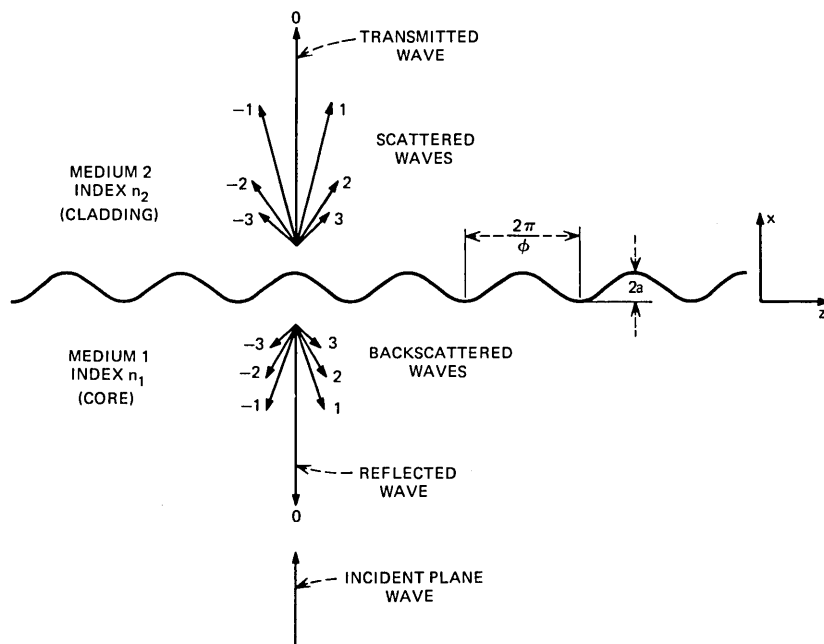


Fig. 1—A sinusoidally deformed dielectric interface functions as a phase grating. The figure shows the propagation vector of an incident plane wave (labeled 0) and reflected as well as transmitted plane waves of the higher-order grating lobes.

and negative integer numbers. These are the lobes of first, second, and higher order of the phase grating. If the incident plane wave make as grazing angle with the interface, the situation is essentially the same; the only difference being that some of the grating lobes form imaginary angles and are consequently evanescent instead of traveling waves. The zero-order lobe does not pass if the incident wave meets the interface at less than the critical angle for total internal reflection. Its angle is thus imaginary and only an evanescent wave exists on the far side of the interface while the incident wave is strongly reflected. The transmitted first-order grating lobes may both have imaginary angles. In this case, there are only scattered reflected waves in the core of the waveguide that can combine to a new guided mode provided that their angles correspond to one of the allowed directions for guided modes. If one of the transmitted first-order side lobes emerges on the far side of the interface with a real angle, it causes power to radiate into the space outside of the waveguide core. This radiation is lost to the guided wave and must be counted as power loss.

For small amplitudes of the sinusoidal core-cladding interface irregularity, the amplitudes of the grating lobes decrease rapidly with increasing grating order. The first-order grating lobe corresponds to the contribution of first-order perturbation theory and is indeed the only contribution of real interest if the core-cladding interface irregularity has a small amplitude. With increasing amplitude of the interface irregularity, the higher-order grating lobes become increasingly important. Their importance is enhanced by the fact that the angles of the second-order grating lobes may be real even when the transmitted first-order lobes both have only imaginary angles. This means that, to first-order of perturbation theory, no scattering loss exists. It is thus necessary to study the higher order grating responses in order to obtain information about scattering losses in case the first-order theory predicts no scattering loss at all.

The study of these higher-order grating lobes and the derivation of power loss coefficients for the higher-order processes is the object of this paper. It is possible that the grating problem has been solved before to the accuracy that is attempted here. Because of the enormous volume of literature that exists on scattering problems, relevant papers may have escaped the author's attention. However, the application of the grating theory to the waveguide loss problem is probably new.

We derive coupled equation systems for the amplitudes of the grating lobes and use these exact equations to obtain first- and second-order approximations in a straightforward way. Perturbation solutions of the exact equation system are particularly appropriate, since each

higher-order solution can be computed from the known lower-order solutions with no need to recompute the lower-order solutions each time the order of perturbation theory is increased by one. Formulae for the first- and second-order scattering loss process from a sinusoidal core cladding interface irregularity are derived. The results of this paper will be used in a companion paper⁵ to calculate the loss penalty for intentional mode mixing in multimode waveguides.

II. PLANE WAVE SCATTERING AT A SINUSOIDAL INTERFACE

We consider the problem of a plane wave that impinges on the interface between two dielectric media. The interface is described by the function

$$f(z) = a \sin \phi z. \quad (2)$$

If only first-order scattering is considered, more general shapes of $f(z)$ can be synthesized by superposition of sinusoidal functions. For higher-order processes, mixing of the sinusoidal terms occurs so that the description of scattering from more general interfaces becomes complicated. The incident plane wave is given by (the time dependence is $e^{i\omega t}$)

$$E_y = A e^{-i(\kappa_i x + \beta_i z)} \quad (3)$$

$$H_x = -\frac{\beta_i}{\omega \mu_0} A e^{-i(\kappa_i x + \beta_i z)} \quad (4)$$

$$H_z = \frac{\kappa_i}{\omega \mu_0} A e^{-i(\kappa_i x + \beta_i z)}. \quad (5)$$

The coordinate system is shown in Fig. 1. It is assumed that no variation of either the field components or the material parameters exists in y direction so that we can symbolically write

$$\frac{\partial}{\partial y} = 0. \quad (6)$$

The remaining three field components E_x , E_z , and H_y vanish. The parameters κ_i and β_i are connected by the following equation

$$\kappa_i = (n_1^2 k^2 - \beta_i^2)^{1/2}. \quad (7)$$

The refractive index n_1 belongs to the medium from which the plane wave approaches the interface, and k is the free space propagation constant.

The reflected and scattered waves are expressed as superpositions of plane waves. We thus have in medium 1

$$E_y = \int_{-\infty}^{\infty} B(\beta)e^{i(\sigma x - \beta z)} d\beta \tag{8}$$

$$H_x = -\frac{1}{\omega\mu} \int_{-\infty}^{\infty} \beta B(\beta)e^{i(\sigma x - \beta z)} d\beta \tag{9}$$

$$H_z = -\frac{1}{\omega\mu} \int_{-\infty}^{\infty} \sigma B(\beta)e^{i(\sigma x - \beta z)} d\beta, \tag{10}$$

with

$$\sigma = (n_1^2 k^2 - \beta^2)^{1/2}, \tag{11}$$

and similarly in medium 2

$$E_y = \int_{-\infty}^{\infty} C(\beta)e^{-i(\rho x + \beta z)} d\beta \tag{12}$$

$$H_x = -\frac{1}{\omega\mu} \int_{-\infty}^{\infty} \beta C(\beta)e^{-i(\rho x + \beta z)} d\beta \tag{13}$$

$$H_z = \frac{1}{\omega\mu} \int_{-\infty}^{\infty} \rho C(\beta)e^{-i(\rho x + \beta z)} d\beta, \tag{14}$$

with

$$\rho = (n_2^2 k^2 - \beta^2)^{1/2}. \tag{15}$$

The boundary conditions at the dielectric interface require the tangential component of the electric field E_y and the tangential component of the magnetic field

$$H_t = \frac{1}{(1 + f'^2)^{1/2}} H_z + \frac{f'}{(1 + f'^2)^{1/2}} H_x \tag{16}$$

to be continuous [f' is the z derivative of the interface function (2)]. The boundary conditions thus lead to the two equations:

$$Ae^{-i(\kappa_i f + \beta_i z)} + \int_{-\infty}^{\infty} B(\beta)e^{i(\sigma f - \beta z)} d\beta = \int_{-\infty}^{\infty} C(\beta)e^{-i(\rho f + \beta z)} d\beta \tag{17}$$

and

$$\begin{aligned} (\kappa_i - \beta_i f')Ae^{-i(\kappa_i f + \beta_i z)} - \int_{-\infty}^{\infty} (\sigma + \beta f')B(\beta)e^{i(\sigma f - \beta z)} d\beta \\ = \int_{-\infty}^{\infty} (\rho - \beta f')C(\beta)e^{-i(\rho f + \beta z)} d\beta. \end{aligned} \tag{18}$$

In order to remove the z dependence from the equation, we multiply

with $\exp(-i\beta'z)$ and integrate over z from $-\infty$ to ∞ . This procedure transforms the equations to the following form:

$$AF(\kappa_i, \beta_i - \beta') + \int_{-\infty}^{\infty} B(\beta)F(-\sigma, \beta - \beta') d\beta = \int_{-\infty}^{\infty} C(\beta)F(\rho, \beta - \beta') d\beta, \quad (19)$$

and

$$AG(\kappa_i, \beta_i, \beta') + \int_{-\infty}^{\infty} B(\beta)G(-\sigma, \beta, \beta') d\beta = \int_{-\infty}^{\infty} C(\beta)G(\rho, \beta, \beta') d\beta \quad (20)$$

with

$$F(\eta, \beta - \beta') = \int_{-\infty}^{\infty} e^{-i[\eta f(z) + (\beta - \beta')z]} dz, \quad (21)$$

and

$$G(\eta, \beta, \beta') = \int_{-\infty}^{\infty} [\eta - \beta f'(z)] e^{-i[\eta f(z) + (\beta - \beta')z]} dz. \quad (22)$$

It is shown in the Appendix that F and G are related in the following way:

$$G(\eta, \beta, \beta') = \frac{\eta^2 + \beta^2 - \beta\beta'}{\eta} F(\eta, \beta - \beta'). \quad (23)$$

Substitution of (2) into (21) yields

$$F(\eta, \beta - \beta') = \int_{-\infty}^{\infty} e^{-ia\eta \sin \phi z} e^{-i(\beta - \beta')z} dz. \quad (24)$$

The first exponential function under the integral sign can be expressed as a series in terms of Bessel functions with the help of the generating function of the Bessel functions. The remaining integration yields delta functions so that we obtain

$$F(\eta, \beta - \beta') = 2\pi \sum_{\nu=-\infty}^{\infty} J_{\nu}(\eta a) \delta[\beta - \beta' + \nu\phi]. \quad (25)$$

Substitution of (23) and (25) transforms the equation systems (19) and (20) to the form

$$\begin{aligned} \sum_{\nu=-\infty}^{\infty} \{-B(\beta_{\nu})J_{\nu}(\sigma, a) + C(\beta_{-\nu})J_{\nu}(\rho_{-}, a)\} \\ = A \sum_{\nu=-\infty}^{\infty} J_{\nu}(\kappa_i a) \delta(\beta_i - \beta' + \nu\phi), \end{aligned} \quad (26)$$

and

$$\sum_{\nu=-\infty}^{\infty} \left\{ \frac{n_1^2 k^2 - \beta_\nu \beta'}{\sigma_\nu} B(\beta_\nu) J_\nu(\sigma_\nu a) + \frac{n_2^2 k^2 - \beta_{-\nu} \beta'}{\rho_{-\nu}} C(\beta_{-\nu}) J_\nu(\rho_{-\nu} a) \right\} = A \sum_{\nu=-\infty}^{\infty} \frac{n_1^2 k^2 - \beta_i \beta'}{\kappa_i} J_\nu(\kappa_i a) \delta(\beta_i - \beta' + \nu\phi). \tag{27}$$

The following abbreviations were used

$$\beta_\nu = \beta' + \nu\phi \tag{28}$$

$$\sigma_\nu = (n_1^2 k^2 - \beta_\nu^2)^{\frac{1}{2}} \tag{29}$$

and

$$\rho_\nu = (n_2^2 k^2 - \beta_\nu^2)^{\frac{1}{2}}. \tag{30}$$

We know from the theory of phase gratings² that only discrete plane waves appear in the reflected and transmitted beams. This and the appearance of the delta functions on the right-hand side of (26) and (27) suggest that the solutions should be of the form

$$B(\beta) = \sum_{\mu=-\infty}^{\infty} b_\mu \delta(\beta - \beta_i - \mu\phi), \tag{31}$$

and

$$C(\beta) = \sum_{\mu=-\infty}^{\infty} c_\mu \delta(\beta - \beta_i - \mu\phi). \tag{32}$$

Substitution of (31) and (32) into (26) and (27) and comparison of the coefficients of the delta functions of equal arguments leads to two infinite equation systems for the unknown coefficients b_μ and c_μ .

$$\sum_{\nu=-\infty}^{\infty} \{-b_{n+\nu} J_\nu(\sigma_{n+\nu} a) + c_{n-\nu} J_\nu(\rho_{n-\nu} a)\} = A J_n(\kappa_i a), \tag{33}$$

and

$$\sum_{\nu=-\infty}^{\infty} \left\{ \frac{\sigma_n^2 - \nu\phi(\beta_i + n\phi)}{\sigma_{n+\nu}} b_{n+\nu} J_\nu(\sigma_{n+\nu} a) + \frac{\rho_n^2 + \nu\phi(\beta_i + n\phi)}{\rho_{n-\nu}} c_{n-\nu} J_\nu(\rho_{n-\nu} a) \right\} = \frac{\kappa_i^2 - n\phi\beta_i}{\kappa_i} A J_n(\kappa_i a). \tag{34}$$

The equation system (33) and (34) is exact. An exact solution appears impossible to obtain. However, the equation system is very convenient for obtaining perturbation solutions of arbitrary order. It is also possible to obtain a solution that is exact in the limit $\phi \rightarrow 0$.

The reflected and transmitted fields follow directly from (8), (12), (31), and (32). In medium 1 we obtain for the reflected field

$$E_y = \sum_{\mu=-\infty}^{\infty} b_{\mu} e^{i(\sigma_{\mu}x - \beta_{\mu}z)}. \quad (35)$$

The field in medium 2 is

$$E_y = \sum_{\mu=-\infty}^{\infty} c_{\mu} e^{-i(\rho_{\mu}x + \beta_{\mu}z)}. \quad (36)$$

The parameters σ_{μ} and ρ_{μ} are defined by (29) and (30). However, now we must use these equations with $\beta_{\mu} = \beta_i + \mu\phi$.

III. SOLUTION FOR $\phi \rightarrow 0$

In the limit $\phi \rightarrow 0$, an exact solution of the equation systems (33) and (34) can be obtained. If $\phi = 0$ is assumed, we use the facts that

$$\sigma_{\nu} = \kappa_i \quad (37)$$

$$\rho_{\nu} = \rho_0 \quad (38)$$

to write (33) and (34) in the form

$$\sum_{\nu=-\infty}^{\infty} \{-b_{n+\nu} J_{\nu}(\kappa_i a) + c_{n-\nu} J_{\nu}(\rho_0 a)\} = A J_n(\kappa_i a) \quad (39)$$

$$\sum_{\nu=-\infty}^{\infty} \{\kappa_i b_{n+\nu} J_{\nu}(\kappa_i a) + \rho_0 c_{n-\nu} J_{\nu}(\rho_0 a)\} = \kappa_i A J_n(\kappa_i a). \quad (40)$$

It is now possible to eliminate c_{ν} from the equations and obtain an equation system for b_{ν} alone,

$$\sum_{\nu=-\infty}^{\infty} b_{n+\nu} J_{\nu}(\kappa_i a) = \frac{\kappa_i - \rho_0}{\kappa_i + \rho_0} A J_n(\kappa_i a). \quad (41)$$

Similarly, we obtain by eliminating b_{ν}

$$\sum_{\nu=-\infty}^{\infty} c_{n-\nu} J_{\nu}(\rho_0 a) = \frac{2\kappa_i}{\kappa_i + \rho_0} A J_n(\kappa_i a). \quad (42)$$

These equations can be solved with the help of the addition theorem for Bessel functions

$$\sum_{\nu=-\infty}^{\infty} J_{n+\nu}(x) J_{\nu}(y) e^{i\nu\theta} = J_n(R) e^{in\theta}, \quad (43)$$

with

$$R = (x^2 + y^2 - 2xy \cos \theta)^{1/2}. \tag{44}$$

With $\theta = 0$, we obtain $R = x - y$ so that we see immediately that

$$b_\mu = \frac{\kappa_i - \rho_0}{\kappa_i + \rho_0} AJ_\mu(2\kappa_i a) \tag{45}$$

is the solution of (41). Since $J_{-\nu}(x) = J_\nu(-x)$ (for integer values of ν) it is also apparent that

$$c_\mu = \frac{2\kappa_i}{\kappa_i + \rho_0} AJ_\mu[(\kappa_i - \rho_0)a] \tag{46}$$

is the solution of (42).

The solutions (45) and (46) are exact for $\phi = 0$. One might expect that $\phi = 0$ describes a plane dielectric interface so that no side lobes should be expected. Even though it is true that all the sidelobes coincide for $\phi = 0$, the solutions (45) and (46) do hold approximately even if $\phi \neq 0$. The sinusoidal shape of the interface is apparently built into the equation system (41) and (42) even though ϕ does not appear explicitly. The solutions (45) and (46) are approximations that hold if $k\phi \ll 1$. These solutions show that the amplitudes of the side lobes are proportional to Bessel functions. This result is well known from the theory of phase gratings.²

IV. PERTURBATION SOLUTIONS

For our purposes, the solution for $k\phi \ll 1$ is of little use. Therefore, we proceed to derive approximate solutions that hold for

$$ka \ll 1. \tag{47}$$

We use the following approximations for the Bessel functions of small argument

$$J_0(x) = 1 - \frac{x^2}{4} \tag{48}$$

$$J_1(x) = -J_{-1}(x) = \frac{x}{2} \left(1 - \frac{x^2}{8}\right) \tag{49}$$

$$J_2(x) = J_{-2}(x) = \frac{x^2}{8}. \tag{50}$$

In addition, we assume that b_o and c_o are zero-order terms, $b_{\pm 1}$ and $c_{\pm 1}$ are of first order, and $b_{\pm 2}$ and $c_{\pm 2}$ are of second order.

By neglecting all but zero-order terms in (33) and (34), we obtain to zero order of approximation for $n = 0$,

$$\left. \begin{aligned} -b_0 + c_0 &= A \\ \sigma_0 b_0 + \rho_0 c_0 &= \kappa_i A \end{aligned} \right\}. \quad (51)$$

Taking $n = 1$, we obtain to first order

$$\left. \begin{aligned} -b_1 + c_1 &= \frac{a}{2} (\kappa_i A - \sigma_0 b_0 - \rho_0 c_0) \\ \sigma_1 b_1 + \rho_1 c_1 &= \frac{a}{2} \{ (\kappa_i^2 - \phi \beta_i) A + [\sigma_1^2 + \phi(\beta_i + \phi)] b_0 \\ &\quad - [\rho_1^2 + \phi(\beta_i + \phi)] c_0 \} \end{aligned} \right\}. \quad (52)$$

The corresponding equation system for b_{-1} and c_{-1} is obtained by using $n = -1$. It has exactly the same form as the equation system (52) and is obtained by replacing b_1 with $-b_{-1}$, c_1 with $-c_{-1}$ and ϕ with $-\phi$.

Finally, we obtain the second-order approximation by setting $n = 2$ and keeping only terms up to second order

$$\left. \begin{aligned} -b_2 + c_2 &= \frac{a^2}{8} (\kappa_i^2 A + \sigma_0^2 b_0 - \rho_0^2 c_0) - \frac{a}{2} (\sigma_1 b_1 + \rho_1 c_1) \\ \sigma_2 b_2 + \rho_2 c_2 &= \frac{a^2}{8} \{ \kappa_i (\kappa_i^2 - 2\phi \beta_i) A - \sigma_0 [\sigma_2^2 + 2\phi(\beta_i + 2\phi)] b_0 \\ &\quad - \rho_0 [\rho_2^2 + 2\phi(\beta_i + 2\phi)] c_0 \} \\ &\quad + \frac{a}{2} \{ [\sigma_2^2 + \phi(\beta_i + 2\phi)] b_1 - [\rho_2^2 + \phi(\beta_i + 2\phi)] c_1 \} \end{aligned} \right\}. \quad (53)$$

The equations for b_{-2} and c_{-2} are obtained by replacing b_2 with b_{-2} , c_2 with c_{-2} , b_1 with $-b_{-1}$, c_1 with $-c_{-1}$, and finally ϕ with $-\phi$ in (53).

It is immediately apparent that each order of approximation follows from the preceding order. We can thus solve all the equations (51), (52), and (53) in succession. Each time we need solve only two equations with two unknowns. The result of the previous approximation is then used to obtain the next higher order of approximation from the next equation system.

The solutions of these equations are listed below.

$$b_0 = \frac{\kappa_i - \rho_0}{\kappa_i + \rho_0} A \quad c_0 = \frac{2\kappa_i}{\kappa_i + \rho_0} A \quad (54)$$

$$b_1 = \frac{\kappa_i (\kappa_i - \rho_0)}{\sigma_1 + \rho_1} a A \quad c_1 = b_1 \quad (55)$$

$$\left. \begin{aligned} b_2 &= \frac{a^2 A}{4} \kappa_i \frac{\kappa_i - \rho_0}{\sigma_2 + \rho_2} \left[\rho_0 + \rho_2 + 2(\kappa_i - \rho_0) \frac{\kappa_i + \rho_0}{\sigma_1 + \rho_1} \right] \\ c_2 &= \frac{a^2 A}{4} \kappa_i \frac{\kappa_i - \rho_0}{\sigma_2 + \rho_2} \left[\rho_0 - \sigma_2 + 2(\kappa_i - \rho_0) \frac{\kappa_i + \rho_0}{\sigma_1 + \rho_1} \right] \end{aligned} \right\} \quad (56)$$

The coefficients b_{-1} and c_{-1} are obtained from (55) by changing the sign in front of the terms and changing the subscripts 1 to -1 on the right-hand side of the equation. The signs of b_{-2} and c_{-2} are the same as those of the coefficients in (56). We obtain these coefficients by changing the signs of the subscripts on the right-hand side of the equations. For $ka \ll 1$ and $k\phi \ll 1$ the equations (45) and (46) can be shown to be identical with (54), (55), and (56).

The amplitudes may belong to plane traveling waves or to evanescent waves. Whether a wave is of the propagating or evanescent type depends on whether the parameters σ_μ and ρ_μ are real or imaginary. For real values we obtain traveling waves while imaginary values indicate that the field decays exponentially with increasing distance from the interface indicating an evanescent wave. The propagation constant in z direction, β_μ , is obtained from (28) by replacing β' with β_i . We thus have for traveling as well as for evanescent waves

$$\beta_\mu = \beta_i + \mu\phi \quad \mu = 0, \pm 1, \pm 2, \dots \quad (57)$$

V. CALCULATION OF SLAB WAVEGUIDE LOSSES

Since the guided modes of the slab waveguide can be expressed as the superposition of two plane waves, whose propagation vectors form equal but opposite angles with the z axis, we can use our present results immediately to calculate the radiation losses suffered by the guided slab waveguide modes.⁴ Our calculation applies to TE modes. However, for slight index differences the losses of TE modes and TM modes are nearly identical. The slab waveguide geometry is shown in Fig. 2.

The radiation losses of slab waveguide modes have a somewhat complicated dependence on either frequency or slab width, since the interference of the waves scattered at one of the two dielectric interfaces with the radiation from the other interface—and also the interference with radiation that is reflected at the opposite interface—has to be taken into account. However, these interference effects cause only fluctuations about an average value. If we content ourselves with establishing only the average loss value, disregarding the fluctuations, the description of radiation losses is greatly simplified. We also gain the advantage of obtaining simpler mathematical expressions. In the spirit of this sim-

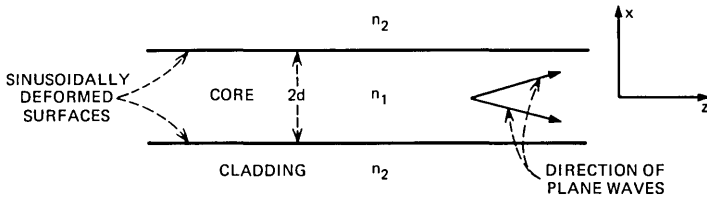


Fig. 2—Cross section of the slab waveguide.

plication, we consider all scattered power as lost—with the exception of those waves that are captured in the core—disregarding incomplete reflection from the other dielectric interface and interference with the directly scattered radiation. The amount of power scattered per unit length of the waveguide is

$$S_x = \frac{1}{2} |(\mathbf{E}_1 \times \mathbf{H}_1^*)_x| + \frac{1}{2} |(\mathbf{E}_2 \times \mathbf{H}_2^*)_x|. \quad (58)$$

The electric and magnetic fields in this expression are only the scattered, untrapped part of the field, exclusive of the incident field. The subscripts 1 and 2 refer to the fields in mediums 1 and 2. With the assumption that only one side lobe is instrumental in dissipating power by radiation we can write

$$S_x = \frac{1}{2\omega\mu} (\sigma_{-\nu} |b_{-\nu}|^2 + \rho_{-\nu} |c_{-\nu}|^2). \quad (59)$$

The subscript ν assumes the values 1 and 2 for first- and second-order light scattering. The negative values must be used (ν is now assumed to be positive) because $\sigma_{-\nu}$ and $\rho_{-\nu}$ must be real since evanescent waves do not carry power.

The power attenuation is

$$\alpha = \frac{S_x}{2S_z d}. \quad (60)$$

S_z is the power per unit length (unit area in the three dimensional case) that is carried by the plane wave in z direction. The total power carried by this one plane wave component in z direction inside of the waveguide core of width $2d$ is thus $2S_z d$. The ratio of the power lost (per unit length along the waveguide axis) divided by the power carried by the wave is the power loss per unit length. Actually, two plane waves are needed to describe the guided mode in the slab waveguide. However, thus far we have considered only the scattering loss from one of the two interfaces. It is sufficient to consider that each of the two plane wave components scatters from one interface. Inclusion of the scattering

from both interfaces introduces a factor of two in the numerator of (60). However, another factor of two is introduced in the denominator by adding the power of the other plane wave to the total power carried by the guided mode. The expression (60) thus holds for the scattering loss of the guided modes provided that both interfaces contribute an equal amount to the power loss. Using the following expression for the power flow density in z -direction,

$$S_z = \frac{\beta_i}{2\omega\mu_0} |A|^2, \quad (61)$$

we obtain the general expression for scattering losses from a dielectric slab waveguide with sinusoidally deformed core-cladding interfaces

$$\alpha_\nu = \frac{1}{2\beta_i |A|^2 d} (\sigma_{-\nu} |b_{-\nu}|^2 + \rho_{-\nu} |c_{-\nu}|^2). \quad (62)$$

The index ν indicates the order of the scattering process. Taking $\nu = 1$, we obtain with the help of (55) the scattering loss contribution from the first-order grating lobes

$$\alpha_1 = \frac{\kappa_i^2 (n_1^2 - n_2^2) k^2}{2\beta_i (\sigma_{-1} + \rho_{-1}) d} a^2. \quad (63)$$

We used the fact that κ_i is real, while ρ_o is imaginary, so that we have $|\kappa_i + \rho_o|^2 = \kappa_i^2 - \rho_o^2 = (n_1^2 - n_2^2) k^2$.

The loss contribution of the second-order side lobes follows from (56) and (62) with $\nu = 2$.

$$\alpha_2 = \frac{a^4 \kappa_i^2 (n_1^2 - n_2^2) k^2}{32\beta_i (\sigma_{-2} + \rho_{-2}) d} [4\sigma_{-1}^2 + \sigma_{-2}\rho_{-2} + (\gamma_i - 2\gamma_{-1})^2]. \quad (64)$$

The γ parameters are defined by the equations

$$\gamma_i = i\rho_o = (\beta_i^2 - n_2^2 k^2)^{1/2} \quad (65)$$

and

$$\gamma_{-1} = i\rho_{-1} = [(\beta_i - \phi)^2 - n_2^2 k^2]^{1/2}. \quad (66)$$

We have tacitly assumed that the first-order side lobe belongs to an evanescent wave if the loss contribution of the second-order side lobe is being considered. The parameters γ_i and γ_{-1} of (65) and (66) are consequently assumed to be real quantities. If the first-order side lobe propagates in medium 2 as a traveling wave, the second-order side lobe also gives rise to a traveling wave. However, since the contribution

from the grating lobe of first order is much stronger than that from the second-order lobe, we can neglect the loss contribution (64). If first-order scattering does occur, it is the predominant effect. Only if first-order scattering contributes only an evanescent wave, and does not cause radiation loss, must the second-order loss process (64) be considered.

VI. DISCUSSION

We have calculated the loss contributions that result from the grating lobes of first and second order when a plane wave impinges on the sinusoidally deformed interface between two different dielectric media. Not all grating orders belong to traveling waves. Some grating orders cause evanescent waves in medium 2 and guided waves in the core so that they do not contribute to radiation loss of a guided wave in medium 1. The first-order radiation loss coefficient (63) is proportional to $(ak)^2$ while the second-order radiation loss coefficient (64) is proportional to $(ak)^4$. If both processes are effective simultaneously, the lower order process is dominant. Whether the second-order process is the only cause of radiation loss, or whether both first- and second-order processes are acting simultaneously depends on the magnitude of the mechanical frequency ϕ of the sinusoidal interface distortion. If γ_{-1} of (66) is real, the first-order side lobe belongs to an evanescent wave, and only the second-order side lobe is causing radiation losses. When γ_{-1} is imaginary, both first- and second-order side lobes carry away real power. Equation (64) is not applicable in this case since it was derived under the assumption that (65) and (66) are both real. However, if first-order radiation losses are possible, the second-order loss coefficient gives only a small contribution to the total radiation loss.

It is interesting to compare the result of our present theory with earlier results obtained from a modal analysis of the slab waveguide problem. From equation (79) of Ref. 1, we obtain in our present notation

$$\alpha = \frac{a^2 k^2 (n_1^2 - n_2^2) \kappa_i^2}{4\beta_i d \left(1 + \frac{1}{\gamma_i d}\right)} \left[\frac{\rho_{-1} \cos^2 \sigma_{-1} d}{\rho_{-1}^2 \cos^2 \sigma_{-1} d + \sigma_{-1}^2 \sin^2 \sigma_{-1} d} + \frac{\rho_{-1} \sin^2 \sigma_{-1} d}{\rho_{-1}^2 \sin^2 \sigma_{-1} d + \sigma_{-1}^2 \cos^2 \sigma_{-1} d} \right]. \quad (67)$$

Equation (67) is modified for the case that both interfaces are sinusoidally distorted, but with a random phase relationship between the two sinusoidal functions. The radiation loss (67) is considered to be an

ensemble average over slab waveguides with all possible phase relationships between the sinusoidal distortions of the two interfaces. In addition, we used the expression

$$\cos^2 \kappa_i d = \frac{\kappa_i^2}{(n_1^2 - n_2^2)k^2} \quad (68)$$

which represents the eigenvalue equation of the even TE modes of the slab waveguide.

It was shown in equations (39) and (40) of Ref. 6 that the average value of the expression in brackets of eq. (67) is given by

$$2/(\sigma_{-1} + \rho_{-1}). \quad (69)$$

Combining (69) with (67) makes it apparent that the average value of the slab waveguide radiation loss is identical to the loss coefficient (63) that was derived from the plane wave model. The only remaining difference can be explained as stemming from the fact that the actual width $2d$ of the slab must be replaced by the effective slab width

$$2d\left(1 + \frac{1}{\gamma_i d}\right) \quad (70)$$

which is caused by the exponential field tail reaching out into the cladding.

The plane wave model used in this paper does not include the interference effects of waves originating from the two interfaces, and from the reflection of the scattered light from the opposite interface. But the averaging processes that were involved in converting the precise scattering loss coefficient of the slab waveguide theory into the loss coefficient derived from the plane wave model may be expected to be effective in a real waveguide. If the phase of the sinusoidal interface distortion varies slowly and randomly along the waveguide, the phase average that is already incorporated in (67) may actually occur. The average over the expression in brackets in (67) would occur either with randomly changing slab half width d , or with randomly varying mechanical frequency ϕ of the sinusoidal interface changes. The loss formula (63) is much simpler than (67), and is of actual practical value for the indicated reasons. The same averaging process involved in the first-order loss coefficient (63) is also implicit in the second-order loss coefficient (64).

Second-order scattering couples two modes ν and μ , if their propagation constants satisfy the condition

$$|\beta_\nu - \beta_\mu| = 2\phi. \quad (71)$$

Mode ν is a guided mode, while mode μ can be either a guided or a radiation mode. Equation (71) replaces eq. (1) for first-order scattering. Higher-order processes have similar coupling laws with the order of the process multiplying the mechanical frequency ϕ .

Second-order scattering losses are of importance for intentionally coupled multimode operation. Mode coupling reduces the pulse delay distortion that is caused by the different group velocities of the modes. It is possible to design the core-cladding interface irregularities in such a way that all guided modes (with the exception of the last) are coupled to each other without coupling to the continuous spectrum of radiation modes by first-order processes.⁷ However, radiation losses via second- and higher-order processes are still possible. The discussion of second-order losses for intentionally coupled multimode operation is the subject of a companion paper.⁵

APPENDIX

Proof of the Relation (23)

We begin by using the fact that a sinusoidal function with an infinite argument can be regarded as zero. This assertion is the basis for the following definition of the delta function:

$$\delta(x) = \lim_{A \rightarrow \infty} \frac{1}{\pi} \frac{\sin xA}{x} \quad (72)$$

vanishes everywhere except at the point $x = 0$. The singularity of the delta function at $x = 0$ is caused by the appearance of x in the denominator. Without this denominator, we are justified to define

$$\lim_{A \rightarrow \infty} \sin xA = 0. \quad (73)$$

Using eq. (73), we write the following identity

$$\begin{aligned} 0 &= \lim_{A \rightarrow \infty} -2i \sin [\eta f(A) + (\beta - \beta')A] \\ &= \lim_{A \rightarrow \infty} \int_{-A}^A \frac{d}{dz} e^{-i[\eta f(z) + (\beta - \beta')z]} dz \\ &= \int_{-\infty}^{\infty} -i[\eta f' + (\beta - \beta')] e^{-i[\eta f + (\beta - \beta')z]} dz. \end{aligned} \quad (74)$$

It was assumed that $f(-z) = -f(z)$. From the last line of (74), we obtain immediately

$$\int_{-\infty}^{\infty} f' e^{-i[\eta f + (\beta - \beta')z]} dz = -\frac{\beta - \beta'}{\eta} \int_{-\infty}^{\infty} e^{-i[\eta f + (\beta - \beta')z]} dz$$

or

$$\begin{aligned} \int_{-\infty}^{\infty} [\eta - \beta f'(z)] e^{-i[\eta f(z) + (\beta - \beta')z]} dz \\ = \frac{\eta^2 + \beta(\beta - \beta')}{\eta} \int_{-\infty}^{\infty} e^{-i[\eta f(z) + (\beta - \beta')z]} dz. \end{aligned} \quad (75)$$

From the definitions (21) and (22) it follows that (75) is identical with (23).

REFERENCES

1. Marcuse, D., "Mode Conversion Caused by Surface Imperfections of a Dielectric Slab Waveguide," *B.S.T.J.*, 48, No. 10 (December 1969), pp. 3187-3215.
2. Goodman, J. W., *Introduction to Fourier Optics*, New York: McGraw-Hill Book Company, 1968.
3. Tien, P. K., "Light Waves in Thin Films and Integrated Optics," *Appl. Phys.*, 10, No. 11 (November 1971), pp. 2395-2413.
4. Marcuse, D., "Hollow Dielectric Waveguides for Distributed Feedback Lasers," *IEEE J. Quantum Elec.*, QE-8, No. 7 (July 1972), pp. 661-669.
5. Marcuse, D., "Higher-Order Loss Processes and the Loss Penalty of Multimode Operation," *B.S.T.J.*, this issue, pp. 1819-1836.
6. Marcuse, D., "Power Distribution and Radiation Losses in Multimode Dielectric Slab Waveguides," *B.S.T.J.*, 51, No. 2 (February 1972), pp. 429-454.
7. Marcuse, D., "Pulse Propagation in Multimode Dielectric Waveguides," *B.S.T.J.*, 51, No. 6 (July-August 1972), pp. 1199-1232.

Higher-Order Loss Processes and the Loss Penalty of Multimode Operation

By D. MARCUSE

(Manuscript received May 5, 1972)

Pulse spreading caused by the different group velocities of the guided modes of a multimode waveguide can be reduced by providing intentional coupling between the modes. Coupling among the guided modes inevitably leads to radiation losses. This loss penalty is discussed for two types of loss processes. We consider that the highest-order mode loses power by second-order coupling to the continuous spectrum of radiation modes. We also consider a loss process that is caused by nonresonant coupling of guided modes to lossy neighboring modes. Both loss processes can cause a substantial loss penalty. However, the loss penalty can always be reduced by limiting the intentional coupling to fewer of the guided modes, allowing the highest-order modes to die out. The discussion is based on a slab waveguide model.

I. INTRODUCTION

Higher-order loss processes have been discussed in a previous paper.¹ The idea of loss processes of different orders is based on perturbation theory. Two modes of a dielectric waveguide are coupled if their propagation constants obey the relation¹

$$|\beta_\nu - \beta_\mu| = m\phi. \quad (1)$$

ϕ is the mechanical frequency of the Fourier spectrum of the coupling function; m is a positive integer that specifies the order of the coupling process. If $m = 1$, mode ν is coupled to mode μ by a first-order process, $m = 2$ indicates a second-order process, etc. For small values of a/λ_0 , (a is the Fourier amplitude that belongs to the mechanical frequency ϕ ; λ_0 is the free space wavelength of the light in the waveguide) the coupling strength is proportional to $(a/\lambda_0)^m$ so that the coupling decreases with increasing order of the coupling process. If mode ν represents a guided mode, mode μ may either be a guided or a radiation mode. In the latter

case, mode ν loses power by radiation. An explanation of the coupling process in terms of diffraction gratings is given in Ref. 1.

There is a different loss process that cannot be understood in terms of higher-order grating lobes. Consider two guided modes. Mode 2 is inherently lossless, while mode 1 suffers high loss. If we couple these two modes by means of a first-order process, a large amount of loss will be transferred from the lossy mode to the hitherto lossless neighbor. However, even if we couple these two modes by means of a sinusoidal coupling function the mechanical frequency of which does not satisfy (1) for any integer m , some loss will be imparted from the lossy mode to the inherently lossless mode. If both modes were lossless, no significant amount of power would be interchanged among them if (1) is not satisfied for $m = 1$ (or any other integer). We call such a coupling process "nonresonant coupling." The small amount of power that flows momentarily from mode 1 to mode 2 is returned in the next instant because the phase relationship required for continuous power flow from one mode to the other does not exist. However, if mode 1 is lossy, mode 2 transfers a small amount of power to mode 1 (even via the nonresonant coupling process) that can not be returned since some of the power is already dissipated in the lossy mode 1. This nonresonant coupling process has the effect of imparting some of the high loss of one mode to a neighboring mode. The attenuation coefficient that results is derived in the Appendix.

In this paper, we calculate the loss penalty that stems from intentional mode coupling in a multimode slab waveguide caused by these higher-order processes. Reference 2 presents the theory of pulse propagation in multimode waveguides in the presence of first-order coupling between the guided modes. The purpose of the coupling is to reduce pulse distortion.³ It was pointed out that it is possible to couple all the guided modes by a first-order process without causing first-order radiation losses.² This possibility arises from the fact that the modes of a slab waveguide are arranged in β space, such that the spacing between neighboring modes increases with increasing mode number. Because of the coupling law (1) with $m = 1$, it is possible to couple all the guided modes, except mode N , by providing a spectrum of mechanical frequencies that has a proper Fourier component for coupling at least the nearest neighbors of all the modes. However, mode N is not coupled to mode $N - 1$ if the Fourier spectrum has an abrupt cutoff so that no mechanical frequency exists that satisfies the relation

$$\beta_{N-1} - \beta_N = \phi. \quad (2)$$

Residual losses result from the fact that it is unrealistic to assume a

Fourier spectrum with an abrupt cutoff. Residual coupling between mode $N - 1$ and mode N is then possible via the tail of the Fourier spectrum. Mode N is necessarily coupled by first-order processes to the continuous spectrum of radiation modes so that mode $N - 1$, being coupled to mode N , suffers loss which causes power loss to the entire ensemble of coupled guided modes. We must now consider the effect of higher-order processes. Even with a Fourier spectrum with perfectly abrupt cutoff, mode $N - 1$ is coupled to the spectrum of radiation modes by second- and higher-order processes. Assuming small amplitudes for the Fourier coefficients, we neglect processes of third- and higher-order and discuss radiation losses caused by the second-order process. We shall see that substantial losses can result even via the second-order loss mechanism. However, luckily, we can readjust the intentional coupling between the guided modes to prevent first-order coupling not only to mode N but also to mode $N - 1$. The distance (in β space) between mode $N - 2$ and the continuum of radiation modes is then greater than 2ϕ so that the second-order loss process is no longer possible. The uncoupled modes (uncoupled from the remaining guided modes) lose power by being coupled to the radiation field. Mode N loses power very rapidly because it is coupled by means of a first-order process. Mode $N - 1$ loses power by means of a second-order process. If the loss caused by coupling of the guided modes to mode $N - 1$ was bothersome, its loss is certainly sufficient to prevent pulse distortion by power flowing along in this mode. It is thus clear that radiation losses can be reduced by limiting the intentional coupling to the lower-order guided modes leaving a few of the higher-order modes to die out because of their high radiation losses.

In a similar manner, nonresonant coupling between the lossy mode N (coupled by a first-order process to radiation modes) and the neighboring guided modes $N - 1$, $N - 2$, etc., influences the loss behavior of the intentionally coupled guided modes. The loss penalty caused by this nonresonant coupling mechanism is considered separately from the higher-order loss process mentioned earlier in order to assess the separate influence of each mechanism. Again, it is advantageous to uncouple some of the higher-order modes from the lower-order guided modes since the nonresonant coupling process decreases in strength with increasing distance (in β space) of the guided modes from the lossy mode N .

II. SUMMARY OF COUPLED POWER THEORY

The coupled power theory presented in Ref. 2 was based on the stochastic partial differential equation for the average power of the modes

$$\frac{\partial P_\nu}{\partial z} + \frac{1}{v_\nu} \frac{\partial P_\nu}{\partial t} = -\alpha_\nu P_\nu + \sum_{\mu=1}^N h_{\nu\mu} (P_\mu - P_\nu). \quad (3)$$

The power loss coefficient α_ν , for mode ν with group velocity v_ν , incorporates heat losses as well as radiation losses. However, heat losses will be ignored in our present discussion. For slab waveguides with random core-cladding interface perturbations, the coupling coefficient assumes the form

$$h_{\nu\mu} = \frac{n_1^2 k^2 \sin^2 \theta_\nu \sin^2 \theta_\mu}{2d^2 \left(1 + \frac{1}{\gamma_\nu d}\right) \left(1 + \frac{1}{\gamma_\mu d}\right) \cos \theta_\nu \cos \theta_\mu} F(\beta_\nu - \beta_\mu). \quad (4)$$

The mode angle θ_ν is defined in terms of the refractive index n_1 of the core, the free space propagation constant k and the propagation constant β_ν of the ν th mode.

$$\cos \theta_\nu = \frac{\beta_\nu}{n_1 k}. \quad (5)$$

The parameter γ_ν appearing in (4) is ($n_2 =$ cladding index)

$$\gamma_\nu = (\beta_\nu^2 - n_2^2 k^2)^{1/2} \quad (6)$$

and d is the slab half width. The function $F(\phi)$ is the ensemble average of the square of the Fourier transform of the core-cladding interface function. It will be referred to as the "power spectrum." For the purpose of this paper, we assume that $F(\phi)$ is constant from zero to the cutoff value ϕ_c of ϕ . For $\phi > \phi_c$ we assume that $F(\phi) = 0$.

For sufficiently large values of z , we obtain the following approximate solution of (3)²

$$P_\nu(z, t) = \frac{2\tau}{\Delta t} k_1 B_{\nu_0}^{(1)} e^{-\alpha_\nu z} \exp \left[-\left(\frac{t - z/v_\nu}{\Delta t/2} \right)^2 \right], \quad (7)$$

with the full width of the Gaussian pulse given by

$$\Delta t = 2(\tau^2 + 4\alpha_\nu^2 z^2)^{1/2}. \quad (8)$$

The input pulse is determined by its half width τ and amplitude G_ν ,

$$P_\nu(0, t) = G_\nu \exp \left(-\frac{t^2}{\tau^2} \right). \quad (9)$$

The coefficient k_1 is given by

$$k_1 = \sum_{\nu=1}^N G_\nu B_{\nu_0}^{(1)}. \quad (10)$$

$B_{\nu_o}^{(1)}$ and $\alpha_o^{(1)}$ are defined as the first eigenvector and eigenvalue of an eigenvalue problem the details of which can be found in Ref. 2. The parameter $\alpha_2^{(1)}$ appearing in (8) is the second-order perturbation of the eigenvalue:

$$\alpha_2^{(1)} = \sum_{j=2}^N \frac{\left\{ \sum_{\nu=1}^N \left(\frac{1}{v_\nu} - \frac{1}{v} \right) B_{\nu_o}^{(1)} B_{\nu_o}^{(j)} \right\}^2}{\alpha_o^{(j)} - \alpha_o^{(1)}}. \tag{11}$$

The superscript j identifies $B_{\nu_o}^{(j)}$ and $\alpha_o^{(j)}$ as the j th eigenvector and eigenvalue of the eigenvalue problem; v is the average group velocity. It is convenient to use the parameter ($\tau \rightarrow 0$ is assumed)

$$R = \frac{\Delta t}{\Delta T} = \frac{4 \sqrt{\alpha_2^{(1)}}}{\left(\frac{1}{v_N} - \frac{1}{v_1} \right) \sqrt{L}} \tag{12}$$

that was introduced in Ref. 2 as a measure of the improvement of the pulse distortion of coupled modes compared to the uncoupled case for a guide of length L . ΔT is the length in time covered by the signal arriving in the many uncoupled modes traveling with different group velocities. It is desirable to make R as small as possible by means of coupling between the guided modes.

Finally, we quote the formulas for the power loss coefficients. From Ref. 1 we obtain for the second-order loss attributable to the second-order grating lobe

$$\alpha = \frac{a^4 k^2 (n_1^2 - n_2^2) k^2}{32 \beta_o (\sigma_{-2} + \rho_{-2}) d} [4\sigma_{-1}^2 + \sigma_{-2} \rho_{-2} + (\gamma - 2\gamma_{-1})^2]. \tag{13}$$

In addition to the parameters already defined earlier, we have

$$\begin{aligned} \beta_o &= \text{propagation constant of the guided mode,} \\ f(z) &= a \sin \phi z, \text{ core-cladding interface distortion,} \end{aligned} \tag{14}$$

$$d = \text{slab half width,} \tag{15a}$$

$$\kappa = (n_1^2 k^2 - \beta_o^2)^{\frac{1}{2}}, \tag{15b}$$

$$\sigma_{-1} = [n_1^2 k^2 - (\beta_o - \phi)^2]^{\frac{1}{2}}, \tag{15c}$$

$$\sigma_{-2} = [n_1^2 k^2 - (\beta_o - 2\phi)^2]^{\frac{1}{2}}, \tag{15d}$$

$$\gamma_{-1} = [(\beta_o - \phi)^2 - n_2^2 k^2]^{\frac{1}{2}}, \tag{15e}$$

$$\rho_{-2} = [n_2^2 k^2 - (\beta_o - 2\phi)^2]^{\frac{1}{2}}, \tag{15f}$$

$$\gamma = [\beta_o^2 - n_2^2 k^2]^{\frac{1}{2}}. \tag{15f}$$

The mode power loss coefficient α_s for nonresonant coupling is derived in the Appendix.

$$\alpha_s = \frac{1}{2} \sum_{\substack{\mu=1 \\ \mu \neq s}}^N \frac{|K_{\mu s}|^2 a^2}{(\alpha_\mu/2)^2 + (\beta_\mu - \beta_s - \phi)^2} \alpha_\mu. \quad (16)$$

The power loss coefficient of the μ th mode is α_μ , $|K_{\mu\nu}|^2$ is the factor of $F(\phi)$ in (4) and a is defined in (17).

III. POWER SPECTRUM OF A SINE WAVE WITH RANDOM PHASE

The second-order radiation loss formula and the loss formula for nonresonant coupling to lossy modes were derived for sinusoidally deformed core-cladding interfaces. We expect that these results remain valid, at least approximately, even if the sinusoidal interface variation has a random phase. This random phase assumption is important to ensure the validity of certain averaging procedures that were involved in deriving equation (13). A purely sinusoidal interface variation with constant phase is not likely to occur in practice. For this reason, we derive the relation between the amplitude a of the sinusoidal interface variation, and the power spectrum $F(\phi)$ of the Fourier spectrum of the sinusoidal process with random phase. We are using the idea of a sinusoidal core-cladding interface distortion and the concept of a flat power spectrum of this function with a definite cutoff "frequency" as though they were compatible with each other. It appears possible that a suitable probability distribution for the random phase of the sinusoidal process could be found that would approximate the desired flat power spectrum. However, we make no effort to investigate the compatibility of these ideas, and use them simultaneously in order to gain an order of magnitude estimate of the loss penalty from higher-order loss process that results from intentional, ideal coupling between the guided modes.

In order to establish the desired relation between the amplitude a of the sinusoidal function $f(z)$, describing the core-cladding interface irregularity and the power spectrum of this function, we introduce

$$f(z) = a \sin [\phi z + \psi(z)], \quad (17)$$

with the mechanical frequency ϕ , amplitude a , and random phase $\psi(z)$. The power spectrum is related to $\langle f^2 \rangle$ by the following equation

$$\langle a^2 \sin^2 (\phi z + \psi) \rangle = \frac{1}{\pi} \int_0^\infty F(\phi') d\phi'. \quad (18)$$

The symbol $\langle \rangle$ indicates an ensemble average. We assume (without justification) that the power spectrum has the shape

$$F(\phi') = \begin{cases} \hat{F} = \text{const} & \text{for } 0 \leq \phi' \leq \phi_c \\ 0 & \text{for } \phi' > \phi_c. \end{cases} \quad (19)$$

Since the ensemble average of the square of the sine function is 1/2, we obtain from (18) and (19)

$$\hat{F} = \frac{\pi a^2}{2\phi_c}. \quad (20)$$

The assumption of the flat power spectrum with cutoff, (19), ensures that to first order of perturbation theory no power loss occurs provided that ϕ_c is chosen such that

$$|\beta_{n-1} - \beta_n| > \phi_c. \quad (21)$$

All modes with mode number $\nu < n$ are coupled to each other, while no first-order coupling to the guided modes $N \geq \nu \geq n$ and to the radiation modes is possible. The residual losses that still exist are thus caused by the higher-order processes that are the object of our study.

We are using $\phi = \phi_c$ in the eqs. (15) through (16).

IV. DISCUSSION OF THE EFFECT OF SECOND-ORDER LOSS

We are now ready to calculate the loss penalty that has to be paid for coupling the guided modes with a coupling function the power spectrum of which is given by (19). We assume that the loss mechanism is second-order coupling of mode N (the highest-order guided mode) to the continuous spectrum of radiation modes. We are ignoring the fact that the highest-order mode is usually coupled to the radiation modes by means of a first-order process. However, our assumption is not unrealistic since we can regard mode N as the last of the guided modes that is still coupled to all the other modes, but which is not the mode nearest to the continuous spectrum of radiation modes. If, for example, mode N is to be taken as being the next to last guided mode, it need not be coupled to the last mode by a first-order process, but can itself be coupled to the continuous spectrum of radiation modes by means of second-order coupling. We are thus using the loss coefficient of equation (13) for α_N appearing in (3) while setting $\alpha_\nu = 0$ for $\nu \neq N$.

Figure 1 shows the loss penalty for the 3, 5, 10 and 20 mode case. Since our model is a slab waveguide, the guided modes are the TE modes of a slab. Both core-cladding interfaces of the slab are considered

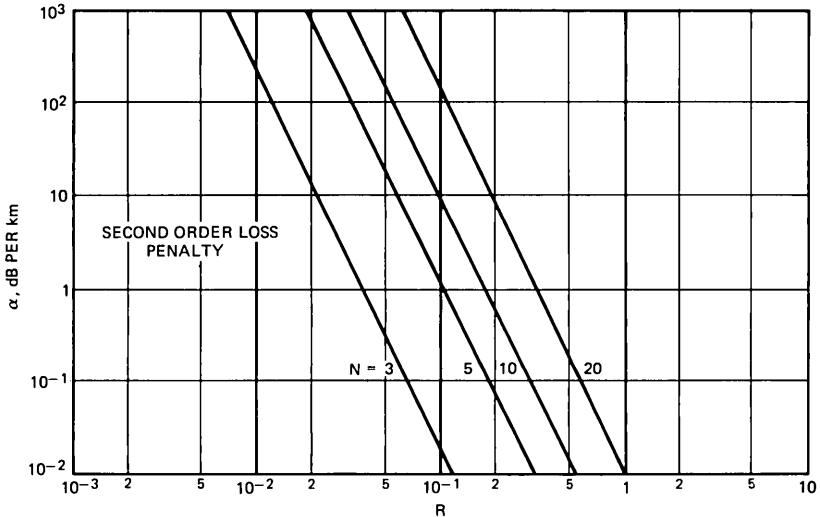


Fig. 1—Loss penalty caused by second-order radiation losses of mode N . α is the power loss coefficient. R is the improvement factor (ratio of pulse width of coupled modes to pulse width of uncoupled modes). N is the number of modes.

to be distorted with the power spectrum of the distortion function given by (19) and (20). We assume in our model that the index ratio of core-to-cladding index is $n_1/n_2 = 1.01$ with $n_1 = 1.5$. The values of kd are

$$kd = 16.5 \quad \text{for 3 modes}$$

$$kd = 35 \quad \text{for 5 modes}$$

$$kd = 70 \quad \text{for 10 modes}$$

$$kd = 145 \quad \text{for 20 modes.}$$

The loss is the steady-state loss per kilometer. We thus assume implicitly that the steady-state distribution is reached, and that the loss is the decrease in power of the steady-state power distribution. (See Ref. 4 for an explanation.) The steady-state loss is plotted as a function of the improvement factor R defined by (12). $R = 0.1$, for example, means that the width of the pulse carried by the coupled guided modes is ten times narrower than it would be in the absence of coupling. The loss penalty increases rapidly with the number of modes. If the third mode of a total of three modes is coupled by the second-order process to the radiation field, an improvement by ten, $R = 0.1$, causes very little radiation loss. However, we see from Fig. 1 that the loss penalty

for the 20-mode case is already more than 100 dB/km. This shows that even the losses caused by second-order coupling of the highest-order mode to the radiation field can cause intolerably high losses if the coupling between the guided modes is strong enough for R to reach $R = 0.1$.

However, to keep the proper perspective, it is important to note that the loss caused by this second-order mechanism would be reduced to zero simply by uncoupling the highest-order mode, and restricting the coupling between the guided modes [by reducing the width of the spectral distribution (19)] to mode 1 through $N - 1$. There are still other losses to contend with. One of these mechanisms will be discussed in the next section. However, second-order losses can be rendered harmless by this device.

It is of interest to know how large an amplitude of the sinusoidal core-cladding interface distortion with random phase is required to cause a given improvement factor R . This question is answered by Fig. 2. The curves of this figure extend below the value $R = 1$, since values of $R > 1$ are of no interest. They are also limited to values of $ka < 1$, because for larger values of ka our perturbation theory becomes meaningless. The figure shows clearly that an improvement factor of $R = 0.1$ can only be reached for fairly large values of ka . In the 20-mode case, we find $ka = 1$ for $R = 0.1$ so that we are approaching the limit of applicability of the second-order perturbation theory used to derive the coupling coefficient (4) and the loss coefficient (13).

V. DISCUSSION OF THE EFFECT OF NONRESONANT COUPLING

We are now considering the nonresonant loss mechanism that led to eq. (16). We are using this equation in the following way. We assume

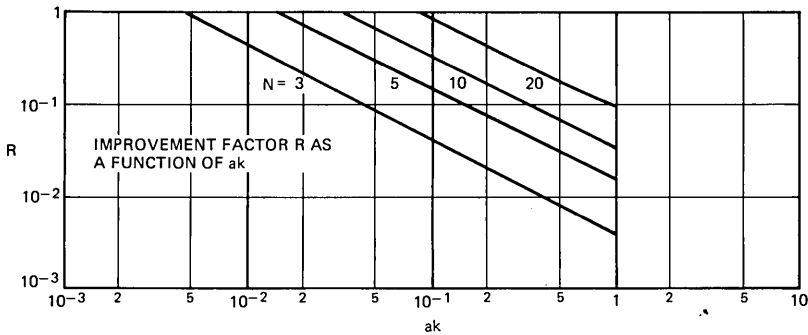


Fig. 2—Improvement factor R as a function of ka . (a = amplitude of sinusoidal core-cladding interface distortion, $k = 2\pi/\lambda_0$).

that the highest-order mode, mode N , is coupled strongly by means of a first-order process to the radiation field. The loss coefficient for this case can be found in equation (63) of Ref. 1. Next, we consider the loss that is transmitted from this high-loss mode to its neighbors. We use eq. (16) to compute the losses of mode $N - 1$, $N - 2$, etc., successively substituting the loss value of each successive iteration to obtain the loss of the next lower mode. We stop at the last mode that is already coupled by a first-order process to the remaining guided modes. The loss penalty that results from coupling this mode to all the other guided modes is being considered here.

Figure 3 shows the loss coefficients for the 10-mode case. The curve on the extreme left is the loss coefficient of mode 10 that loses power via the first-order process to the radiation modes. The modes labeled $s = 9, 8, 7$, etc., suffer loss because of nonresonant coupling to mode 10. The different slopes of these two sets of curves is caused by the fact that the first-order loss process is proportional to $(ak)^2$ while the nonresonant losses are proportional to $(ak)^4$. For a given value of ak , the losses decrease rapidly with decreasing value of s . However, it is surprising how high the losses caused by nonresonant coupling are if $ak = 1$. Figure 4 shows the same data for the 20-mode case.

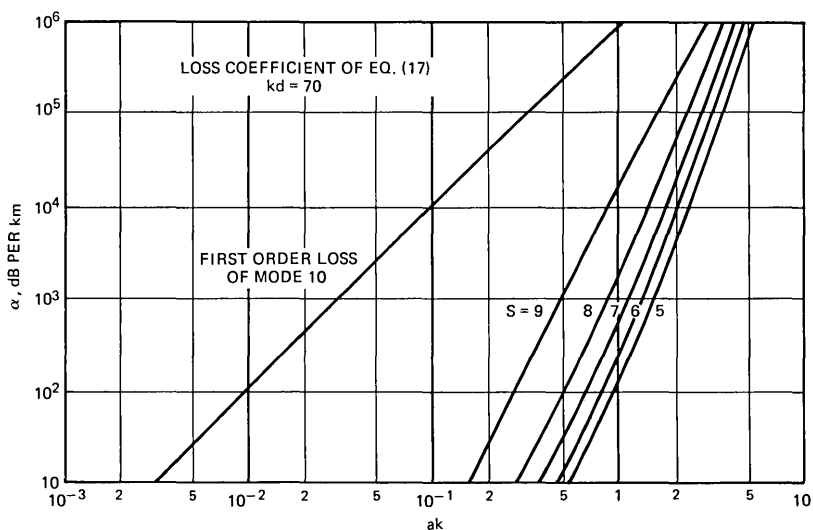


Fig. 3—Mode loss as a function of ak . s is the mode number. The loss is caused by nonresonant coupling of the modes $N - 1$, through mode s to the lossy mode N . The first-order loss of mode N ($N = 10$) is the curve on the left of the figure.

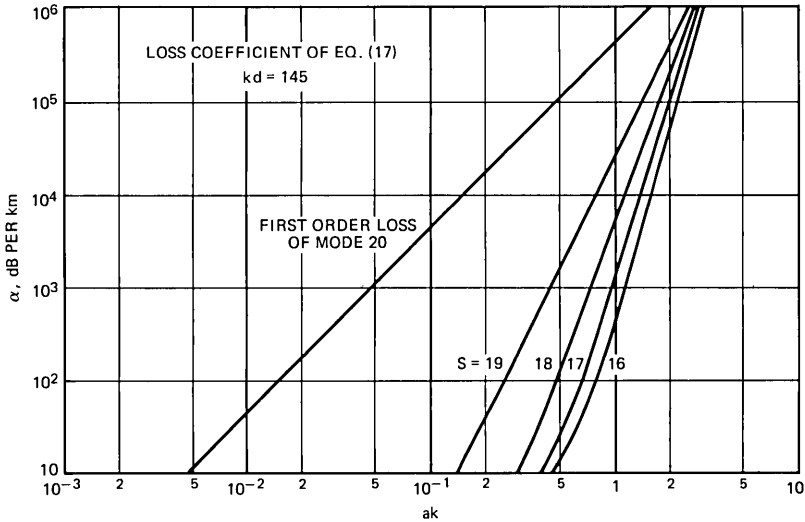


Fig. 4—Same as Fig. 3, with $N = 20$.

Figure 5 shows the loss penalty that results from nonresonant coupling of mode $N - 1$ to the lossy mode N , while modes 1 through $N - 1$ are coupled to each other by the resonant first-order process. The loss penalty is again plotted as a function of the improvement factor R . Figure 5 shows an interesting phenomenon. Whereas the curves for $N = 3$ and $N = 5$ are straight, the curves for $N = 10$ and $N = 20$ are bent. The reason for this difference in behavior can be explained if we consider the shape of the steady-state distribution of mode power P , versus mode number ν . All along the curves for $N = 3$ and $N = 5$, the steady-state power distribution is flat; that means we have equal power in all the modes. On that portion of the curve labeled $N = 10$ that is parallel to the curves with $N = 3$ and 5, we find also that equal power is carried by all the modes in the steady state. However, when the curve begins to bend over, we enter a region where the steady-state power distribution begins to change, favoring the lower-order modes. The loss penalty is correspondingly far less in that region than it would be if the original slope of the curve had been maintained. This is not surprising if we consider that only very little power remains, even in the steady state, in the higher-order modes that couple strongly to the lossy mode $N - 1$. By redistributing the steady-state distribution, the multimode waveguide manages to operate with lower losses. We thus

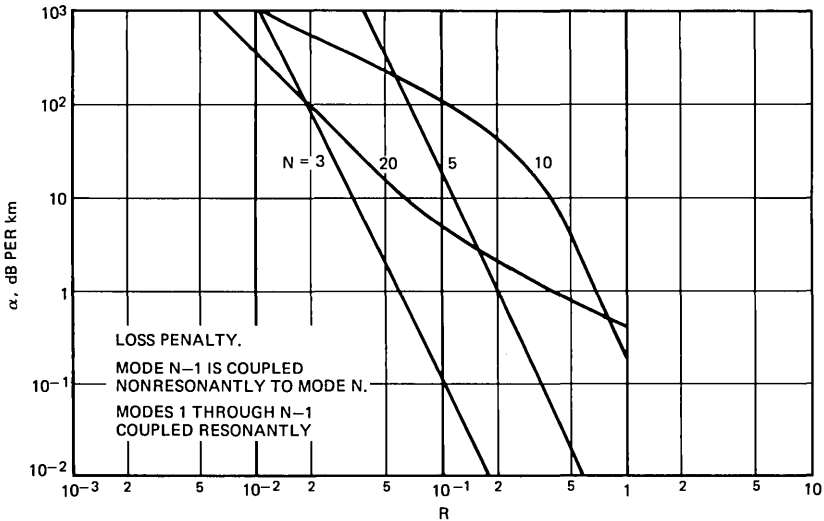


Fig. 5—Loss penalty caused by mode coupling among all $N - 1$ modes and nonresonant coupling of mode $N - 1$ to mode N . The independent variable is the improvement factor R .

find the paradox that, for equal values of R , the 5-mode guide can be lossier than the 10- and 20-mode guide. However, the improvement in the value of R is obtained not by stronger coupling of all the guided modes, but primarily by a reduction in the number of modes that still carry power.

The remaining figures show what happens if we couple fewer guided modes to each other allowing the higher-order modes to die out due to radiation losses. Figure 6 shows the 5-mode case with 4 and 3 guided modes coupled to each other. The improvement in the loss penalty that results from dropping mode 4 from the set of coupled guided modes is substantial.

The same behavior is shown for the 10-mode case in Fig. 7. Again it is apparent how much improvement in the loss penalty can be gained by dropping successively the higher-order modes from the set of coupled guided modes. Only the curve with $n = 9$ behaves anomalously. The change in slope can again be explained by the change in the steady-state distribution. The region with gentler slope corresponds to a steady-state distribution that no longer carries equal power in all the modes but favors the lower-order modes.

This tendency to flip from a steady-state distribution with equal

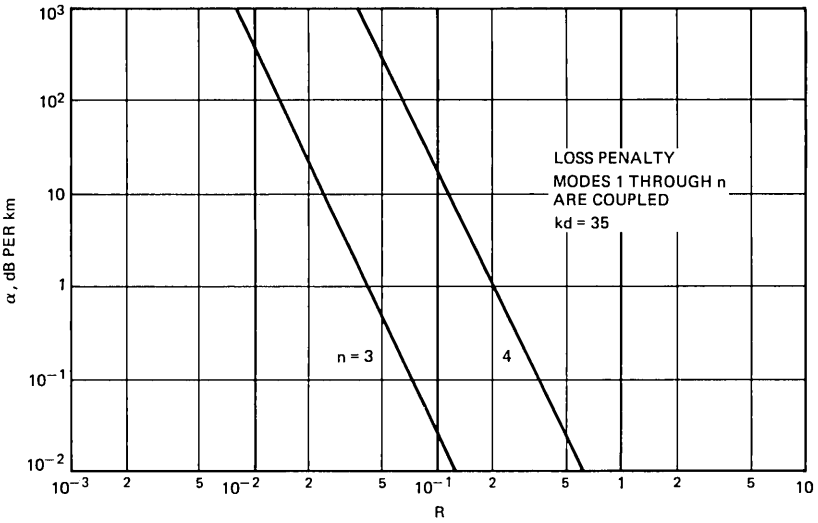


Fig. 6—Loss penalty caused by mode coupling of modes 1 through n and non-resonant coupling of modes $n, n + 1$, etc. to mode N . $N = 5$.

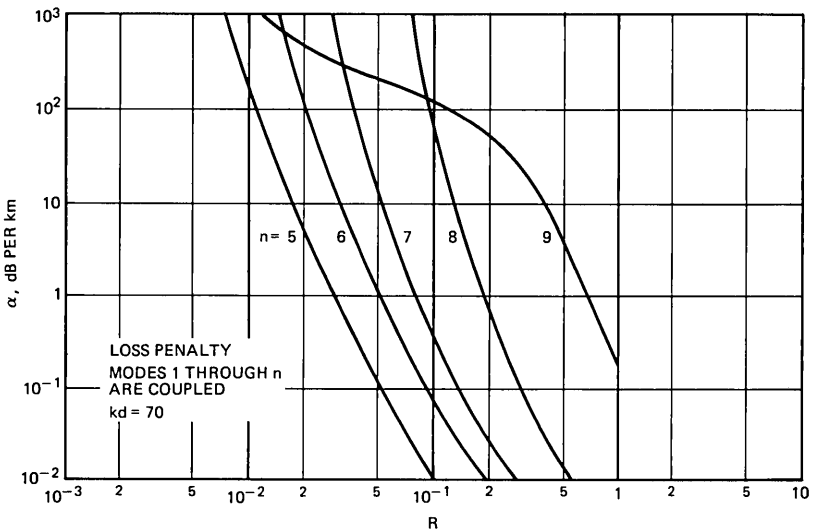


Fig. 7—Same as Fig. 6 for $N = 10$. The curve labeled $n = 9$ departs from the other curves because of a change in the steady-state power distribution.

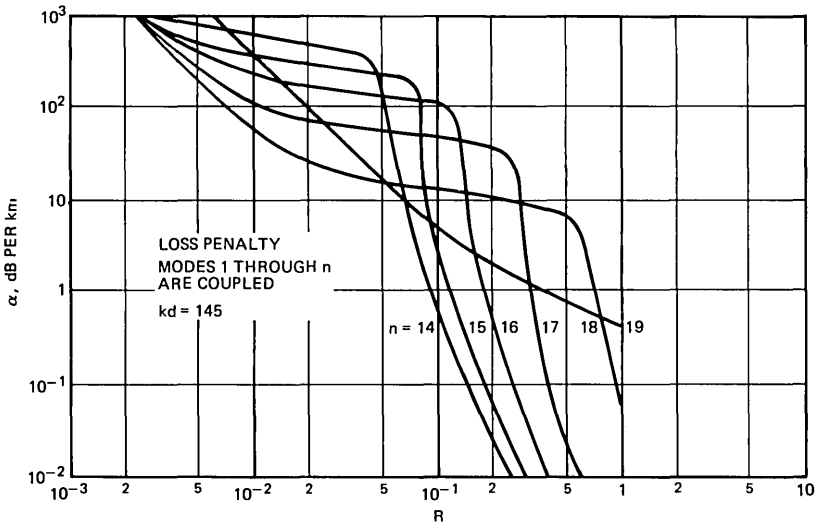


Fig. 8--Same as Figs. 6 and 7 with $N = 20$.

power in all the modes to one that favors lower-order modes is even more apparent in the 20-mode case shown in Figure 8. We also see in this figure that if we want to operate with an improvement factor of $R = 0.1$, and tolerate a loss of 1 dB/km we must uncouple 6 modes from the total of 20 modes allowing only the lowest 14 modes to couple among each other.

VI. CONCLUSIONS

We have studied the loss penalty that results from higher-order loss processes. We have considered two different cases. In both cases, we let most of the guided modes be coupled by a first-order resonant process. In the first case, we assumed that the highest-order mode is coupled to the radiation modes only by means of a second-order process. High losses can still result if we want to achieve a good pulse spreading reduction by means of strong coupling of the guided modes. The loss penalty increases very rapidly with increasing mode number for a fixed value of the improvement factor R . However, by limiting the coupling to one less guided mode, allowing the highest-order guided mode (or more accurately the two highest-order guided modes) to die out, the loss penalty from this second-order process disappears.

There are other processes that still cause a loss penalty even if we drop the two highest-order modes. The lossy modes impart some of their loss to their neighbors via a nonresonant coupling process. The loss penalty from this mechanism can still be high. Again it helps to limit the coupling to fewer of the lower-order modes by reducing the width of the power spectrum of the coupling function. By proper design of the coupling process, the loss penalty for a given improvement factor can be kept in tolerable limits. By uncoupling some of the higher-order modes, we pay an additional loss penalty in the transient before the steady-state distribution has established itself, provided that all modes are excited equally at the beginning of the waveguide.

Our results were obtained by using the model of the slab waveguide. However, they allow an estimate of the performance of the round optical fiber if we keep in mind that the total number of modes of the round fiber is the square of the mode number of the slab waveguide. The 10-mode case of the slab waveguide thus corresponds to a 100-mode round optical fiber. The members within each family of modes with equal circumferential field distribution (the same value of ν in $\cos \nu\phi$) are coupled among each other by diameter changes of the fiber core. Each family of this kind behaves similarly to the modes of the slab waveguide studied here. We are thus able to use the results of the slab waveguide to draw conclusions about the expected behavior of round optical fibers.

APPENDIX

Derivation of Equations (16)

As a starting point we use the coupled wave equations.⁵

$$\frac{da_\nu}{dz} = -i\bar{\beta}_\nu a_\nu + \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^N c_{\nu\mu}(z)a_\mu. \quad (22)$$

The propagation constants $\bar{\beta}_\nu$ are assumed to be complex quantities, $c_{\nu\mu}$ are the coupling coefficients and a_ν are the mode amplitudes. With the slowly varying mode amplitudes A_ν defined by

$$a_\nu = A_\nu e^{-i\bar{\beta}_\nu z} \quad (23)$$

the system of coupled wave equations is transformed into the following form

$$\frac{dA_\nu}{dz} = \sum_{\substack{\mu=1 \\ \mu \neq \nu}}^N c_{\nu\mu}(z)A_\mu e^{i(\bar{\beta}_\nu - \bar{\beta}_\mu)z}. \quad (24)$$

We now assume that only one of the modes, mode s , is strongly excited at $z = 0$, while all the other mode amplitudes vanish initially

$$A_\nu(0) = 0 \quad \nu \neq s. \quad (25)$$

In the vicinity of $z = 0$, we thus obtain approximately

$$\frac{dA_\nu}{dz} = c_{\nu s}(z)A_s e^{i(\bar{\beta}_\nu - \bar{\beta}_s)z} \quad \text{for } \nu \neq s \quad (26)$$

and for $\nu = s$ (since $\bar{\beta}_s$ is assumed to be real we write $\bar{\beta}_s = \beta_s$)

$$\frac{dA_s}{ds} = \sum_{\substack{\mu=1 \\ \mu \neq s}}^N c_{s\mu}(z)A_\mu e^{i(\beta_s - \bar{\beta}_\mu)z}. \quad (27)$$

In analogy with the Wigner-Weisskopf method,⁶ we next assume that the z dependence of the mode amplitude A_s is given by

$$A_s(z) = A_s(0)e^{-\frac{1}{2}\alpha_s z}. \quad (28)$$

The determination of the unknown constant α_s is the objective of the following calculation. Substitution of (28) into (26) and subsequent integration results in

$$A_\mu(z) = A_s(0) \int_0^z c_{\mu s}(x) e^{[i(\bar{\beta}_\mu - \beta_s) - (\alpha_s/2)]x} dx. \quad (29)$$

At this point it becomes necessary to specify the z dependence of the coupling function $c_{\nu\mu}(z)$. We want to determine the effect of nonresonant coupling but are, nevertheless, interested in the influence of a periodic coupling function. For simplicity it is convenient to assume that the coupling coefficients are of the following form

$$c_{\mu s}(z) = \frac{a}{\sqrt{2}} K_{\mu s} e^{-i\phi z}. \quad (30)$$

It is a well-established fact that the coupling coefficient can be decomposed into a constant part $K_{\nu\mu}$ times a function of z . If mode coupling is caused by core-cladding interface irregularities of dielectric waveguides, the z -dependent function describes directly the shape of the core-cladding interface deformation.⁴ Ordinarily, we would expect to see a sine or a cosine function instead of the exponential function that appears in (30) provided that the core-cladding interface distortion is purely sinusoidal. But a sinusoidal function can always be decomposed into two exponential functions. We keep only one of these two terms. This approximation is justified if we consider near-resonant coupling. Only terms with small values of $\bar{\beta}_\nu - \bar{\beta}_\mu - \phi$ will be seen to give a sub-

stantial contribution. The term $\bar{\beta}_\mu - \beta_\mu + \phi$, that would result from the neglected part of the sinusoidal function, is large and therefore makes only a slight contribution to the coupling process.

With the help of (30), we obtain from (29)

$$A_\mu(z) = \frac{a}{\sqrt{2}} K_{\mu s} A_s(0) \frac{e^{[i(\bar{\beta}_\mu - \beta_s - \phi) - (\alpha_s/2)]z} - 1}{i(\bar{\beta}_\mu - \beta_s - \phi) - \frac{\alpha_s}{2}}. \tag{31}$$

It can be shown that for lossless guides the relation $c_{\nu\mu} = -c_{\mu\nu}^*$ is required.⁷ We use this relation in our present case since it must be approximately true even for lossy guides. We then obtain from (30)

$$c_{s\mu} = -\frac{a}{\sqrt{2}} K_{\mu s}^* e^{i\phi z}. \tag{32}$$

Substitution of (28), (31) and (32) into (27) yields

$$\alpha_s = \sum_{\substack{\mu=1 \\ \mu \neq s}}^N |K_{\mu s}|^2 \frac{1 - e^{[i(\beta_s - \bar{\beta}_\mu + \phi) + (\alpha_s/2)]z}}{i(\bar{\beta}_\mu - \beta_s - \phi) - \frac{\alpha_s}{2}} a^2. \tag{33}$$

In order to proceed further, we assume that mode s was inherently lossless prior to being coupled to the other modes. We also assume that the losses of the remaining modes are high. Since it appears reasonable to expect that α_s must be smaller than any of the loss coefficients of the other modes (these loss coefficients are the imaginary parts of $\bar{\beta}_\mu$) we can neglect the exponential term in (33) for large values of z so that we obtain

$$\alpha_s = \sum_{\substack{\mu=1 \\ \mu \neq s}}^N \frac{|K_{\mu s}|^2 a^2}{i(\bar{\beta}_\mu - \beta_s - \phi) - \frac{\alpha_s}{2}}. \tag{34}$$

The coefficient α_s is a complex quantity. Its imaginary part contributes only a slight change to the propagation constant of β_s . We are interested only in its real part. We write

$$\bar{\beta}_\mu = \beta_\mu - i \frac{\alpha_\mu}{2}. \tag{35}$$

In the denominator of (34), we neglect α_s compared to α_μ and obtain finally for the real part of α_s (which we write again α_s for simplicity)

$$\alpha_s = \frac{1}{2} \sum_{\substack{\mu=1 \\ \mu \neq s}}^N \frac{|K_{\mu s}|^2 a^2}{(\alpha_\mu/2)^2 + (\beta_\mu - \beta_s - \phi)^2} \alpha_\mu. \tag{36}$$

Equation (36) is the desired approximation. We assume that

$$|\beta_\mu - \beta_s - \phi| \ll \phi, \quad (37)$$

but require $\beta_\mu - \beta_s \neq \phi$ in the spirit of the nonresonant coupling assumption. Even for high loss modes we assume that the following relation applies

$$\alpha_\mu \ll \phi. \quad (38)$$

It is apparent that replacement of $\beta_\mu - \beta_s - \phi$ in the denominator of (36) with $\beta_\mu - \beta_s + \phi$ would lead to much smaller values of α_s . If we had used the sine or cosine function instead of the exponential function in (30), we would have obtained an additional term with $\beta_\mu - \beta_s + \phi$ (in the denominator) in (36). This additional term is much smaller than the leading term of α_s , so that our approximation (30) appears justified.

REFERENCES

1. Marcuse, D., "Higher-Order Scattering Losses in Dielectric Waveguides," B.S.T.J., this issue, pp. 1801-1817.
2. Marcuse, D., "Pulse Propagation in Multimode Dielectric Waveguides," B.S.T.J., 51, No. 6 (July-August 1972), pp. 1199-1232.
3. Personick, S. D., "Time Dispersion in Dielectric Waveguides," B.S.T.J., 50, No. 3 (March 1971), pp. 843-859.
4. Marcuse, D., "Power Distribution and Radiation Losses in Multimode Dielectric Slab Waveguides," B.S.T.J., 51, No. 2 (February 1972), pp. 429-454.
5. Miller, S. E., "Coupled Wave Theory and Waveguide Applications," B.S.T.J., 33, No. 3 (May 1954), pp. 661-719.
6. Marcuse, D., *Engineering Quantum Electrodynamics*, New York: Harcourt, Brace and World, 1970, p. 200.
7. Marcuse, D., "Derivation of Coupled Power Equations," B.S.T.J., 51, No. 1 (January 1972), pp. 229-237.

A Linear Phase Modulator for a Short-Hop Microwave Radio System

By S. R. SHAH

(Manuscript received May 16, 1972)

A linear phase modulator is a useful component in a short-hop radio system using digital modulation. Such a modulator has been designed, built, and tested for 4-level operation at a carrier frequency of 300 MHz and for a line rate of 20 megabits.

In this paper, the design and the performance of the modulator are presented. Measurements on the phase modulator show that the performance is in agreement with the theory.

I. INTRODUCTION

A linear phase modulator is a useful component in a short-hop radio system using digital modulation.^{1,2} A possible application of this modulator in a digital radio system is shown in a block diagram in Fig. 1. A shared delta modulator multiplex operating at a line rate of 20 megabits is an example of a digital multiplex which could be used in such a system.³ The output binary signal of this multiplex can be converted into a baseband pulse sequence by a 4-level block coder.⁴

The modulator described in this paper satisfies the requirements for the above application. It is based upon the original Armstrong circuit, which is well suited to large baseband bandwidths and is reasonably linear for low modulation indexes.⁵ An analysis of distortion for the type of baseband signal used in this application is discussed by C. L. Ruthroff and W. F. Bodtmann in Ref. 1.

The output of the phase modulator is at an IF frequency of 300 MHz. It can be converted in a linear mixer to any desired RF frequency in the microwave or millimeter-wave range, and can be amplified for transmission by an injection-locked oscillator amplifier.⁶

II. DESCRIPTION OF THE MODULATOR

A block diagram of the modulator is shown in Fig. 2. A quartz crystal oscillator provides a 300-MHz stable carrier frequency. The baseband

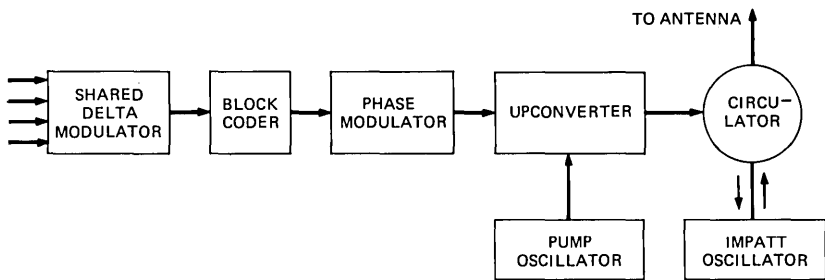


Fig. 1—Block diagram of a solid state transmitter for a digital radio system.

signal is modulated in a double-sideband suppressed-carrier amplitude modulator. At the output of this modulator another carrier, 90 percent out of phase with the first, is added to the sidebands. The combined low-index phase-modulated signal then passes through an amplifier and a times-four frequency multiplier. The output signal from the times-four multiplier is converted in a double-sideband balanced mixer to the original carrier frequency of 300 MHz as described in Ref. 1. The unwanted output frequencies from the mixer are eliminated by a lowpass filter, and the required phase-modulated signal is obtained. A limiter is not used as indicated in Ref. 1 because sufficient amplitude compression occurs in the harmonic generator.

2.1 Double-Sideband Suppressed-Carrier Amplitude Modulator

Carrier suppression by a double-sideband balanced mixer is a well known technique.⁷ By using a conventional mixer, an isolation of about 30 dB between any two ports can be obtained. The spectrum of the

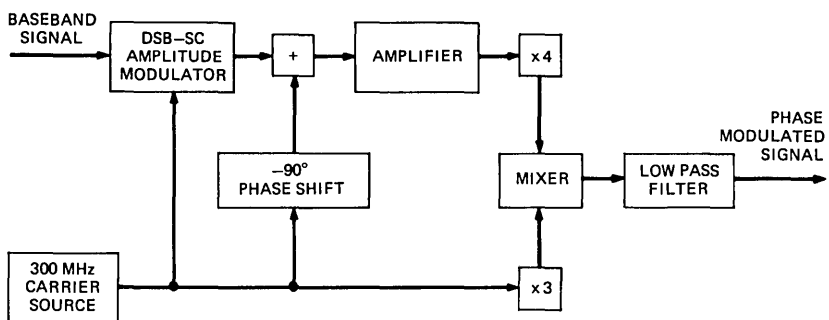


Fig. 2—Block diagram of linear phase modulator.

double-sideband suppressed-carrier amplitude-modulated signal with a modulation index of 0.2 is shown in Fig. 3. The suppression of the carrier is not sufficient, and could cause distortion, particularly if the quadrature carrier is to be reintroduced at the receiving terminal.⁸ This can be explained best by considering the phasor diagrams of Fig. 4. Figure 4a shows the ideal case, in which the carrier is totally suppressed and the quadrature carrier is added without causing any distortion. Figure 4b shows a case in which the carrier is not sufficiently suppressed, and the remaining component of the carrier changes the phase of the quadrature carrier resulting in a phase error. To suppress the remaining component of the carrier, a signal of equal amplitude and 180 degrees out of phase with the carrier is added to the output signal of the mixer. This is illustrated in a block diagram in Fig. 5. The carrier from the local oscillator is divided into two parts, one of which goes to the mixer, and the other, after undergoing the required changes in amplitude and phase, is added to the output signal from the mixer. The spectrum of the resultant signal from the adder is shown in Fig. 6. Note that the carrier is suppressed by 58 dB. The carrier suppression is insensitive to changes in oscillator level; a change in the amplitude level of the crystal oscillator signal by 1 dB results in a change in the amplitude level of the carrier by a fraction of a decibel.

2.2 *Frequency Multipliers*

2.2.1 *Times-Four Multiplier*

A maximum peak deviation of $\pm 3\pi/4$ radians is required in a 4-level PSK system. If the modulator output is multiplied by four, the required peak deviation in the modulator before multiplication is then $3\pi/16$

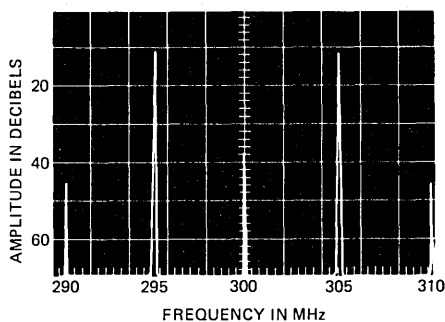


Fig. 3—Carrier suppression of a double-sideband balanced mixer.

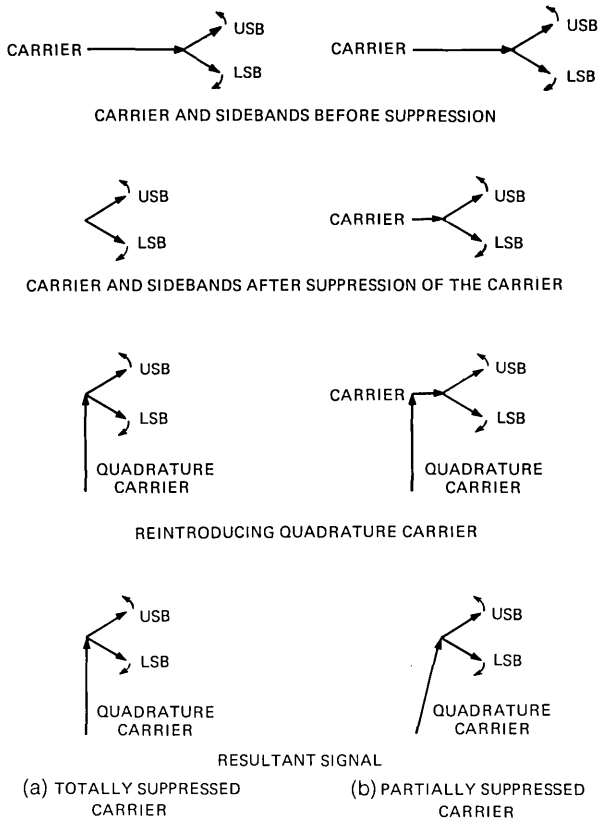


Fig. 4—Phasor diagram of the phase-modulated signal. (a) Totally suppressed carrier. (b) Partially suppressed carrier.

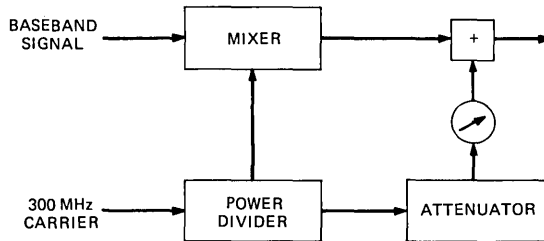


Fig. 5—Block diagram of double-sideband suppressed-carrier amplitude modulator.

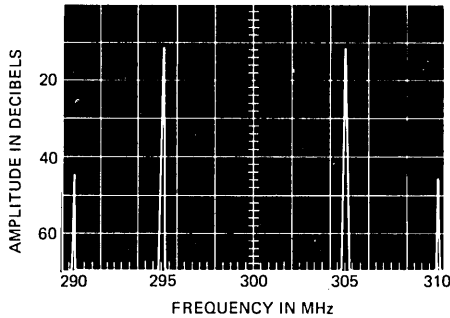


Fig. 6—Carrier suppression by improved double-sideband suppressed-carrier amplitude modulator.

radians. To obtain this magnitude of deviation, a resistive multiplier circuit is used. The conventionally designed circuit employs two Schottky barrier diodes. The input and the output matching sections are five-element 0.1-dB-ripple Tschebyscheff filters.⁹ Variable air trimmer capacitors and hand-wound coils of No. 14 bare copper wire are the elements of the filters and the idler networks.

The performance of the multiplier is measured by applying the phase-modulated sine-wave with a carrier frequency of 300 MHz. The frequency of the baseband signal is varied from 5 MHz to 25 MHz in 5-MHz steps. The peak phase deviation is kept constant at 0.1 radian. The output spectrum of the multiplier corresponding to each baseband multiplier signal frequency is recorded on the same photograph, which is shown in Fig. 7. Ideally, the conversion loss of the multiplier should be the same for each baseband signal frequency. The frequency response of the

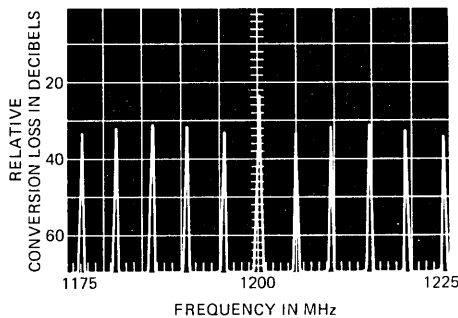


Fig. 7—Frequency response of the times-four multiplier.

mixer used in the double-sideband suppressed-carrier amplitude modulator is not flat for these baseband signal frequencies; this is shown in Fig. 8. The frequency response of the mixer is reflected in the output of the multiplier shown in Fig. 7.

2.2.2 *Times-Three Multiplier*

A times-three multiplier is used to convert the output carrier frequency of the times-four multiplier to the frequency of the carrier source as shown in Fig. 2. This multiplier need not have a broad bandwidth. The technique used for the design and fabrication of the times-three multiplier is the same as that used for the times-four multiplier.

2.3 *90-Degree Phase Shifter*

An accurate 90-degree phase shifter is made from twisted wire distributed elements using the method described in Ref. 10. The coupler, which has a crossover frequency of 300 MHz, is made of twisted pairs of FORMEX wire separated from the ground plane by polyethylene dielectric.

2.4 *Lowpass Filter*

A three-element 0.1-dB-ripple Tschebyscheff filter is used to eliminate the unwanted output frequencies from the mixer. Two variable air trimmer capacitors and a hand-wound coil of No. 14 copper wire are the elements of this lowpass filter.

III. MEASURED PERFORMANCE

The performance of the phase modulator was measured for sine-wave and square-wave baseband signals with various output phase deviations from 0.1 to $3\pi/4$ radians. For any specific input phase deviation, the

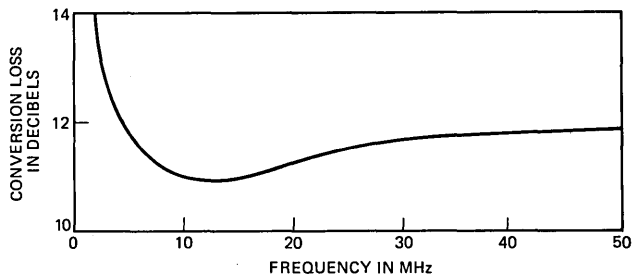


Fig. 8—Frequency response of the mixer.

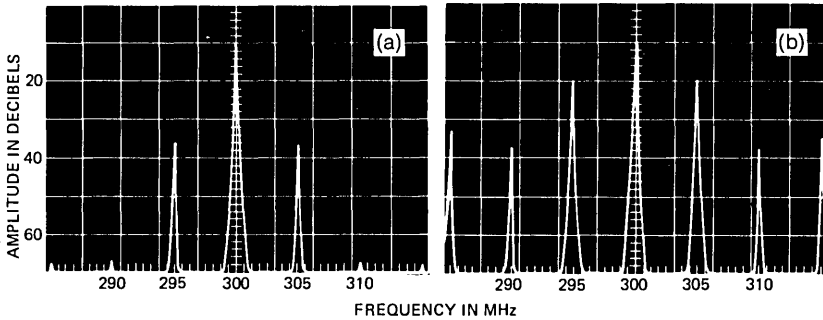


Fig. 9—Output spectrum of the phase-modulated signal from the adder. (a) 0.1 radian phase deviation. (b) $3\pi/16$ radian phase deviation.

amplitude levels of the carrier and the sidebands of the phase-modulated signal can be computed (see Appendix A); the spectrum of the output signal from the adder can then be inspected to verify the input phase deviation as shown in Fig. 9. When the signal passes through the times-four multiplier, the required output phase deviation is obtained. The output spectra of the modulator for a 5-MHz baseband signal with output phase deviations of 0.4 and $3\pi/4$ radians are shown in Fig. 10.

The detected phase-modulated signal can be compared with the input baseband signal by using a phase detector as shown in a block diagram of Fig. 11. The output signal of the detector is calculated as described in Appendix B; Figure 12 shows the calculated results for sinusoidal baseband signal with the phase deviation of $\pi/4$ and $3\pi/4$ radians. The measured results for 5-MHz and 20-MHz baseband signals with 0.4 and $3\pi/4$ radians phase deviation are shown in Fig. 13.

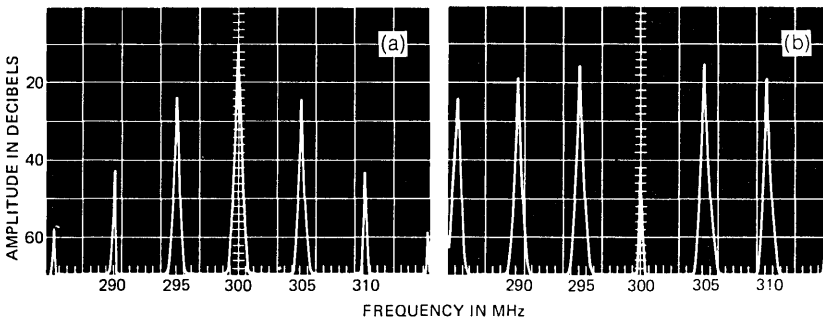


Fig. 10—Output spectrum of the phase modulator for 5 MHz sinusoidal baseband signal. (a) 0.4 radian phase deviation. (b) $3\pi/4$ radians phase deviation.

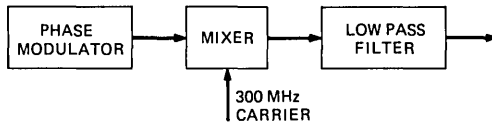


Fig. 11—Block diagram of the phase detector.

IV. CONCLUSIONS

Comparing the calculated and the measured results, particularly of Figs. 12b and 13b, it is clear that the phase modulator is performing as expected. Further, it has been demonstrated that this phase modulator is suitable for a line rate of 20 megabits. The carrier frequency is stable as it is derived from a quartz crystal oscillator. All these aspects make this phase modulator a useful component in a short-hop radio system—especially in coherent phase-shift-keyed PCM systems.

V. ACKNOWLEDGMENT

The author would like to express his thanks to C. L. Ruthroff for his helpful suggestions.

APPENDIX A

The expression for a carrier which is phase-modulated by a sinusoidal signal can be written in the general form:

$$M(t) = A_c \cos(\omega_c t + X_1 \cos \omega_m t). \quad (1)$$

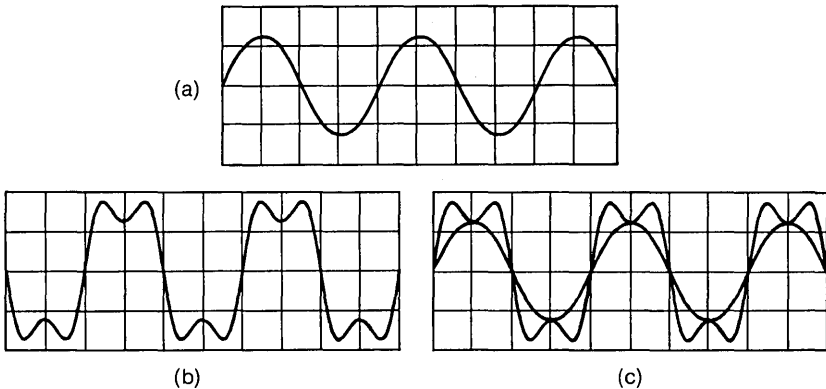


Fig. 12—Calculated output signal from the phase detector. (a) $\pi/4$ radian phase deviation. (b) $3\pi/4$ radians phase deviations. (c) Curve (a) is imposed on curve (b).

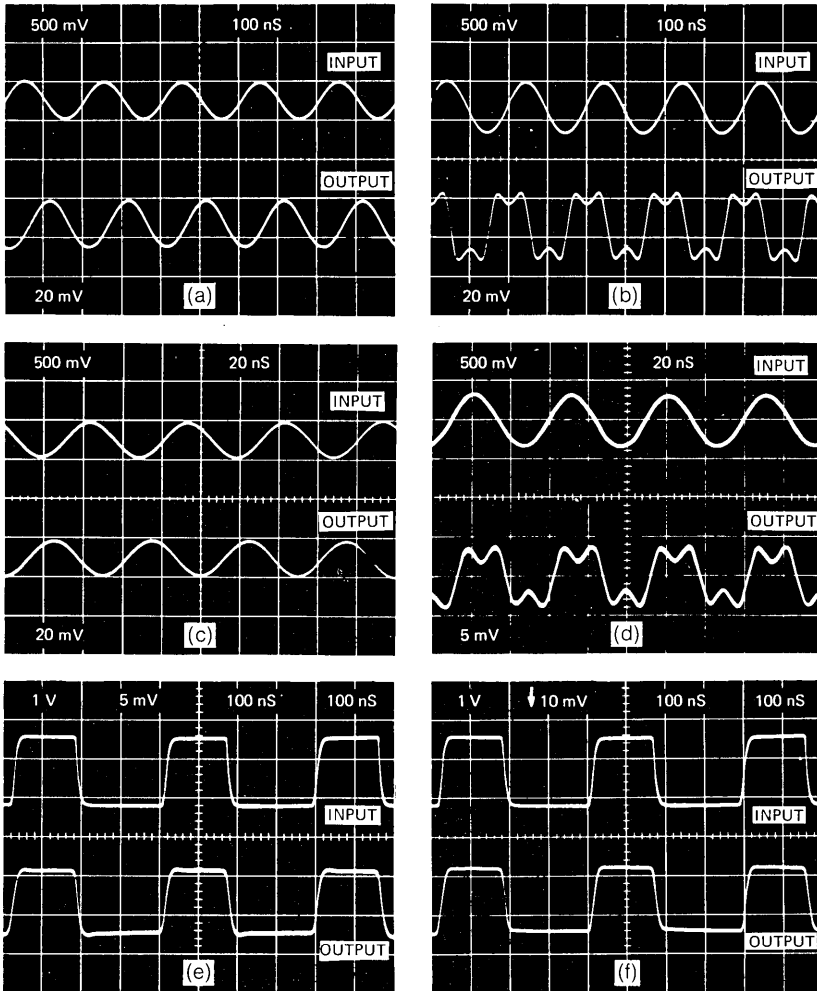


Fig. 13—Photographs of the input baseband signal and the corresponding output signal from the phase detector. (a) 0.4 radian phase deviation. (b) $3\pi/4$ radians phase deviation. (c) 0.4 radian phase deviation. (d) $3\pi/4$ radians phase deviation. (e) 0.4 radians phase deviation. (f) $3\pi/4$ radians phase deviation.

Here, X_1 is the peak phase deviation in radians. Now

$$\begin{aligned}
 M(t) = A_c \left\{ J_0(X_1) \cos \omega_c t + J_1(X_1) \cos \left[(\omega_c + \omega_m)t + \frac{\pi}{2} \right] \right. \\
 + J_1(X_1) \cos \left[(\omega_c - \omega_m)t + \frac{\pi}{2} \right] - J_2(X_1) \cos (\omega_c + 2\omega_m)t \\
 \left. + J_2(X_1) \cos (\omega_c - 2\omega_m)t + \dots \right\}. \quad (2)
 \end{aligned}$$

Equation (2) shows that the magnitudes of the sidebands, relative to the carrier, can be determined by the Bessel coefficients.^{5,8}

For a peak phase deviation of 0.1 radian,

$$J_0(X_1) = 0.9975$$

$$J_1(X_1) = 0.04994$$

$$J_2(X_1) = 0.00125$$

$$\frac{J_0(X_1)}{J_1(X_1)} = 19.97396 \quad \frac{J_0(X_1)}{J_2(X_1)} = 798.$$

Converting in decibels,

$$20 \log \left(\frac{J_0(X_1)}{J_1(X_1)} \right) = 26 \text{ dB}$$

$$20 \log \left(\frac{J_0(X_1)}{J_2(X_1)} \right) = 58 \text{ dB}.$$

In other words, the difference in the energy level between the carrier and the first sidebands is 26 dB, and between the carrier and the second sidebands is 58 dB for 0.1 radian peak phase deviation. This is shown in Fig. 9a.

Similarly, for a peak phase deviation of $3\pi/16$ radians,

$$J_0(X_1) = 0.9149, \quad J_1(X_1) = 0.2823, \quad J_2(X_1) = 0.04226$$

and

$$20 \log \left(\frac{J_0(X_1)}{J_1(X_1)} \right) = 10.2 \text{ dB}$$

$$20 \log \left(\frac{J_0(X_1)}{J_2(X_1)} \right) = 26.65 \text{ dB}.$$

Note that the difference in the energy level between the carrier and the first sidebands is 10.2 dB, and between the carrier and the second

sidebands is 26.65 dB as shown in Fig. 9b. Values of $J_n(x)$ are obtained from Ref. 11.

For the phase deviations of 0.1, 0.4, $3\pi/16$, and $3\pi/4$ radians, the calculated difference in the energy level between the carrier and the sidebands are shown in Table I.

TABLE I—THE CALCULATED DIFFERENCE IN THE ENERGY LEVEL BETWEEN THE CARRIER AND THE SIDEBANDS FOR VARIOUS PHASE DEVIATIONS

Phase deviation in radians	Difference in the energy level between the carrier and the first sidebands in dB	Difference in the energy level between the carrier and the second sidebands in dB
0.1	26	58
0.4	13.8	33.7
$3\pi/16$	10.2	26.7
$3\pi/4$	-26.3	-24.3

APPENDIX B

The phase-modulated signal can be expressed by

$$e_p = \sin [\omega_c t + X_1 \cos \omega_m t]$$

where

$$\omega_c = 2\pi f_c, \quad f_c = \text{carrier frequency}$$

$$\omega_m = 2\pi f_m, \quad f_m = \text{baseband signal frequency}$$

and

X_1 is the peak phase deviation in radians.

In a phase detector, e_p is multiplied by $\cos \omega_c t$, and the low-frequency part of the output is

$$e_o = \sin (X_1 \cos \omega_m t).$$

The values of e_o are calculated with respect to t , and are plotted for the following cases

- (i) $X_1 = \pi/4, f_c = 300 \text{ MHz}$ and $f_m = 5 \text{ MHz}$
- (ii) $X_1 = 3\pi/4, f_c = 300 \text{ MHz}$ and $f_m = 5 \text{ MHz}$.

Figures 12a and b show the curves plotted for case (i) and (ii) respectively. For comparison, the curve of case (i) is imposed on the curve of case (ii), which is shown in Fig. 12c.

REFERENCES

1. Ruthroff, C. L., and Bodtmann, W. F., "A Linear Phase Modulator for Large Baseband Bandwidths," B.S.T.J., 49, No. 8 (October 1970), pp. 1893-1903.
2. Ruthroff, C. L., and Bodtmann, W. F., "Adaptive Coding for Coherent Detection of Digital Phase Modulation," unpublished work.
3. Osborne, T. L., and Michael, S., "Experimental Shared Delta Modulator Multiplex," unpublished work.
4. Bodtmann, W. F., unpublished work.
5. Armstrong, E. H., "A Method of Reducing Disturbances in Radio Signalling by a System of Frequency Modulation," Proc. IRE, 24, No. 5 (May 1936), pp. 689-740.
6. Glance, B., "Low Q Microstrip IMPATT Oscillator at 30 GHz," unpublished work.
7. Caruthers, R. S., "Copper Oxide Modulators in Carrier Telephone Systems," B.S.T.J., 18, No. 2 (April 1939), pp. 315-337.
8. Members of the Technical Staff, *Transmission Systems for Communications*, Bell Telephone Laboratories, Inc., 1970, pp. 96-115, 450-454.
9. Matthaei, G. L., Young, L., and Jones, E. M. T., *Microwave Filters, Impedance Matching Networks and Coupling Structures*, New York: McGraw-Hill, 1964, pp. 83-102, 421-434.
10. Michael, S., "Twisted Wire 3 dB Directional Couplers," unpublished work.
11. Abramowitz, M., and Stegun, I. A., *Handbook of Mathematical Functions*, Washington: National Bureau of Standards, 1964, p. 390.

A Doped Surface Two-Phase CCD

By R. H. KRAMBECK, R. H. WALDEN, and K. A. PICKAR

(Manuscript received May 5, 1972)

The successful operation of an n-channel two-phase charge-coupled device has been achieved. The asymmetry in the surface potential profile necessary to force the charge to move unidirectionally was obtained by ion implanting a nonuniform doping distribution in the Si substrate under each gate. An eight-stage shift register with a length per stage of 80 μm was made, and was operated as both a digital and an analog device. There are two ways to clock the device. Either both clock lines are driven with square waves, out of phase by one half of a period, or one clock is held at a fixed DC potential while the other is driven with a square wave. Using the latter method, the charge transfer efficiency was better than 99.9 percent per transfer over the clock frequency range of 10^3 Hz to 6.5×10^6 Hz.

I. INTRODUCTION

In a charge-coupled device, as described by Boyle and Smith,¹ charge moves successively from the semiconductor region under a given electrode to the region under the next electrode. For information to be transferred from one end of the resulting shift register to the other, it is necessary that the charge always move in the same direction. Until now, the way this directionality has been typically achieved is by the use of three or more clock lines.²⁻⁴ In this type of structure, when charge is transferred from one electrode to the next, the electrode behind the one transferring charge is kept at a potential which repels the free charge and thereby prevents backward flow. The electrode to receive charge, meanwhile, is made more attractive to charge than the one giving up its charge. It can be seen that this arrangement requires three electrodes (at least) for each packet of charge and that each one must be driven by a different clock line.

The use of three clock lines (as opposed to two) has significant topological disadvantages because with three clock lines there must be crossovers. These have been fabricated by a diffusion into the semiconductor surface.² This diffusion must be contacted once for every bit which is undesirable from the standpoint of yield and packing density.

This problem is exacerbated when connection is made to a three-phase shift register in a serpentine layout. Since the serpentine layout typically provides the most compact register, the three-phase device is normally limited to applications in which a straight layout is permissible, as for example, in an imaging device.

The serpentine problem can be overcome if a CCD with four clock lines is employed;^{3,4} however, this device requires two layers of metallization, and fabrication problems may occur due to the present state of development of that technology.

In view of these difficulties, the fabrication of a two-phase CCD is considered an important goal in the development of charge-coupled devices; it requires no crossovers, utilizes only one layer of metallization, and can be easily laid out in a serpentine configuration.

The design, fabrication and operation of a simple eight-bit two-phase shift register will be discussed in the following three sections of this paper. A final section is devoted to conclusions.

II. DEVICE DESIGN

In a two-phase CCD, when an electrode is giving up its charge, both adjacent electrodes are biased to attract that charge. Moreover, they must be equally attractive to the charge because every other electrode is tied together; consequently, directionality must be achieved by associating a potential barrier with each electrode. The directional transfer imposed by such a barrier is shown schematically in Fig. 1. The potential barrier near the left end of the region under each electrode prevents backward flow of charge. While its asymmetrical location under the electrode is important, the shape of the barrier is not critical in determining whether charge flow will be directional or not, but higher speed would be expected if the right side of the barrier is sloped to encourage charge flow in the forward direction.

The desired surface potential barrier can be obtained by implanting into the *p*-type silicon substrate a shallow layer of boron ions in a narrow stripe geometry as shown in Fig. 2. Charge-coupled device operation requires that the entire region under each gate be in deep depletion, so the boron implant must be light enough to be totally ionized and depleted. The resulting negative charge layer assures that, for a given value of applied voltage V_A , the surface potential, ϕ_s , has a lower value in the implanted region than in the unimplanted region. The height of the potential barrier $\Delta\phi_s$ is simply the difference between these values.

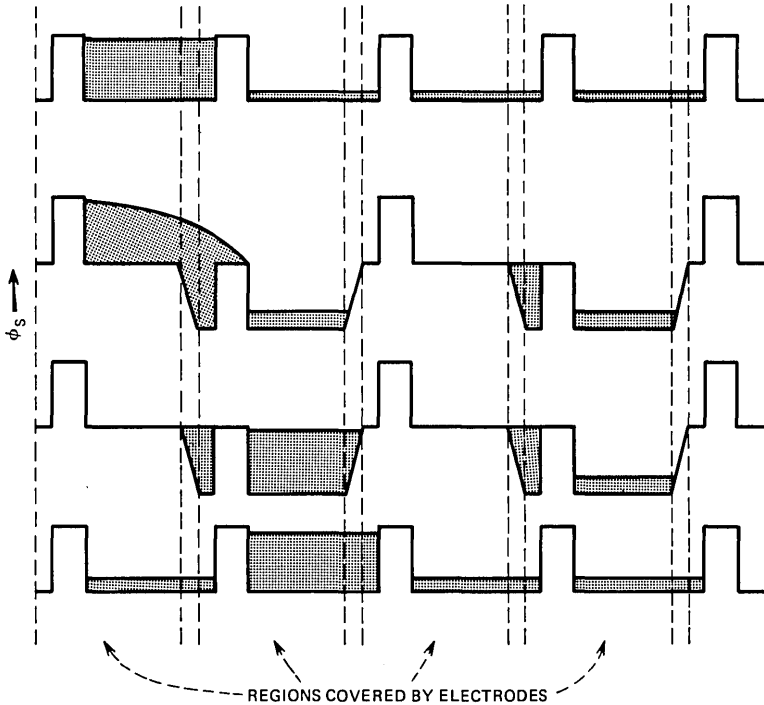


Fig. 1—Sequence of plots of surface potential vs position with free carrier density shown cross hatched. The height of the barrier is typically $\sim 5V$. One complete transfer is shown.

A one-dimensional solution of Poisson's equation can be obtained to give ϕ_s as a function of V_A for the deep depletion condition. It is assumed here that the doping profile is characterized by a constant density N_{A1} highly doped region to a well defined depth x_1 , and the background density, N_{A2} (see Fig. 3). This approximate profile is sufficient to show the major consequences of a shallow charge layer. The results of the calculations are

$$V_A - V_{FB} = \phi_s + V_1 \left(\frac{\phi_s}{\phi_{s1}} \right)^{\frac{1}{2}}, \quad x_d \leq x_1$$

$$V_A - V_{FB} = \phi_s + V_1 - V_2 \left\{ 1 - \sqrt{1 + \frac{N_{A1}}{N_{A2}} \left(\frac{\phi_s}{\phi_{s1}} - 1 \right)} \right\}, \quad x_d \geq x_1 \quad (1)$$

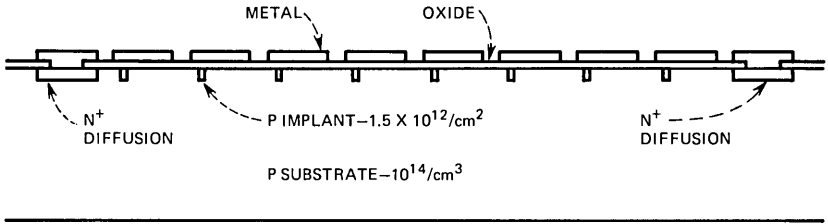


Fig. 2—Cross section of the two-phase charge-coupled device which was fabricated (the actual device made has 8 bits, only 4 are shown here).

where V_{FB} is the flat-band voltage and

$$V_1 = \frac{ex_1 dN_{A1}}{\epsilon_{ox}}$$

$$V_2 = \frac{ex_1 dN_{A2}}{\epsilon_{ox}}$$

$$\varphi_{s1} = \frac{ex_1^2 N_{A1}}{2\epsilon_s}$$

where ϵ_s and ϵ_{ox} are the permittivities of the Si substrate and the SiO₂ film respectively and d is the SiO₂ thickness.

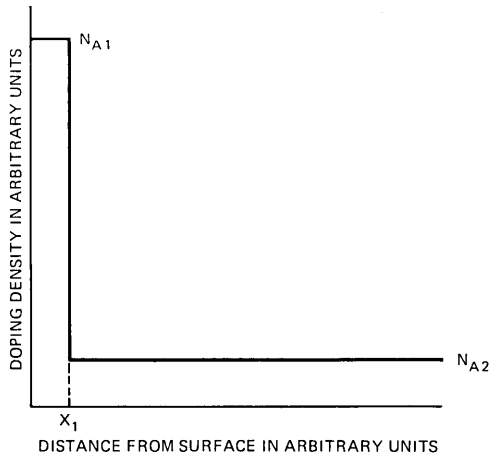


Fig. 3—Doping density vs distance from semiconductor surface.

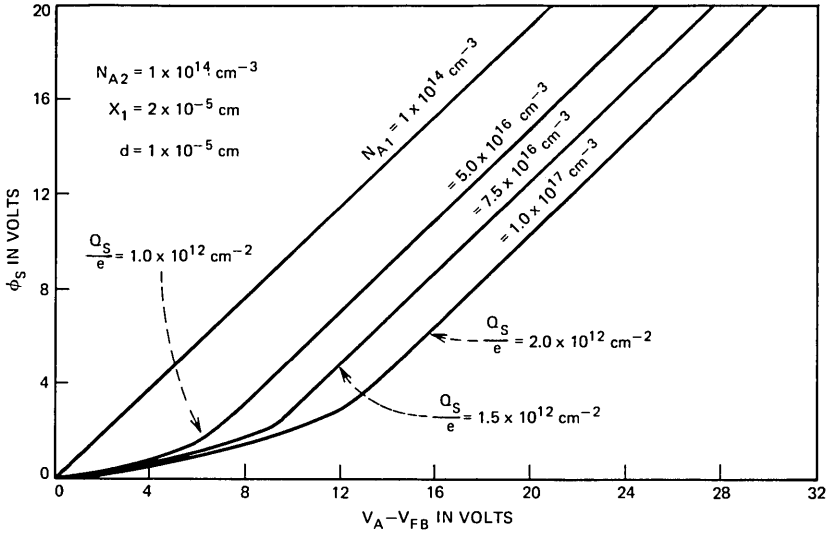


Fig. 4—Curves of surface potential vs applied voltage for various values of surface doping.

The parameters ϕ_{s1} and V_1 are the voltage drops across the Si depletion region and SiO_2 , respectively, when the depletion layer width x_d is the same as x_1 . The results of eq. (1) are plotted in Fig. 4 for several values of N_{A1} ; it was assumed that $N_{A2} = 1 \times 10^{14} \text{ cm}^{-3}$, $x_1 = 2 \times 10^{-5} \text{ cm}$ and $d = 1 \times 10^{-5} \text{ cm}$. Each curve has two regions of distinctly different behavior: the first extends from the origin to the knee, shows ϕ_s varying slowly with respect to V_A , and corresponds to the condition $x_d \leq x_1$; beyond the knee, the curve is nearly linear and corresponds to $x_d > x_1$. The point indicated on each curve corresponds to $x_d = x_1$, and the associated value of $Q_s = ex_1N_{A1}$ is the total charge per unit area in the heavily doped region.

The effect of increasing N_{A1} for a given x is to increase the height of the barrier between implanted and unimplanted regions. The barrier height also depends on applied voltage. It is zero for an applied voltage equal to the flat-band voltage, but it is essentially constant for voltages large enough to insure $x_d > x_1$. Operation with V_A dropping into the region where the barrier is lower than its maximum is permissible, even though an instantaneous decrease in V_A (and consequently V_B) during the transfer process, would permit some of the charge to flow backwards over the shrinking barrier. Actually, V_A has a finite time derivative so some charge is transferred forward before V_A reaches its minimum. This

initial transfer of charge is extremely fast. Using calculations made by Strain and Schryer⁵, an n -channel device with 25μ electrodes, half of the charge is transferred in 5 ns. For a 10μ electrode, only 1 ns is required. Therefore, if the pulse generators driving the clock lines have rise times of several ns, the barrier can shrink by at least a factor of two with no loss of charge. This is because the amount of charge being held back by the barrier decreases faster than the barrier height decreases.

Typical surface potential values during operation are illustrated in Fig. 5. Two curves are shown; one is for the unimplanted region ($1 \times 10^{14} \text{ cm}^{-3}$ p -type), and the other is for the implanted area ($7.5 \times 10^{16} \text{ cm}^{-3}$ p -type, $x_1 = 2 \times 10^{-5} \text{ cm}$). The circled points A, B, C and D determine the operating parameters of the CCD: the peak-to-peak variation in the clock voltage is determined by the voltage difference ΔV_A between points B and C, while the potential well depth $\Delta\phi_s$ is determined by the separation of points A and B or of C and D.

The lower value of applied voltage ($V_A - V_{FB} = 3$ volts) was chosen to make the barrier approximately half of its maximum value. The larger value of V_A ($V_A - V_{FB} = 11$ volts) was chosen to make the surface potential at the top of the barrier under the receiving electrode equal to the surface potential in the unimplanted region under the

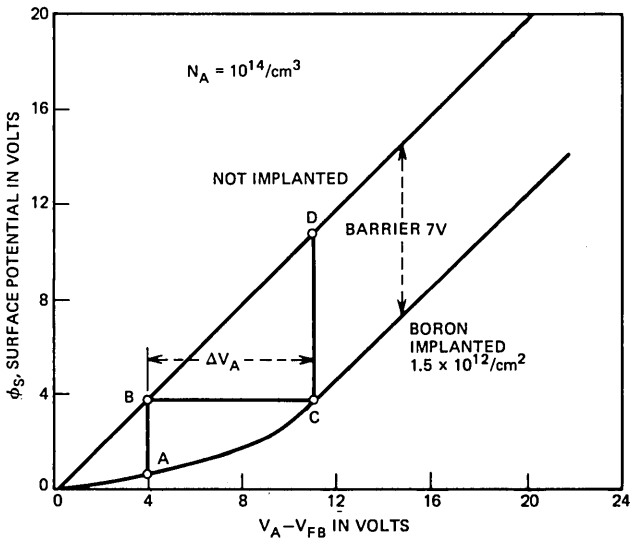


Fig. 5—Curves of surface potential vs applied voltage for both implanted and unimplanted regions (implant $7.5 \times 10^{16}/\text{cm}^3$, depth 0.2μ). Also indicated are the driving voltage used and corresponding surface potentials. These are the parameters of the device fabricated.

electrode giving up charge. Any value of $V_A - V_{FB}$ greater than 11 volts could also be used. Since these values of V_A are convenient, the values of N_{A1} , N_{A2} , and x_1 (10^{14} , 7.5×10^{16} , 2×10^{-5} cm) used for Fig. 5 were also used in the device that was fabricated.

In the design of the device, transfer across the spaces between the electrodes had to be considered also. Calculations made by Krambeck⁶ have shown that for a p substrate with a doping of $10^{14}/\text{cm}^3$ and a thermally grown SiO_2 insulator, the gaps should not interfere with transfer. The interelectrode spacings were made 5μ , and the width of the electrodes was made 35μ . The implanted regions were offset from the left edges of the gates by 5μ and were 5μ wide. These parameters give a device which is practical to fabricate and is capable of moving useful amounts of charge at high rates. The construction, testing and operation of this CCD will be described in the next section.

III. EXPERIMENTAL RESULTS

3.1 Processing

The fabrication of the device involves five steps requiring photolithography. Four of these use standard planar processing: diffusing the input and output diodes, diffusing the channel stop, etching the contact holes and etching the metallization. The fifth photolithographic step is needed to obtain the ion implanted pattern. For this, the photoresist is used as the mask and $1.5 \times 10^{12}/\text{cm}^2$ of boron is implanted. The boron is activated with a high temperature anneal.

A photograph of a completed device is shown in Fig. 6. Along with the shift register itself, several test structures have been fabricated as well. Two MOS capacitors appear in the lower left-hand corner of the photo, one in which the entire region under the gate has been subjected to the ion implantation, and the other which is completely devoid of any implant.

3.2 Pretest

The MOS capacitors were used to determine the operating characteristics of the shift register. The three major measurements that are required are: (i) surface potential ϕ_s as a function of applied voltage V_A ; (ii) generation time T_g for the buildup of an inversion layer in the Si and (iii) surface state density N_{ss} . The first measurement gives information regarding the height of the potential barrier, the second gives an indication of the low frequency limit of device operation, and the third relates to the charge transfer efficiency.⁴

The measurement techniques are based on the use of ramps for

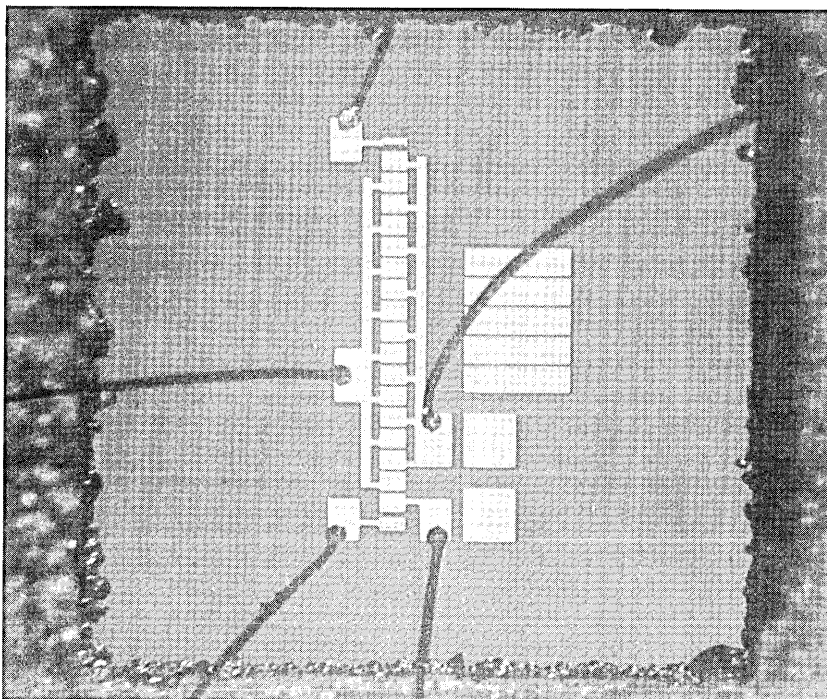


Fig. 6—Photograph of a completed device.

capacitance measurements. Surface potential as a function of applied voltage is obtained by integrating the deep depletion C - V curve as is indicated in the following equation⁷

$$\varphi_s(V_A) - \varphi_s(V_{\text{REF}}) = \int_{V_{\text{REF}}}^{V_A} \left(1 - \frac{C(V)}{C_{\text{ox}}} \right) dV \quad (2)$$

where V_{REF} is some convenient reference voltage and $C(V)$ and C_{ox} are the capacitance at some voltage V and the oxide capacitance respectively. The deep depletion C - V curve is obtained using the "fast ramp" technique⁸ and the integration is done by simply operating the electrometer used to measure the displacement current as a coulombmeter. Surface potential curves for the test structures are shown in the photographs in Fig. 7. The curve in Fig. 7a is for the unimplanted sample, and that in Fig. 7b for the implanted specimen. These curves are similar in appearance to the calculated curves in Fig. 5. At large positive values of V_A , the difference in surface potential $\Delta\varphi_s$ for a given V_A is essentially constant and equal to 6.3 volts. This value for the potential barrier

height corresponds to a total implanted charge density of approximately $1.3 \times 10^{12} \text{ cm}^{-2}$ which is 80 percent of the original dose.

In order to obtain a value for the generation time T_g , the square wave response of the surface potential was measured. The positive-going portion of the square wave causes the MOS capacitor to be driven suddenly into deep depletion, and ϕ_s assumes a large value, then while the applied voltage remains at its peak value, an inversion layer builds up at the silicon surface, causing ϕ_s to decay toward its equilibrium level. The negative-going portion of the square wave drives the sample into accumulation, and ϕ_s is a constant throughout this half cycle at a level which is approximately one volt below the inversion equilibrium value. Photographs of ϕ_s as a function of time are displayed in Fig. 8. The response curves show the decay of ϕ_s followed by an abrupt change to a constant value corresponding to the accumulation condition. The signals for both capacitors to $1/e$ of their initial value in $T_g = 0.1 \text{ s}$ which implies a lower frequency operating limit of a few hundred Hz.

The surface state density N_{ss} in the implanted test capacitor was determined by the method of comparison of high frequency (10^6 Hz) and low frequency (quasistatic) $C-V$ curves.⁹ The derived distribution is shown in Fig. 9. As shown in the Appendix, this surface state density is too low to cause visible loss in an 8-bit CCD.

The above tests showed that while the properties of the oxide-silicon interface were not ideal, they were good enough to permit operation of the CCD. In particular, the tests of ϕ_s , vs V_A showed that a 6-volt potential barrier does exist between implanted and unimplanted regions under the same electrode. Also the -3 -volt flat band voltage indicated that the positive charge in the oxide ($\sim 5 \times 10^{11}/\text{cm}^2$) was in the range

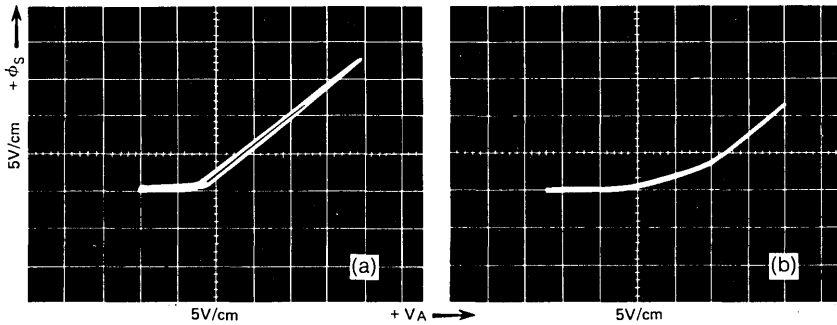


Fig. 7—Measured curves of surface potential vs applied voltage. (a) unimplanted, (b) implanted $1.5 \times 10^{12}/\text{cm}^2$ boron.

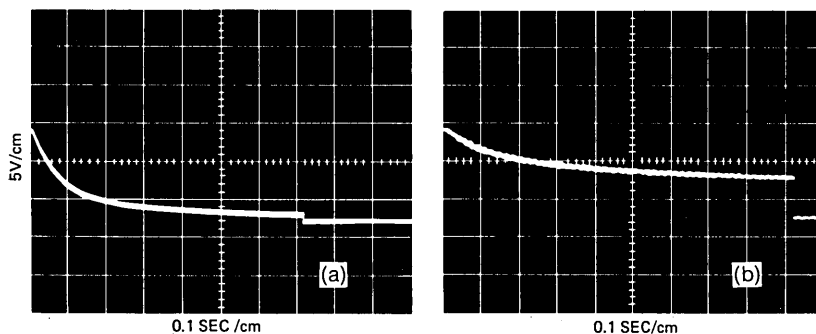


Fig. 8—Surface potential vs time after application of a step in voltage applied. (a) implanted, (b) unimplanted.

which insured transfer across the gap between electrodes.⁶ The value of flat-band voltage is quite convenient since, from Fig. 5, it permits the use of 0 volt as the lower value of V_A .

3.3 Shift Register Operation

To find if the device actually worked as predicted, the circuit of Fig. 10 was used. This circuit includes two diodes; one at each end of the shift register. These provide a convenient way of introducing minority carriers at one end, and of removing and sensing that charge after it has passed through the entire register. The operation of the device is as follows: the two diodes are normally held in reverse bias

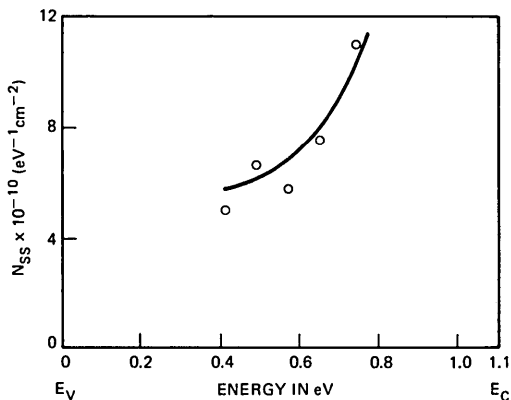


Fig. 9—Calculated values of surface state density vs position in the forbidden gap.

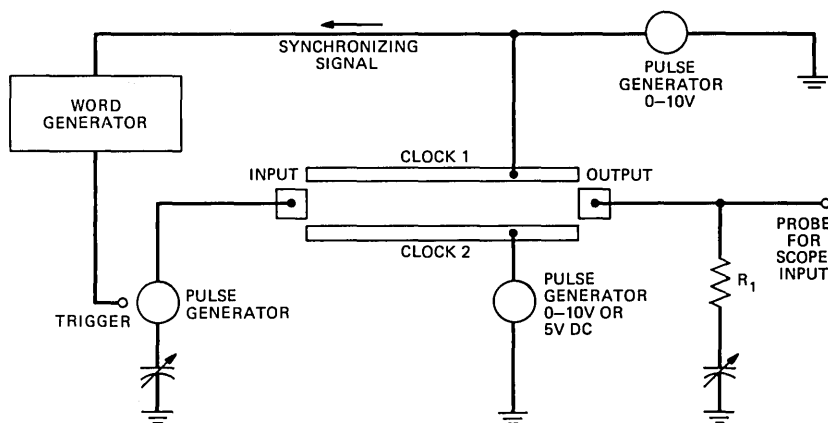


Fig. 10—Circuit used to test the device.

of about 15 volts to prevent injection of charge. The clock bias can be driven in either of two ways. These are shown in Figures 11a and 11b respectively. In one of these modes, both clock bias are driven with square waves, while in the second, only one square wave is used. This latter mode has the obvious advantage that no synchronization is necessary. The two modes give the same operation because charges move in response to differences in potential. Inspection of Figure 11a and 11b shows that the difference between voltages on the two clock lines as a function of time is essentially the same for both. As long as the input diode is held at a fixed voltage of 15 volts, no charge is injected, and ZERO's are sent through the shift register. The only output is that caused by capacitive coupling between the clock lines and the output diode pad.

To provide charge input, a negative pulse is applied to the input diode while a positive pulse is being applied to the first electrode in the shift register. Electrons will flow out of the *n*-diffusion into the region under the electrode at a rate dependent on the relative voltage between the two. The process is illustrated in Fig. 12. The surface potential under the first electrode ultimately becomes equal to that in the diode. When the negative pulse ends, any excess charge flows back into the diode. The resulting packet of charge is then transferred step by step to the other end of the shift register, and is finally transferred into the second diode when the last electrode is driven to the substrate potential. This charge changes the diode voltage by an amount dependent on the total capacitance between the diode and ground. This

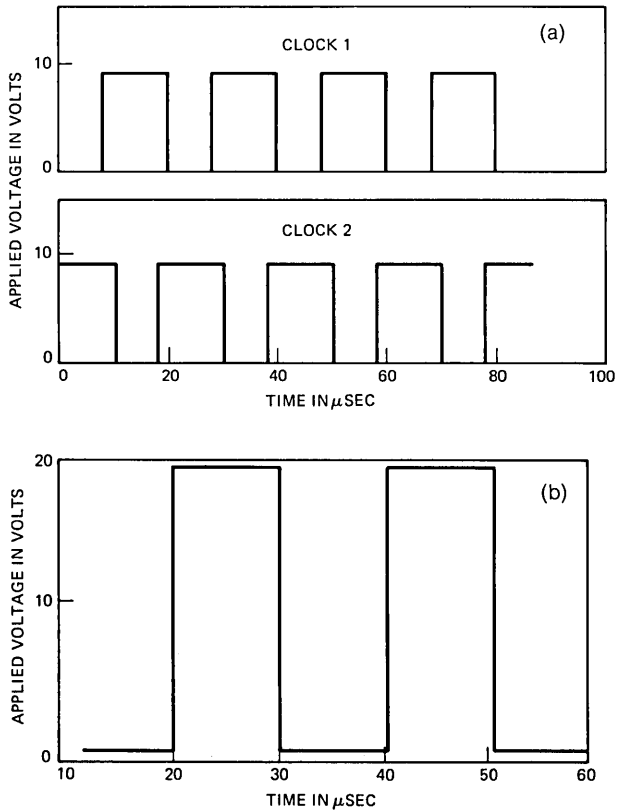


Fig. 11—Voltages used to drive the clock lines. (a) Two clock method and (b) One clock method (this was actually used to drive the device).

voltage decays to zero with a time constant dependent on the resistor R_1 in Fig. 10.

If a series of intermixed ONE's and ZERO's is fed into the shift register, the output should show a temporary voltage drop $\Delta V = Q/C$, where Q is the stored charge and C is the output capacitance. For this device, Q is determined by the activated implant density ($1.2 \times 10^{12}/\text{cm}^2$) and the area of the region which stores charge (10^{-5} cm^2). Therefore, Q is 1.9×10^{-12} coulombs and change in output voltage should be $1.9/C$ volts where C is in pF.

The results, with a clock frequency of 1.5 MHz, are shown in the top photograph in Fig. 13. The lower trace shows the pulses applied to the input diodes, and the upper trace is the output voltage. The output

swing caused by Q is 0.02 volt. It is clear that shift register action has been obtained since the output is delayed 8 cycles as it should be. This confirms the feasibility of guiding charge in a charge-coupled device by the implantation of properly chosen amounts of boron.

Measurements of transfer efficiency were carried out at several frequencies, two of which are shown in Fig. 13, and the results are shown in Fig. 14. The loss is too small to measure at 6.5 MHz, and it is estimated that a loss of 0.1 percent per transfer. Thus the loss increases from less than 0.1 percent per transfer below 6.5 MHz to 2 percent per transfer at 17 MHz.

So far the device has only been discussed as a digital shift register. However, intermediate amounts of charge may also be injected to demonstrate operation as an analog shift register. To accomplish this,

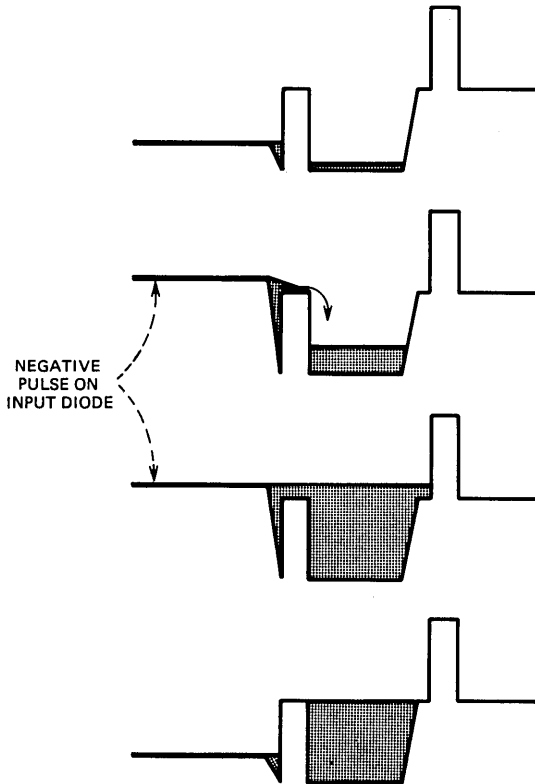


Fig. 12—Surface potential at input of shift register during injection of charge.

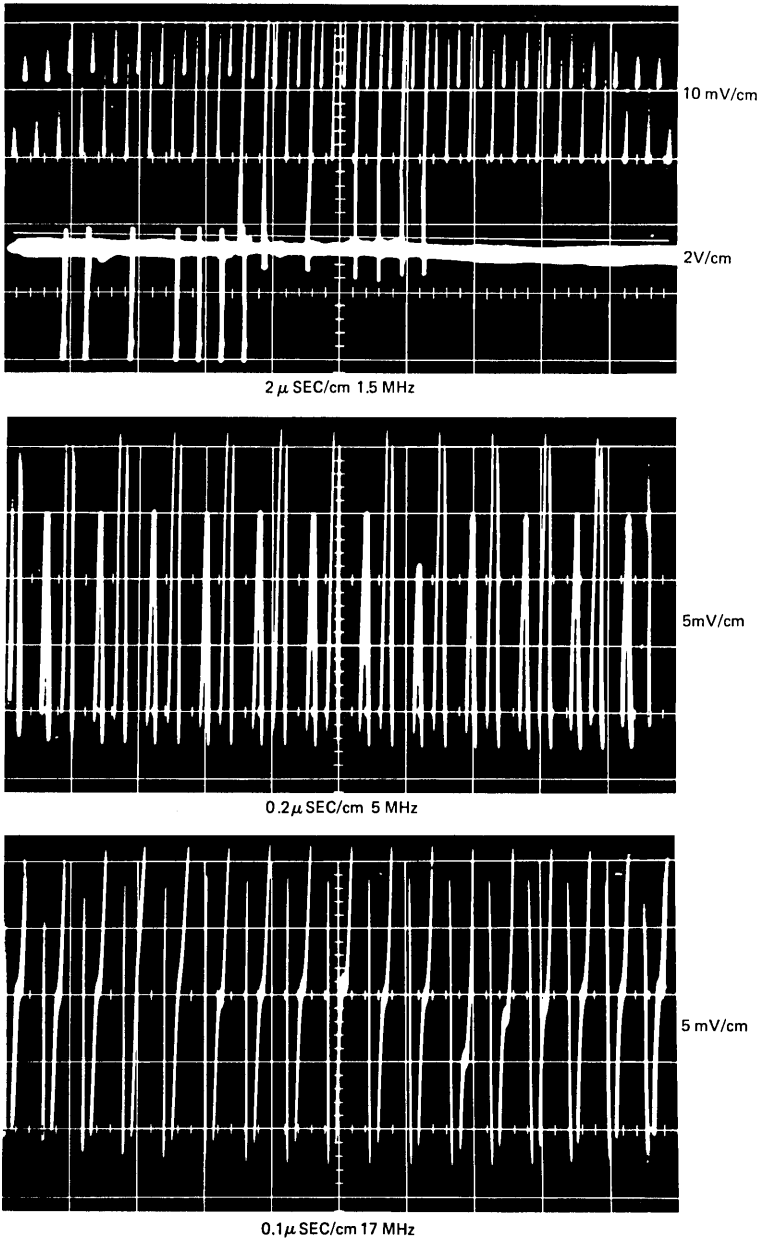


Fig. 13—Top photo: Input (lower trace) and output (upper trace) during operation of the shift register at 1.5 MHz. Middle and lower photos are output of a single ONE at 5 MHz and 17 MHz respectively.

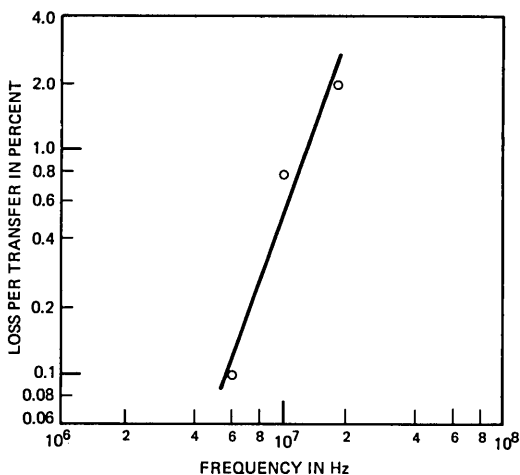


Fig. 14—Loss per transfer vs frequency (measured).

it is necessary to modulate the height of the input pulses so the rate of charge flow and the total amount of charge injected during each cycle have values between those previously used. This modulation was accomplished by placing a ramp generator in parallel with the input pulse generator. The resulting input waveform is shown on the bottom trace of Fig. 15. The output shows a ramp delayed eight periods of the clock frequency.

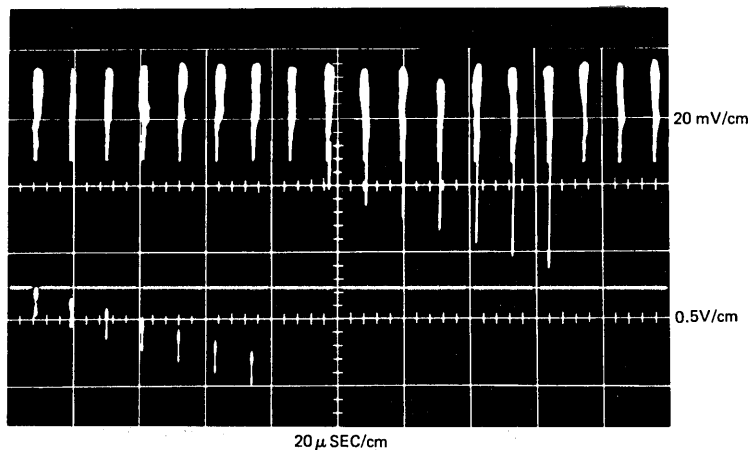


Fig. 15—Input and output (top trace) during operation as an analog device.

IV. CONCLUSION

The successful operation of an 8-bit n -channel two-phase CCD fabricated using ion implantation has been achieved. The device performed in the frequency range of 1 kHz to 6.6 MHz with no discernible loss ($\lesssim 0.1\%$ /transfer). At 17 MHz, it was 2%/transfer. The shift register was also operated as an analog delay line.

The surface potential barrier under a portion of each gate necessary for unidirectional charge transfer was realized by ion implanting 1.5×10^{12} boron ions/cm². This gave rise to a measured barrier height of approximately 6 volts.

The use of the ion-implanted barrier makes possible considerable simplification of the structure. There is only one layer of metallization and no crossovers or crossunders. Moreover, the simplification caused no sacrifice in performance, since the device showed more efficient high-frequency operation than any previously reported CCD.

V. ACKNOWLEDGMENT

The authors are grateful to H. F. Hamilton for his invaluable assistance in fabricating the devices used in this study.

APPENDIX

Previous calculations⁴ have shown that when the charge in a ZERO is a substantial fraction of the charge in a ONE, and square waves are used to drive the clock lines, there is virtually no loss caused by surface states. However, the implanted barrier represents a new feature with respect to surface state loss analysis.

The region that was implanted, and as a result contains a potential barrier, is exposed to a high concentration of free carriers only when transferring a ONE. There is essentially no free charge in this region at other times (refer to Fig. 1). As a result, when a ONE passes the barriers, some of the surface states in that region capture charge. This charge may be divided into three parts. Part one is charge which is re-emitted before the ONE has completed its transfer. This charge rejoins that in the ONE, and has no effect on loss. Part two is charge which is emitted too late to catch up, and is therefore dropped into the ZERO or ZERO's following the ONE. The remaining charge is still trapped when the next ONE arrives. The charge lost by this second ONE is the difference between the charge in the traps when it arrives, and the trapped charge when it leaves. This is the charge referred to as part two above.

Let us follow what happens to the charge which is trapped as a ONE crosses the barrier. The charge rejoining the ONE needs no further consideration. Charge which is emitted during the succeeding half cycle can travel either forward or backward, since the barrier region in question is at a potential maximum. As a result, one half of the charge emitted during this half cycle is lost. The amount of this charge can be determined as follows: Assume that all of the traps, the time constants of which are between $\tau/2$ and τ (where τ is the period of the clock frequency), will emit their captured charge during this half cycle. The number of such traps is $N_{ss}(E) \cdot kT \cdot \ln(\tau/0.5\tau)$, where $N_{ss}(E)$ is the density of traps, per eV and per cm^2 , the position of which in the forbidden gap makes their emission rate $1/\tau$. (It will be assumed that $N_{ss}(E)$ does not vary significantly for a few kT .) Then the charge lost per unit area of potential barrier during this half cycle is

$$\frac{1}{2} \cdot N_{ss}(E) \cdot kT \cdot \ln 2. \tag{3}$$

The formula to convert position in the gap to period of the clock line is

$$f = N_c c_n \exp [-(E_c - E)/kT]$$

or

$$f = 3 \times 10^{10} \exp [-(E_c - E)/kT], \tag{4}$$

where N_c is conduction band effective density of states, c_n is the capture probability, and $E_c - E$ is the depth of the state below the conduction band edge. Formula (4) can be used to determine the appropriate value of E in formula (3).

Suppose n ZERO's now follow the ONE. During the n cycles required for these to pass, no free carriers cross the barrier. However, the surface states continue to emit trapped charge. This goes on from $t = \tau$ to $t = (n + 1)\tau$. The total amount of charge emitted during this period is therefore

$$N_{ss}(E) \cdot kT \cdot \ln ((n + 1)\tau/\tau)$$

or

$$N_{ss}(E) \cdot kT \cdot \ln (n + 1).$$

Not all of this charge is lost however, because the charge can flow backward or forward. The charge which flows backward during the time between $(n + \frac{1}{2})\tau$ and $(n + 1)\tau$ will join the ONE following the n ZERO's and therefore does not contribute to loss. This backward flowing charge is

$$\frac{1}{2} N_{ss}(E) \cdot kT \cdot \ln((n + 1)\tau/(n + \frac{1}{2})\tau).$$

The net loss per unit area of potential barrier is

$$L_n = \frac{1}{2}N_{..}(E) \cdot kT \cdot \ln 2 + N_{..}(E) \cdot kT \cdot \ln(n + 1) \\ - \frac{1}{2}N_{..}(E) \cdot kT \cdot \ln(n + 1/n + 2)$$

or

$$L = \frac{1}{2}N_{..}(E)kT[\ln(2n + 1)(n + 1)]. \quad (5)$$

Equations (4) and (5) can now be used to calculate expected loss per transfer vs frequency. For the experimental measurements of loss per transfer, $n = 1$. Therefore

$$L = 1.3N_{..}(E)kT.$$

As an example at $E_c - E = 0.3 \text{ eV}$, from Fig. 10, $N_{..} = 10^{11}/\text{cm}^2 - \text{eV}$. Therefore $L = 0.33 \times 10^{10}/\text{cm}^2$. From (4), the frequency is $2.7 \times 10^5 \text{ Hz}$. Since the area of the potential barrier is $2.5 \times 10^{-6} \text{ cm}^2$, the charge lost per transfer is $0.13 \times 10^{-14} \text{ coulombs}$. The charge storage capacity is $2 \times 10^{-12} \text{ coulombs}$ which gives a loss per transfer figure of 0.07 percent. This compares to the "measured" loss of less than approximately 0.1 percent. This low loss is directly attributable to the two-phase structure, since only a small fraction of the surface states (those in the implanted region) can cause any loss.

REFERENCES

1. Boyle, W. S., and Smith, G. E., "Charge Coupled Semiconductor Devices," B.S.T.J. 49, No. 4 (April 1970), pp. 587-593.
2. Tompsett, M. F., Amelio, G. F., and Smith, G. E., "Charge Coupled 8-Bit Shift Register," Appl. Phys. Letters, 17 (1970), p. 111.
3. Engeler, W. E., Tiemann, J. J., and Baertsch, R. D., "A Memory System Based on Surface Charge Transport," 1971 IEEE Solid State Circuits Conference.
4. Strain, R. J., "Power and Surface State Loss Analysis of Charge Coupled Devices," 1970 International Electron Devices Meeting.
5. Strain, R. J., and Schryer, N., "A Nonlinear Diffusion Analyses of Charge-Coupled-Device Transfer," B.S.T.J. 50, No. 6 (July-August 1971), pp. 1721-1740.
6. Krambeck, R. H., "Zero Loss Transfer Across Gaps in a CCD," B.S.T.J., 50, No. 10 (December 1971), pp. 3169-3175.
7. Berglund, C. N. IEEE Trans. Elec. Devices, ED-13 (1966), p. 701 ff.
8. Kuhn, M., Int. Elec. Dev. Mtg., Oct. 29-31, 1969, paper 11.2.
9. Castange, R., C. R. Acad. Sci., Paris, 17 (1968), p. 866 ff.

Almost-Coherent Detection of Phase-Shift-Keyed Signals Using an Injection-Locked Oscillator

By M. EISENBERG

(Manuscript received August 18, 1971)

We analyze a proposed scheme of detection of phase-shift-keyed signals using an injection-locked oscillator the bandwidth of which is much less than the modulation rate. The output of the oscillator is a carrier with essentially all of its modulation removed. We analyze the effect of noise and signal modulation on the phase of this reference tone and compute its effect on the probability of detection error. If a suitable encoder and decoder are used for the transmitted signal, this technique can provide nearly ideal coherent demodulation.

I. INTRODUCTION

The two generally recognized methods of detection of phase-shift-keyed (PSK) signals are coherent detection and differential detection. Coherent detection has been shown to be optimum in the presence of Gaussian noise,¹ but, due to the difficulty of storing an absolute phase reference at the receiver, it is seldom used in practice.

In a recent paper,² B. Glance showed that an injection-locked oscillator, the locking bandwidth of which is much less than the modulation rate can, under certain conditions, be used to derive a phase reference from the input signal itself. This is actually a form of a quadrature reference system, where the phase of one quadrature remains essentially unkeyed, and is used to provide the reference tone.^{3,4} In this paper, we examine the effectiveness of this scheme for a two-phase PSK system where the modulation is a random digital signal and additive Gaussian noise is present. We derive an expression for the probability distribution of the reference phase, and from this calculate the average probability of a detection error. We find that if the modulation rate is much greater than the bandwidth of the oscillator, and if a suitable encoder and decoder are used, the method very nearly approaches the ideal performance of coherent detection.

II. LOCKING EQUATION ANALYSIS—ZERO ORDER SOLUTION

A portion of the received signal is used as the input to the injection-locked oscillator. This signal may be represented as

$$x(t) = \sqrt{2} A \cos [\omega t + \theta(t)] + n(t), \quad (1)$$

where A is the signal power, ω is the carrier frequency, and $\theta(t)$ is the phase modulation. The received signal is assumed to be contaminated by additive white Gaussian noise, $n(t)$, with double-sided spectral density $N_0/2$. If ω is sufficiently close to the natural oscillator frequency, ω_0 , locking will occur and the output of the oscillator will be

$$y(t) = \sqrt{2} B \cos [\omega t + \theta(t) - \phi(t)], \quad (2)$$

where B can be assumed to be constant.⁵ $\phi(t)$ is the phase difference between the input and output signals of the oscillator. In the case of interest, the total phase modulation of the oscillator output, $\eta(t) = \theta(t) - \phi(t)$, is small, and $y(t)$ is used as the phase reference in the coherent detection of the remaining portion of the received signal.

In the absence of noise, the phases of the input and output signals of the oscillator are related by the well-known locking equation⁶

$$\frac{d\phi(t)}{dt} + \Delta \sin \phi(t) = \omega - \omega_0 + \frac{d\theta(t)}{dt}, \quad (3)$$

where 2Δ is the locking bandwidth of the oscillator. This equation takes the same form as that for a first order phase-locked loop.⁷ The effect of the noise at the input has been analyzed by Viterbi.⁸ The effect is to add an additional term to eq. (3),

$$\frac{d\phi(t)}{dt} + \Delta \sin \phi(t) = \omega - \omega_0 + \frac{d\theta(t)}{dt} - \frac{\Delta}{A} n'(t), \quad (4)$$

where $n'(t)$ has the same statistics as $n(t)$. We rewrite eq. (4) in terms of $\eta(t)$.

$$\frac{d\eta(t)}{dt} + \Delta \sin [\eta(t) - \theta(t)] = \omega_0 - \omega + \frac{\Delta}{A} n'(t). \quad (5)$$

For the case of zero input phase modulation, i.e., $\theta(t) = 0$, eq. (5) becomes

$$\frac{d\eta(t)}{dt} + \Delta \sin \eta(t) = \omega_0 - \omega + \frac{\Delta}{A} n'(t). \quad (6)$$

Using Fokker-Plank techniques, Viterbi derived from this equation

$p(\eta)$, the steady-state probability density of η . For $\omega = \omega_0$ the solution is

$$p(\eta) = \frac{e^{\alpha^2 \cos \eta}}{2\pi I_0(\alpha^2)}, \quad (7)$$

where $I_0(\cdot)$ is the zeroth order modified Bessel function of the first kind, and $\alpha^2 = 4A^2/N_0\Delta$ is the signal-to-noise ratio in the bandwidth of the oscillator. For $\alpha \gg 1$, this distribution for η small is nearly Gaussian with mean zero and standard deviation $1/\alpha$. For $\omega \neq \omega_0$ the distribution becomes centered about the point $\beta = \sin^{-1}(\omega_0 - \omega)/\Delta$. In the case, $\alpha \gg 1$ and $|(\omega_0 - \omega)/\Delta| \ll 1$, $\beta \approx (\omega_0 - \omega)/\Delta$ and for η small, the distribution is very nearly equal to

$$p(\eta) \approx \frac{e^{\alpha^2 \cos(\eta - \beta)}}{2\pi I_0(\alpha^2)}. \quad (8)$$

We now consider the case $\theta(t) \neq 0$. In a binary PSK signal, $\theta(t)$ is a waveform of the form

$$\theta(t) = \sum_n a_n p(t - nT), \quad (9)$$

where $a_n = \pm 1$, and $p(t)$ is assumed to be a pulse which is nonzero only over the range $0 < t < T$. For the moment we assume zero noise. The resulting equation is

$$\frac{d\eta(t)}{dt} + \Delta \sin[\eta(t) - \theta(t)] = \omega_0 - \omega \quad (10)$$

which we rewrite as

$$\frac{d\eta(t)}{dt} + \Delta \sin \eta(t) \cos \theta(t) - \Delta \cos \eta(t) \sin \theta(t) = \omega_0 - \omega. \quad (11)$$

Our technique for the solution of this case will suggest a method of handling the stochastic problem when the noise term is reintroduced.

The exact solution to eq. (11) is difficult or impossible to obtain for general $\theta(t)$. Let us therefore take advantage of the fact that $T \ll 1/\Delta$ to derive a differential equation, the solution of which approximates that of eq. (11).

We assume that there exists an interval of length τ with the property that $T \ll \tau \ll 1/\Delta$, and we take the average of eq. (11) over τ . Letting $\langle \rangle$ denote this averaging operation, i.e.,

$$\langle \eta(t) \rangle = \frac{1}{\tau} \int_{t-(\tau/2)}^{t+(\tau/2)} \eta(u) du \quad (12)$$

there results

$$\frac{d\langle\eta(t)\rangle}{dt} + \Delta\langle\sin \eta(t) \cos \theta(t)\rangle - \Delta\langle\cos \eta(t) \sin \theta(t)\rangle = \omega_0 - \omega. \quad (13)$$

Choosing $\tau \ll 1/\Delta$ insures that $\eta(t)$ is very nearly constant over the interval of averaging, for, from eq. (10), we have $|d\eta(t)/dt| \leq \Delta + |\omega_0 - \omega| \leq 2\Delta$. Consequently, $\eta(t) \approx \langle\eta(t)\rangle$ and the quantities $\sin \eta(t)$ and $\cos \eta(t)$ may be taken outside the averaging operator.

$$\frac{d\langle\eta(t)\rangle}{dt} + \Delta \sin \langle\eta(t)\rangle \langle\cos \theta(t)\rangle - \Delta \cos \langle\eta(t)\rangle \langle\sin \theta(t)\rangle \approx \omega_0 - \omega. \quad (14)$$

On the other hand, the choice $\tau \gg T$ insures that there will be many data pulses over the interval of averaging. Thus the quantity $\langle\cos \theta(t)\rangle$ is very nearly constant and is equal to

$$C \equiv \langle\cos \theta(t)\rangle = \frac{1}{T} \int_0^T \cos p(t) dt. \quad (15)$$

(Table I lists the value of this quantity for three important pulse shapes.)

Since we usually have no control over the transmitted message, the quantity $\langle\sin \theta(t)\rangle$ will not, in general, be time-independent. However, a scheme has been suggested by C. L. Ruthroff and W. F. Bodtmann⁹ in which, through the use of a simple encoder and decoder, this quantity can be made very nearly equal to zero.

TABLE I—VALUES OF C AND σ FOR THREE PULSE SHAPES

$p(t)$	$C = \frac{1}{T} \int_0^T \cos p(t) dt$	$\sigma = \frac{1}{T} \int_0^T \sin p(t) dt$
raised cosine $\frac{\theta_0}{2} \left[1 - \cos \frac{2\pi t}{T} \right]$	$J_0\left(\frac{\theta_0}{2}\right) \cos \frac{\theta_0}{2}$	$J_0\left(\frac{\theta_0}{2}\right) \sin \frac{\theta_0}{2}$
positive sine $\theta_0 \sin \frac{\pi}{T} t$	$J_0(\theta_0)$	$\frac{4}{\pi} \sum_{k=0}^{\infty} \frac{J_{2k+1}(\theta_0)}{2k+1}$
rectangular $\theta_0, 0 \leq t \leq T$	$\cos \theta_0$	$\sin \theta_0$

The encoder is a device which stores a block of N bits of the message, and counts the net difference of $+1$'s and -1 's contained in the block. It also keeps a separate count of the net difference in $+1$'s and -1 's which have been sent over the entire past history of the message. The entire block is transmitted either with normal polarity, or with reversed polarity, in such a way as to cause the accumulated count to come as close as possible to zero. Preceding each block is a single "code" bit which specifies the polarity of that block. (The code bits are included in the counts.) The decoder decodes the message in an obvious manner.

We assume that $\theta(t)$ is the output signal of such an encoder. Assuming $\tau \gg NT$, we have

$$\langle \sin \theta(t) \rangle \approx 0, \quad (16)$$

so that our approximating differential equation becomes

$$\frac{d\eta_0(t)}{dt} + \Delta C \sin \eta_0(t) = \omega_0 - \omega. \quad (17)$$

We call the solution to this equation the "zero-order approximation to $\eta(t)$ ".

We note that this equation has the identical form as eq. (10) for the case of zero modulation. Thus the "zero-order effect" of the modulation is to reduce the effective value of the locking bandwidth by a factor C .

In Figs. 1 and 2, we demonstrate the above results. In both figures,

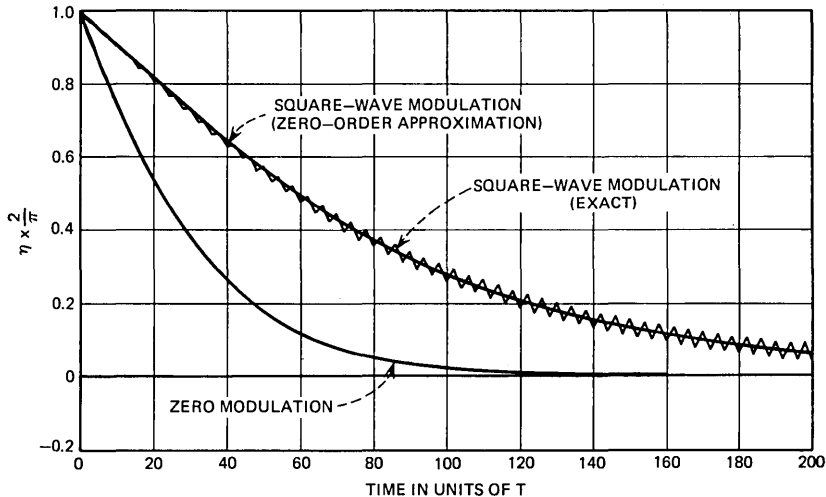


Fig. 1— $\eta(t)$ with square-wave modulation.

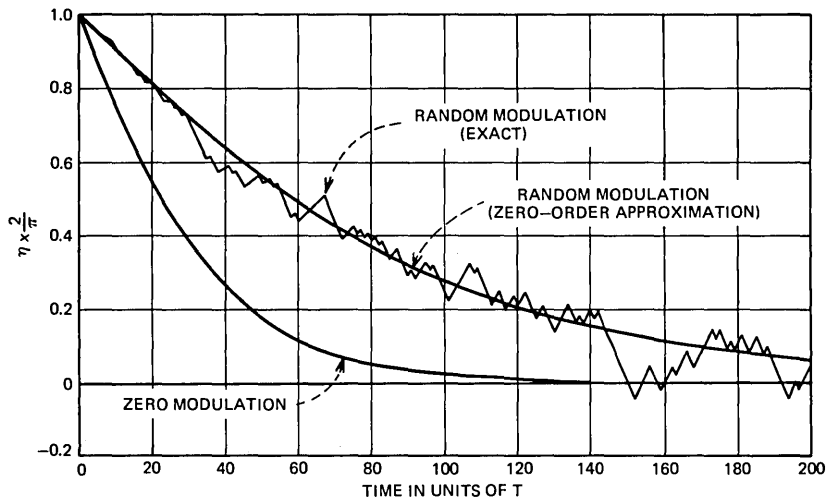


Fig. 2— $\eta(t)$ with random modulation.

$\Delta = \pi/80T$. The lower curve of Fig. 1 shows the behavior of $\eta(t)$ versus time for the case of zero modulation, i.e., $\theta(t) = 0$. The saw-tooth wave shows the behavior of $\eta(t)$, that is, the exact solution to eq. (10), for the case where $\theta(t)$ is a square wave with amplitude $3\pi/8$ and period $4T$. The smooth curve drawn over the saw-tooth wave shows the zero-order approximation obtained from eq. (17). As can be seen, the main effect of the modulation has been an increase in the decay time of the resulting curve. This is a result of the decrease in the effective locking bandwidth caused by the modulation, as predicted above.

In addition to this, the true curve has a "wiggle" which seems to increase in size as the curve approaches zero, and which reaches a maximum magnitude of about $0.02 \times \pi/2$.

Figure 2 shows the true curve and the zero-order approximation where the phase modulation is a random binary signal with $\text{Prob}(+1) = \text{Prob}(-1) = \frac{1}{2}$ which has then been passed through an encoder with $N = 5$; and a rectangular pulse shape of amplitude $3\pi/8$ is used. We note the curve follows the same overall path as in Fig. 1. The "wiggles" are now random in character; we note that their amplitude again appears to grow as the curve approaches zero.

III. FIRST ORDER SOLUTION

To better understand the behavior of these "wiggles" we now derive a correction term, $\eta_1(t)$, which, when added to $\eta_0(t)$, improves the accuracy

of our solution. We call $\eta_0(t) + \eta_1(t)$ the "first-order approximation to $\eta(t)$."

We define $\eta'(t) = \eta(t) - \eta_0(t)$. Subtracting eq. (17) from eq. (10), there results

$$\begin{aligned} \frac{d\eta'(t)}{dt} + \Delta \sin \eta'(t) \cos [\eta_0(t) - \theta(t)] \\ = \Delta C \sin \eta_0(t) - \Delta \cos \eta'(t) \sin [\eta_0(t) - \theta(t)]. \end{aligned} \quad (18)$$

If the assumptions we have made above are valid, the error in the zero-order approximation will be small, $|\eta'(t)| \ll 1$, and we may linearize eq. (18).

$$\begin{aligned} \frac{d\eta'(t)}{dt} + \Delta \eta'(t) \cos [\eta_0(t) - \theta(t)] \\ = \Delta C \sin \eta_0(t) - \Delta \sin [\eta_0(t) - \theta(t)]. \end{aligned} \quad (19)$$

Since $|\eta'(t)| \ll 1$, the second term on the left-hand side may be neglected, and we have approximately

$$\eta'(t) \approx \int_0^t \{ \Delta C \sin \eta_0(t) - \Delta \sin [\eta_0(t) - \theta(t)] \} dt, \quad (20)$$

where we assume the initial conditions are accounted for in $\eta_0(t)$. $\eta_0(t)$ varies slowly compared to $\theta(t)$, so we treat it as a constant in the integral. Our correction term is

$$\begin{aligned} \eta_1(t) = \Delta \sin \eta_0(t) \\ \cdot \int_0^t [C - \cos \theta(t)] dt + \Delta \cos \eta_0(t) \int_0^t \sin \theta(t) dt. \end{aligned} \quad (21)$$

Equation (21) explains the behavior of the "wiggles" in Figs. 1 and 2. If rectangular pulses are used, $\cos \theta(t) = C$, and

$$\eta_1(t) = \Delta \cos \eta_0(t) \int_0^t \sin \theta(t) dt. \quad (22)$$

For a square wave of period $4T$ and amplitude $3\pi/8$, $\eta_1(t)$ is a saw-tooth wave with amplitude $\Delta T \cos \eta_0(t) \sin 3\pi/8 = 0.924 \Delta T \cos \eta_0(t)$. The amplitude of the saw-tooth wave is seen to vary as the cosine of $\eta_0(t)$, and reach a maximum amplitude of $0.924 \Delta T$, which is $0.023 \times \pi/2$ for Fig. 1. This behavior agrees with our earlier observations. The accuracy of our approximation in this case is quite good. A plot of $\eta_0(t) + \eta_1(t)$ superimposed on the figure can not be distinguished by eye from the true curve $\eta(t)$.

We remark that the first term of eq. (21), which is zero for a rectangular pulse shape, may in general be ignored relative to the second term. The first term equals zero at the beginning and end of each pulse, and never achieves magnitude greater than $\Delta T \sin \eta_0(t)$. The second term, however, increases monotonically during a positive pulse and decreases monotonically during a negative pulse. The magnitude of this term can reach a maximum of $\frac{3}{2} N \Delta T \sigma \cos \eta_0(t)$, where

$$\sigma = \frac{1}{T} \int_0^T \sin p(t) dt. \quad (23)$$

(The factor $\frac{3}{2}$ arises because of the possibility of having a message block consisting of N +1's followed by a block consisting of $N/2$ +1's and $N/2$ -1's). For the case of interest, $\eta_0(t)$ will be near zero and N will be greater than about 10 or 20. This means that the second term of eq. (21) will predominate, and thus eq. (22) may be used for arbitrary pulse shapes. (Table I lists the value of σ for three important pulse shapes.)

IV. THE EFFECT OF NOISE

We now consider a system where both noise and modulation are present. We saw in eq. (17) that the zero-order effect of the modulation was simply to reduce the locking bandwidth by a factor C .

Accordingly, we take as our zero-order approximation the solution to the stochastic differential equation

$$\frac{d\eta_0(t)}{dt} + \Delta C \sin \eta_0(t) = \omega_0 - \omega + \frac{\Delta}{A} n'(t). \quad (24)$$

In this case, we are interested in the steady-state probability density for η_0 , $p(\eta_0)$. By comparison with eqs. (6) and (8), the solution can be written down immediately.

$$p(\eta_0) \approx \frac{e^{\alpha_e^2 \cos(\eta_0 - \beta_e)}}{2\pi I_0(\alpha_e^2)}, \quad (25)$$

where α_e^2 , the effective signal-to-noise ratio in the presence of modulation, is higher than the zero-modulation signal-to-noise ratio by a factor of $1/C$.

$$\alpha_e^2 = \frac{\alpha^2}{C} = \frac{4A^2}{N_0 \Delta C}. \quad (26)$$

The average phase shift due to frequency offset is increased.

$$\beta_e = \sin^{-1} \frac{\omega_0 - \omega}{\Delta C} \approx \frac{\omega_0 - \omega}{\Delta C} = \frac{\beta}{C}. \quad (27)$$

To get a correction term in our solution, we proceed as before, by taking the difference between eqs. (24) and (5). We notice that the noise term cancels, and we again obtain eq. (18) and the approximate solution, eq. (22). For the case of large signal-to-noise ratio, α_e , η_0 will, with high probability, be in the vicinity of zero (assuming $\beta_e \approx 0$). Since η_0 affects η_1 only as the cosine, η_1 is essentially independent of η_0 in this case and is equal to

$$\eta_1(t) = \Delta \int_0^t \sin \theta(t) dt. \quad (28)$$

$\eta_1(t)$ depends on the particular digital signal being transmitted. However, the use of the encoder described earlier insures that

$$|\eta_1(t)| < \frac{3}{2} N\Delta T\sigma. \quad (29)$$

Thus if each of the components is small, the output phase is seen to consist of the sum of three essentially independent parts:

- (i) A constant $(\omega_0 - \omega)/\Delta C$ resulting from the difference between the carrier frequency and the natural oscillator frequency.
- (ii) A time varying part which depends upon the digital modulation, and which has a maximum magnitude of $\frac{3}{2} N\Delta T\sigma$.
- (iii) A random part, due to noise, the distribution of which has a standard deviation of $1/\alpha_e = \sqrt{N_0\Delta C}/2A$.

Thus we can write the probability distribution of the reference phase approximately as

$$p(\eta) = \frac{e^{\alpha_e^2 \cos(\eta - \psi)}}{2\pi I_0(\alpha_e^2)}, \quad (30)$$

where ψ , the nonrandom portion of the phase, has a magnitude less than or equal to $\frac{3}{2} N\Delta T\sigma + |(\omega_0 - \omega)/\Delta C|$.

V. CALCULATION OF ERROR PROBABILITY

If η has a known value, the probability of a decoding error, assuming equal likelihood detection, is

$$P_e = \frac{1}{2} \operatorname{erfc}(\rho \cos \eta), \quad (31)$$

where ρ^2 is the signal-to-noise ratio in the bandwidth of the signal.¹⁰ If the receiving filter has a bandwidth $2W$, then

$$\rho^2 = \frac{A^2}{N_0 2W}, \quad (32)$$

where A is the rms signal amplitude, and $N_0/2$ is the double-sided spectral density of the noise. The average error probability \bar{P}_e , is obtained by averaging the quantity in eq. (31) over the possible values of η .

$$\bar{P}_e = \int_{-\pi}^{\pi} \frac{1}{2} \operatorname{erfc}(\rho \cos \eta) \frac{e^{\alpha_e^2 \cos(\eta-\psi)}}{2\pi I_0(\alpha_e^2)} d\eta. \quad (33)$$

This integral was performed using an expansion technique similar to that described in Ref. 10.

VI. DEMONSTRATION OF RESULTS

In order to reduce the noise as much as possible, the bandwidth of the receiving filter in a PSK system is usually set at the value which allows the signal to pass essentially undistorted. This bandwidth $2W$ is roughly given by

$$2W \approx \frac{1.6}{T}. \quad (34)$$

Thus the signal-to-noise ratio in the bandwidth of the loop, $\alpha_e^2 = (4A^2/N_0 \Delta C)$, is related to the signal-to-noise ratio of the received signal, $\rho^2 = (A^2/N_0 2W)$, as

$$\alpha_e^2 = \frac{6.4}{\Delta TC} \rho. \quad (35)$$

The effect of noise on the output phase is thus reduced by a factor of $\Delta TC/6.4$ from its input value. We have already seen that the effect of the phase modulation on the output phase was proportional to $\Delta T\sigma$. Thus the size of the quantity ΔT is important in determining the performance of the system: it should be kept as small as possible. The extent to which this can be done, however, is limited by the need to keep the difference between the carrier frequency and the natural oscillator frequency small relative to ΔC . This frequency difference may be reduced by the use of a negative feedback loop.¹¹ A reasonable value of ΔT presently obtainable in the laboratory for which these conditions can be satisfied is $\Delta T \approx 10^{-3}$.

Using the value $\Delta T \approx 10^{-3}$, and the relationship of eq. (35), the methods of the preceding section were employed to compute the average error probability. Under the assumption that the nonrandom portion of the output phase is 10 degrees, the results of this computation for a raised cosine pulse shape are plotted in Fig. 3. Also plotted in the figure are the curves representing true coherent detection, $\bar{P}_e = \frac{1}{2} \operatorname{erfc}(\rho)$, and

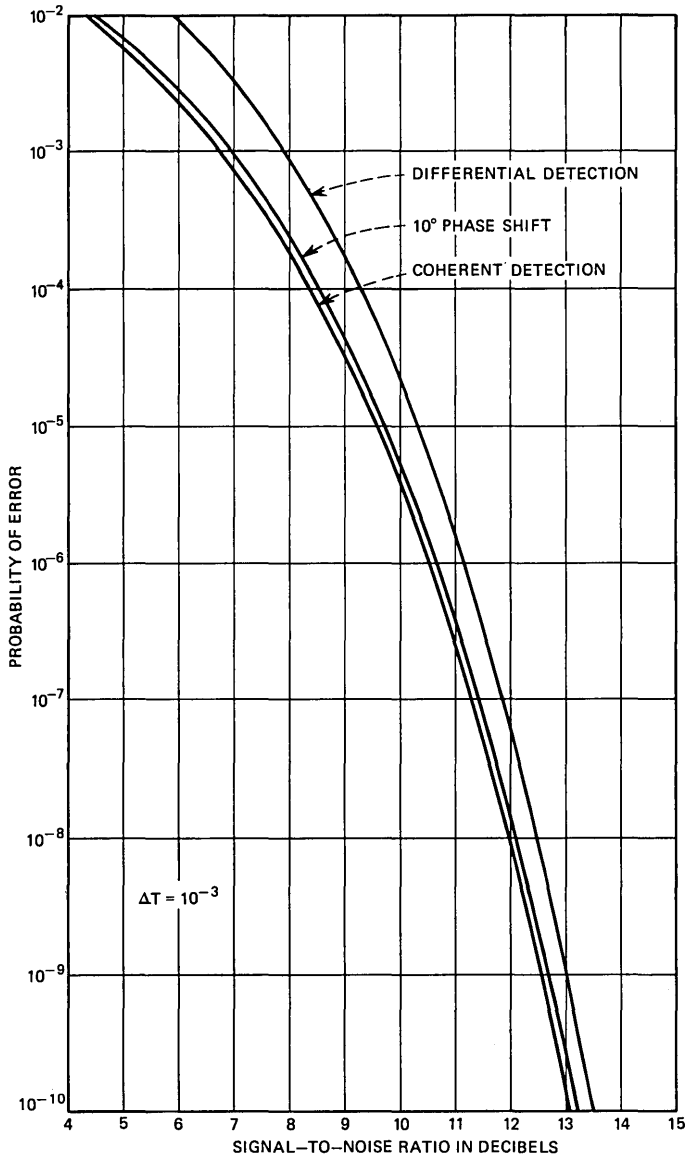


Fig. 3—Probability of error as a function of signal-to-noise ratio.

differential detection $\bar{P}_e = \frac{1}{2}e^{-\rho^2}$.¹² As can be seen, the curve comes quite close to the ideal of coherent detection. The curve corresponding to a nonrandom phase shift of 5 degrees was also computed but is not plotted in Fig. 3 because it comes so close to the coherent curve that the two cannot be distinguished on the scale of the drawing.

By taking the nonrandom part of the output phase to be zero, it is possible to determine the increase in error probability (over the coherent case) which is due to noise. We found the error probability to be virtually identical to coherent detection in this case. This indicates that the effect of the noise on the output phase shift is negligible.

For reasonably small values of $\Delta T (\leq 10^{-3})$, and for the range of signal-to-noise ratios usually of interest (> 5 db) the output phase shift resulting from noise may be safely ignored: the increase in error probability of the proposed system over coherent detection is almost entirely due to the effects of modulation and the offset in the carrier frequency. Consequently, under these conditions an approximate expression for the error probability, which is very nearly correct is $\bar{P}_e = \frac{1}{2} \operatorname{erfc}(\rho \cos \psi)$, where ψ is the phase shift resulting from the modulation and carrier offset. In Fig. 3, for example, the plotted curve is almost identical to $\frac{1}{2} \operatorname{erfc}(\rho \cos 10 \text{ degrees})$.

For $\Delta T = 10^{-3}$, a total of 289 consecutive raised-cosine pulses of maximum amplitude $\pi/2$ with the same polarity are needed to shift the output phase by 10 degrees. For $\Delta T = 10^{-2}$, this number is reduced to 29. For positive sine pulses of maximum amplitude $\pi/2$, the corresponding numbers are 241 and 24 respectively. We remark that increasing ΔT from 10^{-3} to 10^{-2} also increases the effect of the noise on the output phase, but that this effect remains negligible.

In order to demonstrate the effect of the noise on the output phase, we must consider an extremely high-noise example. Figure 4 plots the average error probability versus the signal-to-noise ratio over a range of from -7 to $+7$ db, for the case $\Delta T = 10^{-1}$. Nonrandom phase shifts of zero and 10 degrees were assumed respectively. As can be seen, the zero-phase shift curve almost coincides with the coherent curve for $\rho \geq 4$ db, even for this large value of ΔT .

VII. CONCLUSIONS

We have analyzed the proposed system of PSK detection for the case of random modulation and additive Gaussian noise. If the modulation rate is much greater than the bandwidth of the oscillator ($\Delta T \ll 1$), and if a suitable encoder and decoder are used, we have shown that the

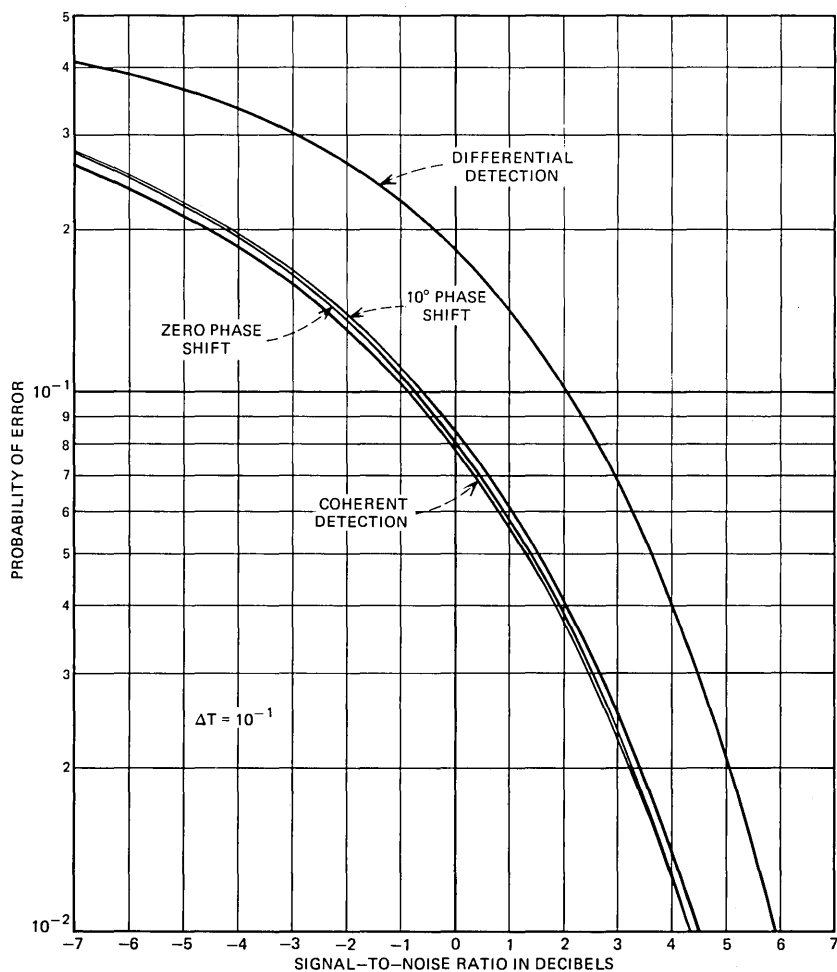


Fig. 4—Probability of error as a function of signal-to-noise ratio.

system will perform as almost a perfect coherent detector. For a binary system, this means a power savings of about $\frac{1}{2}$ db over the presently employed method of differential detection. This is not very great. However, with only a slightly more complex encoder and decoder, this same technique may be utilized for higher level systems. (The analysis requires only minor modifications.) For a 4-level system, for example, the power savings over differential detection is about 3 db which is significant.

Construction of an experimental encoder and decoder is presently being carried out in the laboratory, and tests of the system are planned to confirm the theoretical results.

VIII. ACKNOWLEDGMENTS

I would like to thank V. K. Prabhu for considerable help and encouragement during the course of this work. The problem was suggested by C. L. Ruthroff, who also suggested the general method of approach taken in this paper.

REFERENCES

1. Hancock, J. C., and Lucky, R. W., "Performance of Combined Amplitude- and Phase-Modulated Communication Systems," IRE Trans. Commun. Syst., CS-8 (December 1960), pp. 232-237.
2. Galance, B., "Digital Phase Demodulator," B.S.T.J., 50, No. 3 (March 1971), pp. 933-949.
3. Bussgang and Leiter, "Phase Shift Keying with a Transmitted Reference," IEEE Trans. COMTECH (February 1966), pp. 14-22.
4. Lindsey, "Phase-Shift-Keyed Signal Detection with Noisy Reference Signals," IEEE Trans. AES (July 1966), pp. 393-401.
5. Osborne, T. L., "Amplitude Behavior of Injection-Locked Oscillators," unpublished work.
6. Adler, R., "A Study of Locking Phenomena in Oscillators," Proc. IRE, 34, No. 6 (June 1946), pp. 351-357.
7. Ruthroff, C. L., and Stover, H. L., "The Similarity of the Phase-Locked Loop and the Injection-Locked Oscillator," unpublished work.
8. Viterbi, A. J., *Principles of Coherent Communication*, New York: McGraw-Hill Book Company, 1966.
9. Ruthroff, C. L., private communication.
10. Prabhu, V. K., "Error Rate Considerations for Coherent Phase-Shift Keyed Systems with Co-Channel Interference," B.S.T.J., 48, No. 3 (March 1969), pp. 743-767.
11. Ruthroff, C. L., "Injection-Locked Oscillator FM Receiver Analysis," B.S.T.J., 47, No. 8 (October 1968), pp. 1653-1661.
12. Cahn, C. R., "Performance of Digital Phase-Modulation Communication Systems," IRE Trans. on Commun. Syst., CS-7, No. 1 (May 1959), pp. 3-6.

Analysis of a Dual Mode Digital Synchronization System Employing Digital Rate-Locked Loops

By R. W. CHANG

(Manuscript received April 4, 1972)

We examine a data-rate synchronization system capable of operating in two modes: (i) in master-to-slave mode when the data stations are connected by digital transmission facilities, and (ii) in slave-to-slave mode when the data stations are connected by analog transmission facilities. The first part of this paper determines the steady-state behavior and the transient response in the master-to-slave mode. The results show that the system is well behaved in the transient stage, and that the steady-state behavior is satisfactory. From the transient analysis, the buffer size requirements of the system and the counter size requirements of the rate-locked loops are determined. Formulas are developed from which the start-up time of the system can be estimated.

The second part of this paper examines the behavior of the system in the slave-to-slave mode. It is shown that the data stations can settle to the same steady-state signaling rate, and this signaling rate is determined. The dependence of this signaling rate on other system parameters is examined. It is shown that the system can be easily designed such that the steady-state signaling rate will lie within desired limits. (This is so regardless of the starting sequence, the initial system conditions, and time delays in the communication channels.)

I. INTRODUCTION

When data stations are connected by wholly digital transmission facilities, it is most efficient to slave the clocks at the data stations to a master clock. To perform this operation, hereafter referred to as master-to-slave operation, an interface unit at the data station extracts timing pulses from the incoming data stream. These timing pulses are passed through a phase-locked loop to eliminate noise and jitter. The output of the phase-locked loop controls the signaling rate of the data station.

Unfortunately, a technical problem arises when data stations are synchronized in the above manner. Before digital systems evolve into a well-connected network, data stations are also often connected by wholly analog transmission facilities. When two data stations equipped to operate in the master-to-slave mode are connected by analog facilities, each station will regard the clock at the other station as the master clock, and the two stations will attempt to mutually synchronize each other. This mode of synchronization can be called "slave-to-slave." Conventional phase-locked loops¹ which perform well in the master-to-slave mode may not perform well in the slave-to-slave situation, being unusually sensitive to path-length delays and other system parameters. This technical problem can be solved by avoiding the slave-to-slave situation in the following manner:

- (i) Informing the data stations when analog transmission facilities are used. This will permit the stations to break up the slaving paths in the data sets and use their own clocks as the timing source.
- (ii) Providing a looped connection within the analog system containing appropriate buffers, and a clock of sufficient accuracy to serve as the master for the data stations.

Unfortunately, these schemes reduce the economic attractiveness of the system. Consequently, there is a need for a synchronization scheme capable of operating in both the master-to-slave and the slave-to-slave modes.

We analyze a synchronization system which employs digital rate-locked loops to determine if it can operate successfully in both modes. The phase detector in the rate-locked loop is a multistage counter that counts the difference between the number of zero crossings of the input signals. Because of this nonlinear counting process, the operation of the synchronization system is determined by nonlinear differential-integral equations. Such equations do not appear in earlier synchronization studies²⁻⁴ which considered different phase detectors. As will be shown, a digital rate-locked loop locks to neither the phase nor the frequency of the timing signal, but to the zero-crossing intervals. This difference complicates the analysis. We have examined the problem without making linear approximations.⁵⁻⁸ In a previous paper,⁹ we analyzed, in a rigorous fashion, the steady-state behavior of the system in the master-to-slave mode and proved that, in the absence of filtering in the rate-locked loop, the slave oscillator will lock to the master oscillator exactly. In this paper, it is proved rigorously that if the

filter in the loop satisfies a simple condition, the system will reach equilibrium (that is, data stations cannot add or delete bits from a customer's data stream). Based on this, it is demonstrated that in the presence of RC filtering in the loop, the slave oscillator will lock to the master oscillator exactly. Following these analyses, this paper determines the transient response of the system in the master-to-slave mode, and examines the behavior of the system in the slave-to-slave mode. Sections II and III of this paper examine the master-to-slave mode. Transient response, buffer-size requirements of the system, and counter-size requirements of the rate-locked loop are determined. Section IV considers the slave-to-slave mode. Steady-state signaling rate of the system is determined, with its dependence on the other system parameters examined. A simple method of designing the system to ensure satisfactory steady-state signaling rate is presented. Section V summarizes the results of this paper and may be read next.

II. MATHEMATICAL MODEL

In this and the following two sections, we examine the master-to-slave mode. Consider two communication stations as depicted in Fig. 1. Station 1 (with slave clock) represents a data station. Station 2 (with master clock) represents a station in the digital transmission facility. The master clock at Station 2 emits a timing signal which controls the transmission of data from Station 2 to Station 1 (for example, Station 2 transmits a digit to Station 1 at every second zero crossing of this timing signal). Station 2 transmits to Station 1 at some standard rate, say, f_2 digits per second.

Station 1 receives data from Station 2, and derives from the received data a timing signal $s_2(t) = \sin(\omega_2 t + \theta_2)$, where $\omega_2 = 2\pi f_2$ and θ_2 is

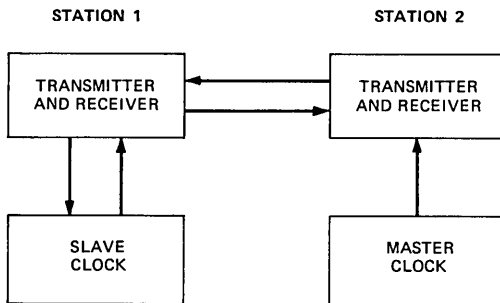


Fig. 1—Master-to-slave operation, block diagram.

an arbitrary phase angle. The signal $s_2(t)$ and the output $s_1(t)$ of a local oscillator are compared in a digital phase detector (Fig. 2). The digital-phase detector is a counter which counts the zero crossings of $s_2(t)$ and $s_1(t)$, and produces an output proportional to the difference between these two counts. Mathematically, this operation can be specified as follows. Let it be assumed that the digital phase detector is activated at $t = 0$. Let $N_1(t)$ and $N_2(t)$ be, respectively, the number of zero crossings (both upward and downward zero crossings) of $s_1(t)$ and $s_2(t)$ in the time interval 0 to t ; then the output of the digital phase detector is

$$u_1(t) = e_1[N_2(t) - N_1(t)] \quad (1)$$

where e_1 is a positive constant (volts/count) and may be called the gain of the counter. As depicted in Fig. 2, $u_1(t)$ is passed through a filter, and the filter output $v_1(t)$ controls the frequency of a voltage-controlled oscillator (VCO_1). Let $\omega_1 = 2\pi f_1$ be the free-running radian frequency of VCO_1 , then the output of VCO_1 is

$$s_1(t) = \sin \left[\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_1 \right] \quad (2)$$

where α_1 is the gain of VCO_1 (radians/volt \times second). The signal $s_1(t)$ is used to control the transmission of data from Station 1 to Station 2 (for example, Station 1 transmits a digit to Station 2 at every 2nd zero crossing of $s_1(t)$). Note that θ_1 in (2) represents the phase of $s_1(t)$.

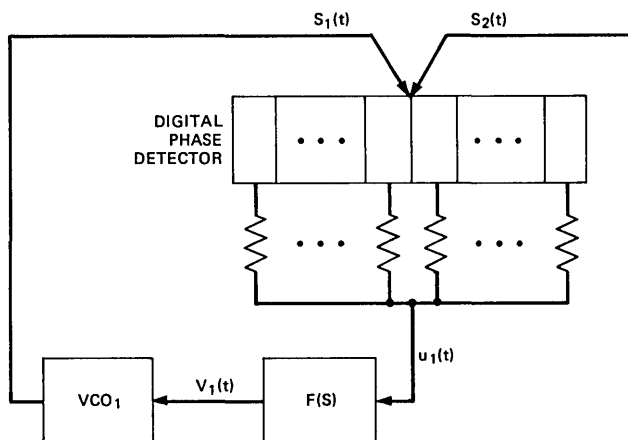


Fig. 2—Digital phase detector and the rate-locked loop at Station 1.

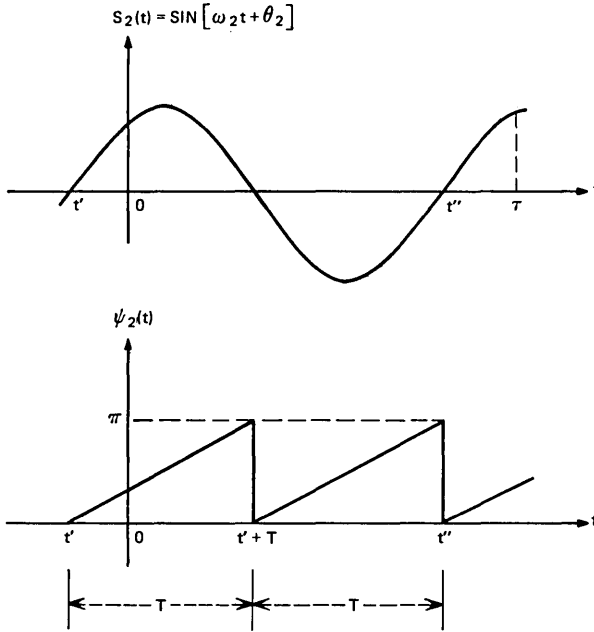


Fig. 3—Illustration of $S_2(t)$, $N_2(t)$, and $\psi_2(t)$.

at $t = 0$. Without loss of generality, we may assume $0 \leq \theta_1 \leq \pi$, and $0 \leq \theta_2 \leq \pi$.

Let us derive an analytic expression for the number of zero crossings of $s_2(t)$ from $t = 0$ to a particular time instant τ . As illustrated in Fig. 3, t' is the time instant at which the last zero crossing prior to $t = 0$ takes place, and t'' is the time instant at which the last zero crossing prior to $t = \tau$ takes place. It is obvious from Fig. 3 that the number of zero crossings in the time interval 0 to τ is

$$N_2(\tau) = \frac{\omega_2 t'' + \theta_2}{\pi} \tag{3}$$

Note from (3) that the phase cumulated from t'' to τ does not contribute to the value of $N_2(\tau)$. This residual phase (or phase quantization error) will be designated $\psi_2(\tau)$, i.e.,

$$\psi_2(\tau) = \omega_2 \tau - \omega_2 t'' \tag{4}$$

Equations (3) and (4) hold for all $\tau > 0$; therefore, we can replace their τ by the time variable t . The variation of $\psi_2(t)$ with t is illustrated in

Fig. 3. Note that $\psi_2(t)$ increases from 0 to π . When $\psi_2(t)$ reaches π radians, a zero crossing takes place; and $\psi_2(t)$ drops to zero and increases from zero again. Clearly, $0 \leq \psi_2(t) \leq \pi$. Since $s_2(t)$ is a pure sine wave, $\psi_2(t)$ is a sawtooth wave.

From (3) and (4), we have

$$N_2(t) = \frac{\omega_2 t + \theta_2 - \psi_2(t)}{\pi}. \quad (5)$$

Similarly, one can write the number of zero crossings of $s_1(t)$ as

$$N_1(t) = \frac{\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_1 - \psi_1(t)}{\pi} \quad (6)$$

where $\psi_1(t)$ is the residual phase as illustrated in Fig. 4. As can be seen, $0 \leq \psi_1(t) \leq \pi$. Note that $\psi_1(t)$ is not shown as a sawtooth wave in Fig. 4 because $s_1(t)$ is not a pure sine wave in the transient stage after $t = 0$.

In this paper, the filter $F(s)$ in Fig. 2 is assumed to be the usual RC filter (Fig. 5). Thus, its transfer function $F(s)$ is $1/(1 + sCR)$. Substituting (5) and (6) into (1), and rearranging the equation, we obtain

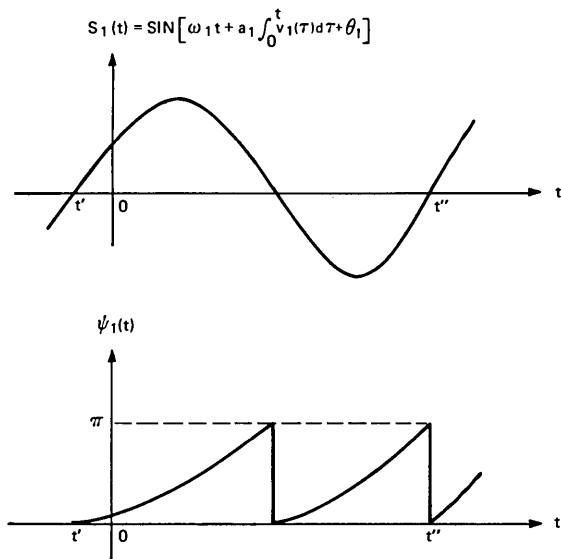


Fig. 4—Illustration of $S_1(t)$, $N_1(t)$, and $\psi_1(t)$.

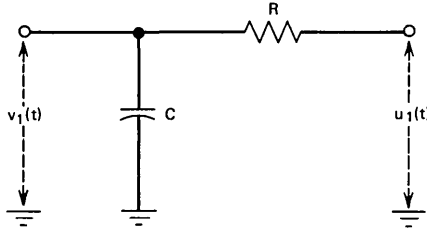


Fig. 5—RC filter in the rate-locked loop.

$$u_1(t) = k \left[\delta t - \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_2 - \theta_1 + \psi_1(t) - \psi_2(t) \right] \tag{7}$$

where

$$k = \frac{e_1}{\pi} \tag{8}$$

$$\delta = \omega_2 - \omega_1 . \tag{9}$$

We shall use one-sided Laplace transform in the analysis (the words one-sided will be omitted). As usual, the Laplace transform of a time function will be consistently denoted by the appropriate capital letter. For instance, $U_1(s)$ will denote the Laplace transform of $u_1(t)$. The symbol $\mathcal{L}[f(t)]$ denotes the Laplace transform of $f(t)$, and the symbol $\mathcal{L}^{-1}[F(s)]$ denotes the inverse Laplace transform of $F(s)$. Taking the Laplace transform of (7), we obtain

$$U_1(s) = k \left[\frac{\delta}{s^2} - \alpha_1 \frac{V_1(s)}{s} + \frac{\theta_2}{s} - \frac{\theta_1}{s} + \Psi_1(s) - \Psi_2(s) \right]. \tag{10}$$

Multiplying both sides of (10) by $F(s)$, using $F(s)U_1(s) = V_1(s)$, and rewriting the resulting equation in time domain we obtain

$$v_1(t) = \mathcal{L}^{-1} \left[H(s) \frac{\delta}{s^2} \right] + \mathcal{L}^{-1} \left[H(s) \frac{\theta_2}{s} \right] - \mathcal{L}^{-1} \left[H(s) \frac{\theta_1}{s} \right] \\ - \mathcal{L}^{-1} [H(s)\Psi_2(s)] + \mathcal{L}^{-1} [H(s)\Psi_1(s)], \quad t > 0 \tag{11}$$

where

$$H(s) = \frac{ks}{CR(s + r_1)(s + r_2)}$$

$$r_1 = \frac{1 + \sqrt{1 - 4CRk\alpha_1}}{2CR}$$

$$r_2 = \frac{1 - \sqrt{1 - 4CRk\alpha_1}}{2CR}.$$

The two roots r_1 and r_2 are real numbers when $1 - 4CRk\alpha_1 > 0$, and are complex numbers when $1 - 4CRk\alpha_1 < 0$. It can be shown that in the second case the frequency of $s_1(t)$ overshoots that of $s_2(t)$ before it finally settles. Such an overshooting should be avoided because Station 1 may be required to operate in the slave-to-slave mode (see Section V). Therefore, throughout this paper we assume

$$1 - 4CRk\alpha_1 > 0. \quad (12)$$

III. STEADY-STATE AND TRANSIENT ANALYSES

In the master-to-slave mode, we have to consider the following questions:

- (i) Can the signaling rate of Station 1 lock to that of Station 2 in the presence of phase quantization errors?
- (ii) What is the steady-state frequency of VCO₁?
- (iii) During the transient stage after $t = 0$, the signaling rate of Station 1 can be higher than that of Station 2. Therefore, data can be transmitted from Station 1 to Station 2 faster than it can be transmitted out of Station 2. Consequently, a buffer storage is required at Station 2. What should be the size of this buffer?
- (iv) The digital phase detector is a counter that counts the difference between the number of zero crossings of $s_1(t)$ and $s_2(t)$. How many stages are required in the counter to avoid overflow (i.e., to ensure pulling in)?

We shall first determine the transient response of the system in Section 3.1, and then consider these questions in Sections 3.2 to 3.4.

3.1 Transient Response and Settling Time

We evaluate the first three inverse transforms in (11) to obtain

$$\mathcal{L}^{-1} \left[H(s) \frac{\delta}{s^2} \right] = \frac{k\delta}{CRr_1r_2(r_1 - r_2)} [(r_1 - r_2) + r_2e^{-r_1t} - r_1e^{-r_2t}] \quad (13)$$

$$\mathcal{L}^{-1}\left[H(s)\frac{\theta_2}{s}\right] - \mathcal{L}^{-1}\left[H(s)\frac{\theta_1}{s}\right] = \frac{(\theta_2 - \theta_1)k}{CR(r_1 - r_2)} [e^{-r_2t} - e^{-r_1t}]. \tag{14}$$

The fourth inverse transform $\mathcal{L}^{-1}[H(s)\Psi_2(s)]$ is more difficult to evaluate. We obtain after lengthy manipulations

$$\mathcal{L}^{-1}[H(s)\Psi_2(s)] = a(t) + p(t) \tag{15}$$

where

$$a(t) = \frac{k}{\sqrt{1 - 4CRk\alpha_1}} \left[\frac{\omega_2}{r_1} - \theta_2 + \frac{\pi e^{r_1t}}{1 - e^{r_1T}} \right] e^{-r_1t} + \frac{k}{\sqrt{1 - 4CRk\alpha_1}} \left[-\frac{\omega_2}{r_2} + \theta_2 - \frac{\pi e^{r_2t}}{1 - e^{r_2T}} \right] e^{-r_2t}. \tag{16}$$

$p(t)$ = a periodic function of period T , identical with $p_0(t)$ in the time period $0 \leq t \leq T$ (17)

$$p_0(t) = \frac{\omega_2}{\alpha_1} + \frac{k\pi e^{r_1t}}{\sqrt{1 - 4CRk\alpha_1} (1 - e^{r_1T})} [-1 + (1 - e^{r_1T})u(t - t_s)] e^{-r_1t} + \frac{k\pi e^{r_2t}}{\sqrt{1 - 4CRk\alpha_1} (1 - e^{r_2T})} [1 - (1 - e^{r_2T})u(t - t_s)] e^{-r_2t}. \tag{18}$$

Note that $a(t)$ is the sum of two decaying exponential terms. When t increases, $a(t)$ approaches zero, and only the periodic steady-state response $p(t)$ remains. It can be shown that $p(t)$ has zero mean.

Now consider the last inverse transform $\mathcal{L}^{-1}[H(s)\Psi_1(s)]$ in (11). Since $\Psi_1(s)$ is the Laplace transform of $\psi_1(t)$, $\mathcal{L}^{-1}[H(s)\Psi_1(s)]$ can be evaluated if $\psi_1(t)$ can be determined. As illustrated in Fig. 4, $\psi_1(t)$ depends on the positions of the zero crossings of $s_1(t)$. Furthermore, because we are dealing with a closed-loop control system, $\psi_1(t)$ and the phase of $s_1(t)$ must satisfy the integral equation (7). In order to determine $\psi_1(t)$, one must simultaneously consider (7) and the zero crossings of $s_1(t)$. The mathematical problem is extremely complex and it is impossible to obtain a closed-form expression of $\psi_1(t)$ for all t . Consequently, the inverse transform $\mathcal{L}^{-1}[H(s)\Psi_1(s)]$ cannot be evaluated in closed form. However, we have obtained a tight upper and lower bound for its value as follows:

$$|\mathcal{L}^{-1}[H(s)\Psi_1(s)]| < 2e_1. \tag{19}$$

We have obtained the closed-form expression of the first four components of $v_1(t)$, and tightly bounded the fifth component of $v_1(t)$. This gives the transient response of the system. Note from (13) to (19) that transients in $v_1(t)$ either decay exponentially or can be bounded by a small number. Thus, the system is well behaved in the transient stage. Furthermore, from these equations, one can plot $v_1(t)$ vs t , and easily estimate the settling time of VCO_1 . The settling time of VCO_1 can be rather long when CR is large. For example, consider the first term $\mathcal{L}^{-1}[H(s)\delta/s^2]$ in $v_1(t)$ (this is usually the dominating term in $v_1(t)$). From (13) it can be rewritten as

$$\mathcal{L}^{-1}\left[H(s)\frac{\delta}{s^2}\right] = \frac{\delta}{\alpha_1} \left[1 - e^{-r_1 t} + \frac{r_1(e^{-r_2 t} - e^{-r_1 t})}{r_2 - r_1} \right].$$

Since $r_1 > r_2 > 0$, the last term in the right-side bracket is negative for all t . Thus, the convergence of $\mathcal{L}^{-1}[H(s)\delta/s^2]$ is even slower than the convergence of the time function $1 - e^{-r_1 t}$. This clearly shows that $v_1(t)$ converges slowly when the filter time constant CR is large.

3.2 Steady-State Frequency of VCO_1

Now we answer the first two questions at the beginning of Section III. First, we have found that the signaling rate of Station 1 will lock to that of Station 2 in the presence of phase quantization errors. The proof of this is complicated, and is given in the Appendix. In this section, we examine the steady-state frequency of VCO_1 , and point out an important difference between digital and analog phase detectors. The instantaneous frequency of VCO_1 is $[\omega_1 + \alpha_1 v_1(t)/2\pi]$. In order to see if it approaches a fixed steady-state value, we evaluate $\lim_{t \rightarrow \infty} v_1(t)$, which can be found by evaluating the limits of the five inverse transforms in (11). As shown in the Appendix, when signaling rate of Station 1 locks to that of Station 2, the zero crossings of $s_1(t)$ and $s_2(t)$ will alternate with probability one, and $\psi_1(t)$ will be a periodic function of period T . This means that $\mathcal{L}^{-1}[H(s)\psi_1(s)]$ also approaches a periodic function of period T . Let this periodic function be denoted by $q(t)$. Then, one can show from (11) that

$$\left(\begin{array}{l} \text{Instantaneous fre-} \\ \text{quency of } VCO_1 \end{array} \right) = f_2 - \frac{\alpha_1}{2\pi} p(t) + \frac{\alpha_1}{2\pi} q(t). \quad (20)$$

When signaling rate of Station 1 locks to that of Station 2, zero crossings of $s_1(t)$ and $s_2(t)$ alternate. Therefore, $\psi_1(t) \neq \psi_2(t)$ and $p(t) \neq q(t)$. Thus, from (20), instantaneous frequency of VCO_1 does

not lock to the master clock frequency f_2 . Instead, it is a periodic function $f_2 - (\alpha_1/2\pi)p(t) + (\alpha_1/2\pi)q(t)$. The output $s_1(t)$ of VCO_1 is a periodic wave with the same period as $s_2(t)$; however, it is not a pure sine wave as one would expect from experience with analog-phase lock loops. *The digital loop locks to neither the instantaneous frequency nor the phase of the incoming signal $s_2(t)$. It locks only to the rate of zero crossings of $s_2(t)$.* For this reason, it should be referred to as a digital rate-locked loop, rather than a digital frequency-locked loop or a digital phase-locked loop. This difference between digital and analog loops should be noted in the applications.

3.3 Size of the Data Buffer at Station 2

As described in Section II, Station 2 transmits to Station 1 at a standard rate of f_2 digits per second. In general, Station 1 is also required to transmit to Station 2 at this standard rate. To achieve this, Station 1 transmits a digit to Station 2 at every second zero crossing of $s_1(t)$ [this enables station 1 to transmit also at the standard rate when $s_1(t)$ is synchronized to $s_2(t)$].

Usually, Station 2 relays the data it receives from Station 1 to another station at the standard rate of f_2 digits per second. Thus, when the system is in synchronization, data is transmitted to Station 2 at the same rate as it is transmitted out of Station 2. However, when Station 1 is first synchronized (that is, during the transient stage of synchronization), the transmission rate of Station 1 can be higher than f_2 . Consequently, during the transient stage, data can be transmitted from Station 1 to Station 2 faster than it can be transmitted out of Station 2. This means a data buffer is required at Station 2. In this section, we determine the size of this buffer.

As defined in Section II, $N_1(t)$ is the number of zero crossings of $s_1(t)$ in the time interval 0 to t . Since Station 1 transmits a digit to station 2 at every second zero crossing of $s_1(t)$, the number of digits transmitted from Station 1 to Station 2 in the time interval 0 to t is $N_1(t)/2$ or $N_1(t) - 1/2$, depending on whether $N_1(t)$ is even or odd. To simplify our discussions, we shall use the following definition throughout this paper.

Definition: For any positive number a , $\langle a \rangle$ denotes the integer immediately less than a when a is not an integer. $\langle a \rangle = a$ when a is an integer.

Using this definition, the number of digits transmitted from Station 1 to Station 2 in the time interval 0 to t is $\langle N_1(t)/2 \rangle$. The number of digits Station 1 should transmit in this time interval is $\langle N_2(t)/2 \rangle$. If

$\langle N_1(t)/2 \rangle$ is larger than $\langle N_2(t)/2 \rangle$, a buffer would be required at Station 2 and the buffer size is $\langle N_1(t)/2 \rangle - \langle N_2(t)/2 \rangle$ digits. It can be shown from the previous equations that

$$-\left[\frac{\omega'}{e_1\alpha_1} + 3 \right] < N_1(t) - N_2(t) < \frac{\omega'}{e_1\alpha_1} + 3. \quad (21)$$

Since the buffer size is $\langle N_1(t)/2 \rangle - \langle N_2(t)/2 \rangle$, we obtain from (21)

$$\text{Buffer Size} < \frac{\omega'}{2e_1\alpha_1} + 2. \quad (22)$$

Equation (22) gives an upper bound for the buffer size. It can also be shown that in order to prevent overflow the buffer size must be greater than $(\omega'/2e_1\alpha_1) - \frac{1}{2}$. Combining this with (22), we have

$$\frac{\omega'}{2e_1\alpha_1} - \frac{1}{2} < \text{Buffer Size} < \frac{\omega'}{2e_1\alpha_1} + 2. \quad (23)$$

Let us define $B = \langle \omega'/2e_1\alpha_1 + 2 \rangle$. It can be seen from (23) that the buffer size is B , $B - 1$, or $B - 2$. Thus, the buffer size is determined to within two digits. Since the two-digit difference is negligible, one may use the simple formula

$$\text{Buffer Size} = B = \left\langle \frac{\omega'}{2e_1\alpha_1} + 2 \right\rangle \text{ digits}. \quad (24)$$

As explained at the end of Section II, in this paper we use the constraint $1 - 4CRk\alpha_1 > 0$. From this constraint, we can rewrite (24) as

$$\text{Buffer Size} = \left\langle \frac{2CR\omega'}{\beta_1\pi} + 2 \right\rangle \text{ digits} \quad (25)$$

where $\beta_1 = 4CR e_1\alpha_1/\pi$. Clearly, $0 < \beta_1 < 1$. Equation (25) will be used in later discussions.

3.4 Counter Size of the Digital Phase Detector

The counter in the digital phase detector counts the difference between $N_2(t)$ and $N_1(t)$. Now we determine the counter size so that the counter will not overflow when the maximum positive count or negative count is reached. Consider first the case of negative counts. It can be shown that $N_1(t) - N_2(t)$ can be larger than $\omega'/e_1\alpha_1$. It has been shown in the preceding subsection that $N_1(t) - N_2(t)$ must be less than $\omega'/e_1\alpha_1 + 3$. Thus, if we define

$$N = \left\langle \frac{\omega'}{e_1\alpha_1} + 3 \right\rangle = \left\langle \frac{4CR\omega'}{\beta_1\pi} + 3 \right\rangle,$$

the integer $N_1(t) - N_2(t)$ can be as large as $N - 2$, but will not exceed N . Therefore, the counter will not overflow if it can count $N_1(t) - N_2(t)$ up to $N_1(t) - N_2(t) = N$.

Next, consider positive counts. One can show that the counter will not overflow if it can count $N_2(t) - N_1(t)$ up to $N_2(t) - N_1(t) = N$. Combining the two cases, we see that the counter will not overflow if

$$\text{Counter Size} = \pm N \text{ counts} = \pm \left\langle \frac{4CR\omega'}{\beta_1\pi} + 3 \right\rangle \text{ counts} \quad (26)$$

where β_1 is defined after (25).

IV. SLAVE-TO-SLAVE SYNCHRONIZATION USING DIGITAL RATE-LOCKED LOOPS

In this section, we consider mutual synchronization between two data stations where each station regards the clock at the other station as the master clock. Such a mutual synchronization is usually called slave-to-slave synchronization.

A mathematical model of slave-to-slave synchronization is depicted in Fig. 6. The local oscillator at Station 1 (VCO_1 in Fig. 6) emits a timing signal $S_{11}(t)$ which controls the transmission of data from

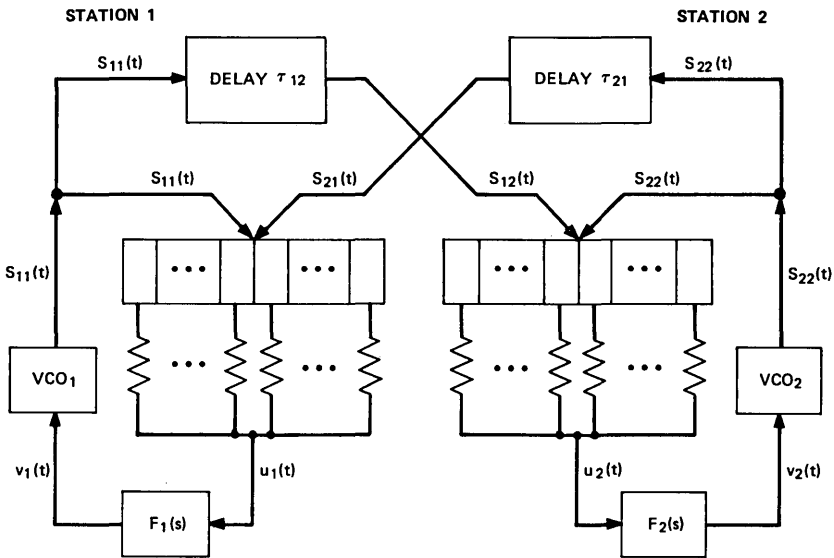


Fig. 6—Slave-to-slave synchronization with digital rate-locked loops at both stations.

Station 1 to Station 2. [For example, Station 1 may transmit a digit to Station 2 at every second zero crossing of $S_{11}(t)$.] Station 2 derives from the received data a timing signal $S_{12}(t)$ and compares $S_{12}(t)$ with its local oscillator output $S_{22}(t)$ at the digital phase detector. The digital phase detector is essentially a counter which counts the zero crossings of $S_{12}(t)$ and $S_{22}(t)$ and produces an error signal $u_2(t)$ proportional to the difference between these two counts. The error signal $u_2(t)$ is passed through a filter $F_2(s)$ to control the frequency of VCO₂. Thus, in this fashion, Station 2 adjusts its clock rate toward that of Station 1. Similarly, as depicted in Figure 6, Station 1 regards the clock at Station 2 as the master clock and adjusts its clock rate toward that of Station 2.

Practically, it is impossible to activate the two counters at the two different stations at the same time instant. Therefore, in this study, we consider an arbitrary starting sequence as follows:

- (i) At an arbitrary time instant t_1 , either the counter at Station 1 or the counter at Station 2 is activated.
- (ii) The other counter is activated at an arbitrary later time instant t_2 ($t_2 > t_1$).

For analytical purpose, we shift the time origin such that $t_2 = 0$. We shall analyze the behavior of the system for $t > 0$.

Let ω_1 be the free-running radian frequency of VCO₁, then we can write

$$\begin{aligned} s_{11}(t) &= \sin \rho_{11}(t) \\ &= \sin \left[\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_{11} \right] \end{aligned} \quad (27)$$

and

$$s_{12}(t) = \sin [\rho_{11}(t - \tau_{12})] \quad (28)$$

where τ_{12} is the time delay introduced by the channel. Similarly, the free-running frequency of VCO₂ is denoted ω_2 and

$$\begin{aligned} s_{22}(t) &= \sin \rho_{22}(t) \\ &= \sin \left[\omega_2 t + \alpha_2 \int_0^t v_2(\tau) d\tau + \theta_{22} \right] \end{aligned} \quad (29)$$

and

$$s_{21}(t) = \sin [\rho_{22}(t - \tau_{21})]. \quad (30)$$

We define $N_{i_j}(t)$ as the number of counts from $s_{i_j}(t)$ in the time interval 0 to t . Let us first derive an analytical expression for $N_{11}(t)$. From (27), we can write

$$\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_{11} = M\pi N_{11}(t) + \psi_{11}(t) \tag{31}$$

where the parameter M is defined by: M equals one when the counters at the two stations count both the upward and downward zero crossings; M equals two when the counters count only the upward (or downward) zero crossings. The last term $\psi_{11}(t)$ in the above equation represents phase quantization errors and $0 \leq \psi_{11}(t) < M\pi$. From the above equation we have

$$N_{11}(t) = \frac{1}{M\pi} \left[\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_{11} - \psi_{11}(t) \right]. \tag{32}$$

Similarly, one can write analytical expressions for $N_{12}(t)$, $N_{21}(t)$, and $N_{22}(t)$.

At Station 1, the digital counter output is

$$u_1(t) = u_1(0) + e_1[N_{21}(t) - N_{11}(t)] \tag{33}$$

where $u_1(0)$ is the initial count at $t = 0$. Let us define

$$k_1 = \frac{e_1}{M\pi} \tag{34}$$

$$\delta = \omega_2 - \omega_1 \tag{35}$$

$$\theta_1 = \theta_{21} - \theta_{11} \tag{36}$$

$$R_{21}(s) = e^{-s\tau_{21}} \int_{-\tau_{21}}^0 v_2(t)e^{-st} dt. \tag{37}$$

The filters $F_1(s)$ and $F_2(s)$ in Fig. 6 are assumed to be the usual RC filters, i.e.,

$$F_1(s) = \frac{1}{1 + sC_1R_1} \tag{38}$$

$$F_2(s) = \frac{1}{1 + sC_2R_2}. \tag{39}$$

At Station 2, the digital filter output is

$$u_2(t) = u_2(0) + e_2[N_{12}(t) - N_{22}(t)] \tag{40}$$

where $u_2(0)$ is the initial count at $t = 0$. Furthermore, we define

$$k_2 = \frac{e_2}{M\pi} \quad (41)$$

$$\theta_2 = \theta_{12} - \theta_{22} \quad (42)$$

$$R_{12}(s) = e^{-s\tau_{12}} \int_{-\tau_{12}}^0 v_1(t)e^{-st} dt. \quad (43)$$

From the previous equations, we can determine $V_1(s)$ and $V_2(s)$. The results are:

$$V_1(s) = \frac{A_1(s)}{B(s)} + \frac{A_2(s)}{B(s)} + \frac{A_3(s)}{B(s)} + \frac{A_4(s)}{B(s)} + \frac{A_5(s)}{B(s)} \quad (44)$$

where

$$A_1(s) = \left[\frac{u_1(0)F_1(s)}{s} + \frac{k_1\theta_1F_1(s)}{s} + v_1(0)C_1R_1F_1(s) + \frac{k_1\alpha_2R_{21}(s)F_1(s)}{s} \right] \cdot \left[1 + k_2\alpha_2 \frac{F_2(s)}{s} \right]$$

$$A_2(s) = [-k_1\Psi_{21}(s)F_1(s) + k_1\Psi_{11}(s)F_1(s)] \left[1 + k_2\alpha_2 \frac{F_2(s)}{s} \right]$$

$$A_3(s) = \left[\frac{u_2(0)F_2(s)}{s} + \frac{k_2\theta_2F_2(s)}{s} + v_2(0)C_2R_2F_2(s) + \frac{k_2\alpha_1R_{12}(s)F_2(s)}{s} \right] \cdot \left[k_1\alpha_2 \frac{e^{-s\tau_{21}}}{s} F_1(s) \right]$$

$$A_4(s) = [-k_2\Psi_{12}(s)F_2(s) + k_2\Psi_{22}(s)F_2(s)] \left[k_1\alpha_2 \frac{e^{-s\tau_{21}}}{s} F_1(s) \right]$$

$$A_5(s) = k_1 \delta \frac{F_1(s)}{s^2} \left[1 + k_2\alpha_2 \frac{F_2(s)}{s} (1 - e^{-s\tau_{21}}) \right]$$

$$B(s) = 1 + k_1\alpha_1 \frac{F_1(s)}{s} + k_2\alpha_2 \frac{F_2(s)}{s} + k_1\alpha_1 k_2\alpha_2 \frac{F_1(s)}{s} \frac{F_2(s)}{s} [1 - e^{-s(\tau_{12} + \tau_{21})}].$$

Similarly, we obtain

$$V_2(s) = \frac{A_6(s)}{B(s)} + \frac{A_7(s)}{B(s)} + \frac{A_8(s)}{B(s)} + \frac{A_9(s)}{B(s)} + \frac{A_{10}(s)}{B(s)} \quad (45)$$

where

$$A_6(s) = \left[\frac{u_2(0)F_2(s)}{s} + \frac{k_2\theta_2F_2(s)}{s} + v_2(0)C_2R_2F_2(s) + \frac{k_2\alpha_1R_{12}(s)F_2(s)}{s} \right] \cdot \left[1 + k_1\alpha_1 \frac{F_1(s)}{s} \right]$$

$$A_7(s) = [-k_2\Psi_{12}(s)F_2(s) + k_2\Psi_{22}(s)F_2(s)] \left[1 + k_1\alpha_1 \frac{F_1(s)}{s} \right]$$

$$A_8(s) = \left[\frac{u_1(0)F_1(s)}{s} + \frac{k_1\theta_1F_1(s)}{s} + v_1(0)C_1R_1F_1(s) + \frac{k_1\alpha_2R_{21}(s)F_1(s)}{s} \right] \cdot \left[k_2\alpha_1 e^{-s\tau_{12}} \frac{F_2(s)}{s} \right]$$

$$A_9(s) = [-k_1\Psi_{21}(s)F_1(s) + k_1\Psi_{11}(s)F_1(s)] \left[k_2\alpha_1 e^{-s\tau_{12}} \frac{F_2(s)}{s} \right]$$

$$A_{10}(s) = -k_2 \delta \frac{F_2(s)}{s^2} \left[1 + k_1\alpha_1 \frac{F_1(s)}{s} (1 - e^{-s\tau_{12}}) \right]$$

Note that our problem is not solved. Equation (44) is not a closed-form solution of $V_1(s)$ because $\Psi_{i,j}(s)$, which appears in $A_2(s)$ and $A_4(s)$, depends on $V_1(s)$ and $V_2(s)$ (the phase quantization errors $\psi_{i,j}(t)$ depends on $v_1(t)$ and $v_2(t)$). Similarly, (45) is not a closed-form solution of $V_2(s)$. These equations will, however, enable us to examine the steady-state behavior of the system in the following sections.

4.1 A Steady-State Solution of Signaling Rates

As described previously, the zero crossings of $s_{11}(t)$ are used to control the transmission from Station 1 to Station 2 (for example, Station 1 may transmit a digit to Station 2 at every second zero crossing of $s_{11}(t)$). Similarly, the zero crossings of $s_{22}(t)$ are used to control the transmission from Station 2 to Station 1. Therefore, to determine the steady-state signaling rates of these two stations, it suffices to determine the steady-state distribution of the zero crossings of $s_{11}(t)$ and $s_{22}(t)$. To facilitate our discussion, let us first introduce the following definition. *Definition:* $s_0(t)$ denotes a sine wave $\sin \omega_0 t$ with

$$\omega_0 = \frac{1}{k_1\alpha_1 + k_2\alpha_2 + k_1\alpha_1k_2\alpha_2(\tau_{12} + \tau_{21})} \left[\omega_1k_2\alpha_2 + \omega_2k_1\alpha_1 + \omega_1k_1\alpha_1k_2\alpha_2\tau_{12} + \omega_2k_1\alpha_1k_2\alpha_2\tau_{21} + [u_1(0) + k_1\theta_1 + k_1\alpha_2R_{21}(0)]k_2\alpha_1\alpha_2 + [u_2(0) + k_2\theta_2 + k_2\alpha_1R_{12}(0)]k_1\alpha_1\alpha_2 \right]. \tag{46}$$

Based on (44) and (45), a steady-state solution of the zero-crossing distribution has been obtained.¹⁰ In order to conserve space, let us omit the lengthy derivations and write only the results as follows: *A Steady-State Solution:* When the counters at the two stations count both the upward and downward zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, the upward and downward zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, are uniformly spaced when $t \rightarrow \infty$ and the time interval between each two consecutive zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, is identical with the time interval between each two consecutive zero crossings of $s_0(t)$.

If the counters count only the upward (or downward) zero crossings, the above solution should be modified: When the counters at the two stations count only the upward (or downward) zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, the upward (or downward) zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, are uniformly spaced when $t \rightarrow \infty$, and the time interval between each two consecutive upward (or downward) zero crossings of $s_{ij}(t)$, $i, j = 1, 2$, is identical with the time interval between each two consecutive upward (or downward) zero crossings of $s_0(t)$.

4.2 Analysis of the Steady-State Signaling Rate

In this section, we show that the system can be easily designed such that the steady-state signaling rate lies within desired limits.

Before the two stations are mutually synchronized, $s_{11}(t)$ is $\sin \omega_1 t$ and the signaling rate of Station 1 is $h\omega_1$ digits/second. (h is a proportionality constant. For example, $h = 1/2\pi$ when Station 1 transmits a digit at every second zero crossing of $s_{11}(t)$.) Similarly, before the two stations are synchronized, $s_{22}(t)$ is $\sin \omega_2 t$ and the signaling rate of Station 2 is $h\omega_2$ digits/second. When the two stations are mutually synchronized, $s_{11}(t)$ and $s_{22}(t)$ have the same zero-crossing distribution as $s_0(t) = \sin \omega_0 t$ and the signaling rates of the two stations are $h\omega_0$ digits/second. The synchronization is satisfactory if $h\omega_0$ is sufficiently close to $h\omega_1$ or $h\omega_2$. More specifically, the steady-state signaling rate is satisfactory if

$$h\omega_1 - \epsilon < h\omega_0 < h\omega_2 + \epsilon \quad (47)$$

when $\omega_1 < \omega_2$, and if

$$h\omega_2 - \epsilon < h\omega_0 < h\omega_1 + \epsilon \quad (48)$$

when $\omega_2 < \omega_1$. The number ϵ is a prescribed small number.

As can be seen from (46), ω_0 depends on ω_1, ω_2 , and the following parameters: gains e_1 and e_2 of the two counters, gains α_1 and α_2 of the two oscillators, initial counter outputs $u_1(0)$ and $u_2(0)$, initial phases

θ_1 and θ_2 , initial filter outputs $v_1(t)$ and $v_2(t)$, and the time delays τ_{12} and τ_{21} in the communication channels. Since ω_o depends on so many parameters, it is not immediately clear whether ω_o satisfies the specifications in (47) and (48). In the following, we derive simple design constraints such that when these constraints are satisfied, ω_o will satisfy the specifications in (47) and (48).

So far, we have considered the arbitrary starting sequence described at the beginning of this section. Since we may always designate the station that is started first as Station 1, we need to consider only the following starting sequence in the sequel: At an arbitrary time $t_1 < 0$, the counter at Station 1 is activated. The counter at Station 2 is activated at $t = 0$.

There are two cases to be considered: $\omega_1 \leq \omega_2$ and $\omega_1 > \omega_2$. Our analyses of these two cases yield the same design constraint; hence, we describe only the case $\omega_1 \leq \omega_2$.

Note from the starting sequence that for $t \leq 0$, Station 2 is the master and Station 1 is the slave. We therefore can use the results in Section III to bound $v_1(t)$ for $t \leq 0$. From this, we can show that ω_o always satisfies the following inequalities:

$$\omega_o > \left[\begin{aligned} &\omega_1 + \frac{(k_1\alpha_1 + k_1\alpha_1k_2\alpha_2\tau_{21})(\omega_2 - \omega_1)}{k_1\alpha_1 + k_2\alpha_2 + k_1\alpha_1k_2\alpha_2(\tau_{12} + \tau_{21})} \\ &- \frac{[\frac{5}{6}k_2\alpha_2 + k_1\alpha_1k_2\alpha_2\tau_{12}]6e_1\alpha_1}{k_1\alpha_1 + k_2\alpha_2 + k_1\alpha_1k_2\alpha_2(\tau_{12} + \tau_{21})} \end{aligned} \right] \quad (49)$$

and

$$\omega_o < \omega_2 + \frac{[\frac{5}{6}k_2\alpha_2 + k_1\alpha_1k_2\alpha_2\tau_{12}]6e_1\alpha_1}{k_1\alpha_1 + k_2\alpha_2 + k_1\alpha_1k_2\alpha_2(\tau_{12} + \tau_{21})}. \quad (50)$$

It should be clear from (49) and (50) that, regardless of the values of the time delays τ_{12} and τ_{21} , one can easily select the gain $e_1\alpha_1$ of the first station so that ω_o will satisfy the constraint in (47). To show this more explicitly, we further simplify (49) and (50) (this simplification will, however, make the constraint on $e_1\alpha_1$ slightly more stringent). Since ω_o satisfies (49) and (50), ω_o will definitely lie in the following broader range

$$\omega_1 - 6e_1\alpha_1 < \omega_o < \omega_2 + 6e_1\alpha_1. \quad (51)$$

Comparing (51) with (47) shows that ω_o satisfies the specification in (47) if

$$e_1\alpha_1 < \frac{\epsilon}{6h}. \quad (52)$$

From (52), one can easily determine the value of $e_1\alpha_1$. Since we have designated the station that is started first as Station 1, and since either station can be started first, (52) should be applied to both stations. To emphasize this, we replace (52) with the following two constraints

$$e_1\alpha_1 < \frac{\epsilon}{6h} \quad (53)$$

$$e_2\alpha_2 < \frac{\epsilon}{6h}. \quad (54)$$

Now, to summarize this section: we have shown that if the gains of the two stations are designed to satisfy the simple constraints in (53) and (54), ω_o will satisfy (47) and the steady-state signaling rates will be satisfactory. Since (53) and (54) can be easily satisfied, and are independent of all the other parameters in (46), we conclude that the steady-state signaling rate can be easily made satisfactory regardless of the starting sequence, the initial system conditions, and the time delays in the communication channels.

V. SUMMARY AND CONCLUSIONS

Sections II and III examine the behavior of the system in the master-to-slave mode. The station with the slave clock (Station 1 in Fig. 1) represents a data station, while the station with the master clock (Station 2) represents a station in the digital transmission facility. The slave clock at Station 1 employs a digital rate-locked loop which consists of a digital counter, an RC filter, and a slave oscillator (Fig. 2). The counter is not restricted to have only one stage. A mathematical model of the system is formulated in Section II. Transient response of the system is determined in Section 3.1. It is shown that, under the condition $1 - 4CRk\alpha_1 > 0$ in (12), the signaling rate of Station 1 approaches that of Station 2 in a monotone fashion (transients either decay exponentially as shown in (13), (14), (15) and (16), or can be tightly bounded as shown in (19)).

From the transient response, settling time of the slave oscillator can be easily estimated. As discussed at the end of Section 3.1, this settling time can be rather long when the RC filter has a large time constant. For fast start-up purpose, it may be desirable for Station 1 to transmit data before the slave oscillator is completely settled. Thus, during the start-up period, data can be transmitted from Station 1 to Station 2 faster than it can be transmitted out of Station 2. Con-

sequently, a buffer storage is required at Station 2. This buffer size is determined and is given in (25). Section 3.4 examines the size of the counter in the rate-locked loop (counter size determines the pull-in range of the rate-locked loop). In order to avoid counter overflow (that is, to ensure pulling in), the counter must have a certain minimum size. This minimum size is determined and is given in (26).

As emphasized in Section 3.2, the slave oscillator in the rate-locked loop locks to neither the instantaneous frequency nor the phase of the master oscillator. It locks only to the rate of zero crossings of the master oscillator. For this reason, we refer to this control loop as a rate-locked loop, instead of a frequency-locked loop or a phase-locked loop. This difference, while immaterial in the present application, should be carefully noted in other applications.

Section IV examines the behavior of the system in the slave-to-slave mode. The two stations to be mutually synchronized represent two data stations connected by analog transmission facilities. A rate-locked loop is used at each station, and an RC filter is included in each loop. A random starting sequence is considered where either station can be started first, with the other station activated at an arbitrary later time. When the two stations are mutually synchronized, the two stations settle to the same steady-state signaling rate $h\omega_o$ (h is a proportionality constant and ω_o is given in (46)). Equation (46) shows that ω_o depends on the gains of the counters and oscillators, the initial conditions of the counters, filters, and oscillators, and the time delays in the communication channels. It is shown that, although ω_o depends on so many parameters, the steady-state signaling rate $h\omega_o$ will lie within desired limits if the simple design constraints in (53) and (54) are satisfied (these conditions can be relaxed by using the more complicated equations (49) and (50)). These results show that the steady-state signaling rate of the system can easily be made satisfactory regardless of the starting sequence, the initial system conditions, and the time delays in the communication channels. Therefore, there is no need to attempt to activate the two stations simultaneously or to equalize the delays and gains of the communication channels.

In conclusion, the detailed transient and steady-state analyses show that a synchronization system employing digital rate-locked loops can be designed to operate successfully both in the master-to-slave mode and in the slave-to-slave mode. Such a synchronization system, therefore, is useful in applications where both digital and analog transmission facilities are utilized in connecting data stations or other types of terminals.

VI. ACKNOWLEDGMENTS

The author gratefully acknowledges many helpful discussions with J. Salz, J. E. Mazo, R. R. Anderson, D. Hirsch, R. D. Fracassi, F. W. Lescinsky, and C. E. Young.

APPENDIX

In this appendix, we first introduce the concept of equilibrium. The system is said to be in equilibrium if, corresponding to every digit received from Station 2 (the station with master clock), Station 1 (the station with slave clock) also transmits a digit back to Station 2. Then we prove a general theorem which states that the system will reach equilibrium if the arbitrary filter $F(s)$ (not necessarily an RC filter) satisfies a simple condition. Based on this general theorem, we then show that when an RC filter is used, the signaling interval of Station 1 will lock to that of Station 2 exactly.

For brevity, we define $\rho_1(t)$ as $\omega_1 t + \alpha_1 \int_0^t v_1(\tau) d\tau + \theta_1$. The lowpass filter transfer function $F(s)$ can always be normalized such that $F(0) = 1$. Clearly, any useful lowpass filter must cut off as frequency approaches infinity; therefore, we can write $F(\infty) = 0$. By changing units, we can and shall set $e_1 = \pi$ and $\alpha_1 = 1$. Without loss of generality, we assume that $\omega_2 - \omega_1 > 0$, and that the counter counts both upward and downward zero crossings. The zero crossings of $s_1(t)$ and $s_2(t)$ control the signaling rate of Station 1 and Station 2, respectively. Let T be the time interval between each two consecutive zero crossings of $s_2(t)$, that is, $T = \pi/\omega_2$. When the time interval between each two consecutive zero crossings of $s_1(t)$ also becomes T , signaling rate of Station 1 locks to that of Station 2. Thus, to determine the locking behavior, we need only to examine $N_1(t)$ when $t \rightarrow \infty$. Since $N_1(t)$ can be deduced from $\rho_1(t)$, $v_1(t)$, or $u_1(t)$, we shall examine either $\rho_1(t)$, or $v_1(t)$, or $u_1(t)$ in the following analysis (depending on which one is the most convenient).

The behavior of the system is governed by the equation

$$\rho_1(t) = \omega_1 t + \int_0^t [f * u_1(\rho_1)] d\tau + \theta_1 \quad (55)$$

where $*$ denotes convolution, and the symbol $u_1(\rho_1)$ indicates that u_1 is a function of ρ_1 . Since u_1 depends on ρ_1 through the nonlinear zero-crossing counting process, (55) is a nonlinear differential-integral equation. It is impossible to solve this equation for all t , so we shall first examine $v_1(t)$ and $u_1(t)$ to obtain a steady-state solution of this

equation. Then we shall consider the uniqueness of this steady-state solution.

From the mathematical formulation in text,

$$u_1(t) = \left[(\omega_2 - \omega_1)t - \int_0^t v_1(\tau) d\tau + \theta_2 - \theta_1 + \psi_1(t) - \psi_2(t) \right]. \quad (56)$$

Let $U_1(s)$, $V_1(s)$, $\Psi_1(s)$, and $\Psi_2(s)$ be the Laplace transforms of $u_1(t)$, $v_1(t)$, $\psi_1(t)$, and $\psi_2(t)$, respectively. From (56) and $V_1(s) = F(s)U_1(s)$, we obtain

$$U_1(s) = \frac{\omega_2 - \omega_1}{s[s + F(s)]} + \frac{\theta_2 - \theta_1}{s + F(s)} + \frac{s}{s + F(s)} \Psi_1(s) - \frac{s}{s + F(s)} \Psi_2(s) \quad (57)$$

and

$$V_1(s) = \frac{F(s)(\omega_2 - \omega_1)}{[s + F(s)]s} + \frac{F(s)(\theta_2 - \theta_1)}{s + F(s)} + \frac{F(s)s}{s + F(s)} \Psi_1(s) - \frac{F(s)s}{s + F(s)} \Psi_2(s). \quad (58)$$

We wish to determine the $N_1(t)$ that satisfies the system equation (55) when $t \rightarrow \infty$. For brevity, such a solution is called a steady-state solution. From (58), a steady-state solution is obtained, and is stated in the following theorem.

Theorem 1: At steady-state (that is, when $t \rightarrow \infty$), (55) is satisfied if

$$N_1(t) = N_2(t - \tau_0) \quad (59)$$

where τ_0 is such that the mean value of $u_1(t)$ is $\omega_2 - \omega_1$.

Proof: Since $\psi_1(t)$ and $\psi_2(t)$ do not approach a limit when $t \rightarrow \infty$, one cannot apply final value theorem to the last two terms in (58). However, final value theorem can be applied to the first two terms. This yields

$$v_1(t) = \omega_2 - \omega_1 + \mathcal{L}^{-1} \left[\frac{F(s)s}{s + F(s)} \Psi_1(s) \right] - \mathcal{L}^{-1} \left[\frac{F(s)s}{s + F(s)} \Psi_2(s) \right], \quad t \rightarrow \infty. \quad (60)$$

The condition $t \rightarrow \infty$ applies to the rest of the proof. Clearly, (59) is equivalent to the statement that

$$\rho_1(t) = \omega_2 t + \tilde{\rho}_1(t) \quad (61)$$

where $\tilde{\rho}_1(t)$ is a periodic function of period T . Thus, to prove Theorem 1, one needs only to show that the right side of (55), when computed

from (59), is identical with the right side of (61). From (59), $\psi_1(t)$ is a periodic function of period T . Consequently, $\mathcal{L}^{-1}\{[F(s)s/s + F(s)]\Psi_1(s)\}$ is a zero-mean periodic function of period T . Since $\psi_2(t)$ is a period function with period T , $\mathcal{L}^{-1}\{[F(s)s/s + F(s)]\Psi_2(s)\}$ is also a zero-mean periodic function of period T . Thus, from (60)

$$v_1(t) = \omega_2 - \omega_1 + \bar{v}_1(t) \quad (62)$$

where $\bar{v}_1(t)$ is a zero-mean periodic function with period T . Substituting the $v_1(t)$ in (62) for the integrand $[f * u_1(\rho_1)]$ in (55), we see that the right side of (55) is identical with the right side of (61). This proves Theorem 1.

Equation (59) in Theorem 1 implies that signaling rate of Station 1 locks to that of Station 2. Now we consider the problem of uniqueness (that is, whether (59) is the only steady-state solution). We first prove that, under a simple condition, Station 1 cannot add or delete bits from a customer's data stream.

As described in Section II, the zero crossings of $s_1(t)$ and $s_2(t)$ control the signaling rates of Station 1 and Station 2, respectively. More specifically, Station 2 transmits the m th digit to Station 1 at the mn_0 th zero crossing of $s_2(t)$; and Station 1 transmits the m th digit to Station 2 at the mn_0 th zero crossing of $s_1(t)$ (in practice, $n_0 \geq 1$). Thus, Station 2 transmits a digit to Station 1 every n_0T seconds. We say that the system is in equilibrium if, corresponding to every digit received from Station 2, Station 1 also transmits a digit back to Station 2. More precisely, the system is in equilibrium if we can partition the time axis into N_0T -second time intervals such that Station 1 will transmit a digit back to Station 2 in each of the n_0T -second time intervals.

Theorem 2: The system will reach equilibrium if

$$-\pi < \mathcal{L}^{-1}\left[\frac{s}{s + F(s)} \Psi_1(s)\right] < \pi. \quad (63)$$

Proof: The condition $t \rightarrow \infty$ is implied throughout this proof. Using the final-value theorem, one can show from (57) that when $t \rightarrow \infty$,

$$u_1(t) = \omega_2 - \omega_1 - \sigma_2(t) + \sigma_1(t) \quad (64)$$

where

$$\sigma_2(t) = \mathcal{L}^{-1}\left[\frac{s}{s + F(s)} \Psi_2(s)\right] \quad (65)$$

$$\sigma_1(t) = \mathcal{L}^{-1}\left[\frac{s}{s + F(s)} \Psi_1(s)\right] \quad (66)$$

and \mathcal{L}^{-1} denotes inverse Laplace transform.

Since $\psi_2(t)$ is periodic with period T , $\sigma_2(t)$ is a zero-mean periodic function of period T . Let $\max \sigma_2(t)$ and $\min \sigma_2(t)$ be the maximum and minimum value of $\sigma_2(t)$, respectively. We now determine $\max \sigma_2(t) - \min \sigma_2(t)$. Note that $\psi_2(t)$ can be written as

$$\psi_2(t) = (\omega_2 t + \theta_2) - \sum_i \pi u[t - (t' + iT)] \tag{67}$$

where $u(t)$ is the unit step function defined by

$$\begin{aligned} u(t) &= 0, & t < 0 \\ &= 1, & t > 0. \end{aligned} \tag{68}$$

When $\psi_2(t)$ is applied to a network with transfer function $[s/s + F(s)]$ (hereafter called network A), the output is $\sigma_2(t)$. Clearly, when the first term $\omega_2 t + \theta_2$ in (67) is applied to network A , the output is a continuous time function for $t > 0$. The second term in (67) consists of unit step functions. It can be shown that, when a unit step function $u(t)$ is applied to network A , the output is unity when $t = 0^+$, approaches zero when $t \rightarrow \infty$, and is continuous for $0 < t < \infty$. From these results, it is clear that

$$\max \sigma_2(t) - \min \sigma_2(t) \geq \pi. \tag{69}$$

We have set $e_1 = \pi$. Therefore, $u_1(t)$ is a multiple of π . We are considering the case $\omega_2 - \omega_1 > 0$. As illustrated in Fig. 7, let n be an integer such that

$$n\pi \leq \omega_2 - \omega_1 < n\pi + \pi. \tag{70}$$

It is clear from (70) and (69) that there is a t at which $\omega_2 - \omega_1 - \sigma_2(t)$ equals $n\pi$ or $(n + 1)\pi$ (let this t be denoted by t_1). Note that $\omega_2 - \omega_1 - \sigma_2(t)$ may intersect only the level $n\pi$, or only the level $(n + 1)\pi$, or both the levels. For this proof, we need to consider only the first case. Since $\sigma_2(t)$ is periodic with period T , $\omega_2 - \omega_1 - \sigma_2(t)$ is also periodic

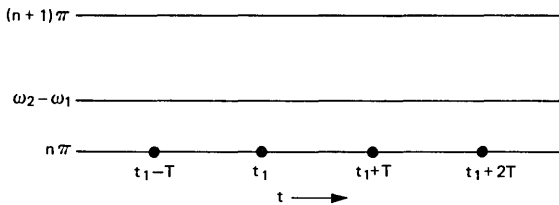


Fig. 7—Illustration for the proof of Theorem 2 (showing the definition of n and the partition of the time axis into successive T -second intervals).

with period T . Thus, $\omega_2 - \omega_1 - \sigma_2(t)$ must intersect the level $n\pi$ also at time instants $t_1 + iT$, $i = \pm 1, \pm 2, \dots$. These intersections are illustrated in Fig. 7.

Now consider the value of $u_1(t)$ at an intersection $t_1 + iT$, $i = 0, \pm 1, \pm 2, \dots$. From (63), we have, at $t = t_1 + iT$

$$n\pi - \pi < u_1(t) < n\pi + \pi. \quad (71)$$

Since $u_1(t)$ must be a multiple of π , (71) implies that

$$u_1(t) = n\pi \quad (72)$$

at $t = t_1 + iT$, $i = 0, \pm 1, \pm 2, \dots$. Since $N_2(t)$ increases by one every T seconds, (72) requires that $N_1(t)$ increase by one in each of the T -second intervals illustrated in Fig. 7. This proves Theorem 2.

Now we consider the case where the filter $F(s)$ is the usual RC filter. We first prove that eq. (63) is satisfied in this case (consequently, the system will reach equilibrium).

When RC filter is used,

$$F(s) = \frac{1}{1 + sCR} \quad (73)$$

$$\frac{s}{s + F(s)} \Psi_1(s) = \Psi_1(s) - \frac{1}{CR(s + r_1)(s + r_2)} \Psi_1(s) \quad (74)$$

where

$$r_1 = \frac{1 + \sqrt{1 - 4CR}}{2CR} \quad (75)$$

$$r_2 = \frac{1 - \sqrt{1 - 4CR}}{2CR}. \quad (76)$$

The system is designed such that

$$1 - 4CR > 0. \quad (77)$$

Therefore r_1 and r_2 are real numbers and

$$r_1 > r_2 > 0. \quad (78)$$

Let $[1/CR(s + r_1)(s + r_2)]$ be denoted by $G(s)$, then

$$g(t) = \mathcal{L}^{-1}[G(s)] = \frac{1}{CR(r_1 - r_2)} [e^{-r_2 t} - e^{-r_1 t}]. \quad (79)$$

From (78) and (79),

$$g(t) > 0, \quad t > 0. \quad (80)$$

From (80), we can write

$$\begin{aligned} & \mathcal{L}^{-1}[\Psi_1(s)G(s)] \\ &= \int_0^t \psi_1(\tau)g(t - \tau) d\tau < \int_0^t \pi g(t - \tau) d\tau < \int_0^\infty \pi g(\tau) d\tau. \end{aligned} \tag{81}$$

Clearly, $\int_0^\infty g(\tau)d\tau = G(0) = 1$. From this and (81), we have

$$0 < \mathcal{L}^{-1}[\Psi_1(s)G(s)] < \pi. \tag{82}$$

From (82) and (74)

$$-\pi < \mathcal{L}^{-1}\left[\frac{s}{s + F(s)}\Psi_1(s)\right] < \pi. \tag{83}$$

Hence, (63) is satisfied and the system will reach equilibrium.

Next, we examine the detailed behavior of the rate-locked loop. Note that there are two basic variables in the rate-locked loop, namely, $u_1(t)$ and $v_1(t)$. Let the $u_1(t)$ and $v_1(t)$ corresponding to the steady-state solution in (59) be denoted by $u_1^*(t)$ and $v_1^*(t)$, respectively. First, we sketch $u_1^*(t)$ and $v_1^*(t)$. From Section II in text, $s_2(t) = \sin[\omega_2 t + \theta_2]$. To simplify our graphs, let us omit θ_2 . Then $N_2(t)$ jumps by 1 at $t = lT, l = 0, 1, 2, \dots$. From this and (59), we see that $u_1^*(t)$ is as sketched in Fig. 8, where l denotes an arbitrary integer. The pulse width y^* in Fig. 8 is such that the mean-value of $u_1^*(t)$ is $\omega_2 - \omega_1$. Therefore,

$$y^* = \frac{T}{\pi} [\omega_2 - \omega_1 - n\pi]. \tag{84}$$

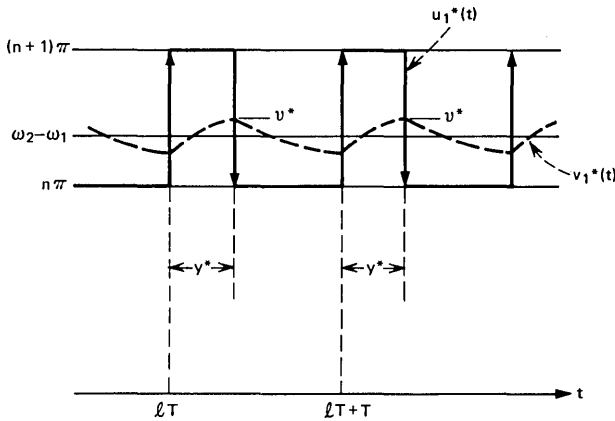


Fig. 8—Sketch of $u_1^*(t)$ and $v_1^*(t)$.

Since $u_1^*(t)$ is periodic with period T , $v_1^*(t)$ is also periodic with period T . As can be seen from $v_1^*(t)$ in Fig. 8, $u_1^*(t)$ charges the capacitor C in the time interval lT to $lT + y^*$, and the capacitor C discharges in the time interval $lT + y^*$ to $lT + T$. Let v^* denote the value of $v_1^*(t)$ at $t = lT + y^*$. Clearly, v^* must have such a value that $v_1^*(t)$ has a mean-value of $\omega_2 - \omega_1$.

In order to show that $u_1^*(t)$ and $v_1^*(t)$ are the only steady-state solution, we begin by assuming different $u_1(t)$ and $v_1(t)$, and demonstrate that they must approach $u_1^*(t)$ and $v_1^*(t)$ as t increases. We have proved that the system must reach equilibrium. From (72), when the system reaches equilibrium, $u_1(t) = n\pi$ at $t = t_1 + iT$, $i = 0, \pm 1, \pm 2, \dots$. [As can be seen from the discussion after (70), $u_1(t)$ may assume the other value $(n + 1)\pi$ at such time instants. However, these two cases are similar and we need to consider only the first case.] Therefore, $u_1(t)$ can assume only one of the two forms in Fig. 9 in each of the time intervals $t_1 + iT$ to $t_1 + iT + T$. The first form is illustrated in the time interval t_1 to $t_1 + T$ in Fig. 9, while the second form is illustrated in the time interval $t_1 + T$ to $t_1 + 2T$. In the first form, the zero crossing of $s_1(t)$ (represented by the downward arrow) takes place prior to the zero crossing of $s_2(t)$ (represented by the upward arrow). The order is reversed in the second form. Note that, if $u_1(t)$ always assumes the first form, one would have $v_1(t) < n\pi$. From this, one can easily show that $u_1(t)$ cannot always assume the first form in the successive T -second intervals. Next, consider the width of the pulse when $u_1(t)$

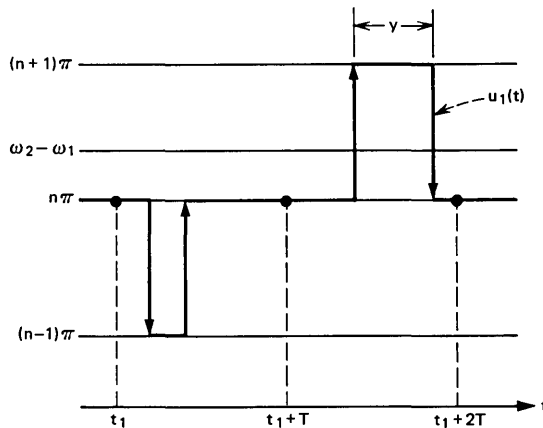


Fig. 9—Illustration of the two forms of $U_1(t)$.

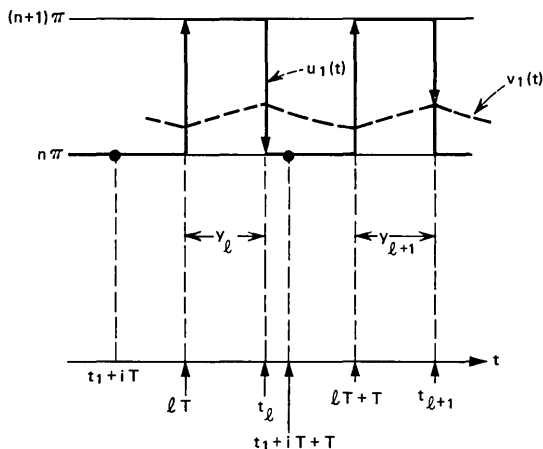


Fig. 10—Sketch of $u_1(t)$ and $v_1(t)$.

assumes the second form. This width, designated by y in Fig. 9, may vary from one T -second interval to the next. If this width were always less than y^* , $v_1(t)$ would be less than ν^* for all t . From this, one can show that this width cannot always be less than y^* . From these results, there must be some T -second intervals in which $u_1(t)$ assumes the second form and the pulse width y is equal to or greater than y^* . We shall select one such time interval (say, the time interval $t_1 + iT \leq t < t_i + iT + T$ illustrated in Fig. 10) and examine $u_1(t)$ and $v_1(t)$ for $t > t_1 + iT$. We need to consider only two cases (refer to Fig. 10):

Case 1: $v_1(t_i) < \nu^*$

Case 2: $v_1(t_i) \leq \nu^*$.

The instantaneous radian frequency of VCO_1 is $\omega_1 + \alpha_1 v_1(t)$, where ω_1 is the free-running radian frequency and $\alpha_1 v_1(t)$ is the correction term. In data communications, ω_1 is very close to the radian frequency ω_2 of the master clock. (For example, it may be specified that the maximum difference between ω_1 and ω_2 be limited to 0.005 percent of ω_2 .) Consequently, only a very small correction term $\alpha_1 v_1(t)$ is needed. For this reason, the time interval between each two consecutive zero crossings of $s_1(t)$ is essentially determined by the term $\omega_1 t$ in $\rho_1(t)$. Therefore, the pulse width y changes only very slightly from one pulse to the next (in other words, in Fig. 10 y_{l+1} is very close to y_l).

For the purpose of illustration, in Fig. 10 $v_1(t)$ is shown to increase

and decrease quite rapidly in each T -second interval. In practice, the filter time constant RC is several orders larger than the time interval T (for example, $RC = 10^{-1}$ seconds, $T \cong 10^{-5}$ seconds). Thus, $v_1(t)$ is essentially a constant in each T -second time interval.

Now consider y_{l+i} and $v_1(t_{l+i})$, $j = 1, 2, 3, \dots$. It can be shown rigorously that if there is an h such that

$$\begin{aligned} y_{l+h} &= y^* \\ v_1(t_{l+h}) &= \nu^* \end{aligned}$$

then $y_{l+i} = y^*$ and $v_1(t_{l+i}) = \nu^*$ for all $j > h$. Therefore, to show that $u_1(t)$ and $v_1(t)$ approach $u_1^*(t)$ and $v_1^*(t)$, we need only show that y_{l+i} and $v_1(t_{l+i})$ approach y^* and ν^* , respectively.

Now consider Case 1; after $t_l + iT$, y_{l+i} and $v_1(t_{l+i})$ approach y^* and ν^* in three stages. Immediately after $t_l + iT$, $v_1(t_{l+i})$ is less than ν^* . Consequently, the time interval between each two consecutive zero crossings of $s_1(t)$ is slightly larger than T , and y_{l+i} increases slowly with j . (Note from Theorem 2 that y_{l+i} must remain less than T .) Since $v_1(t_{l+i})$ is less than ν^* and y_{l+i} remains larger than y^* , $v_1(t_{l+i})$ must increase slowly with j . The second stage starts when $v_1(t_{l+i})$ reaches ν^* . Since the pulse width y_{l+i} is larger than y^* , $v_1(t_{l+i})$ keeps increasing with j . (Note from Theorem 2 that $v_1(t_{l+i})$ cannot exceed $(n+1)\pi$.) This implies that $v_1(t_{l+i})$ will be larger than ν^* . Consequently, y_{l+i} must decrease with j . Clearly, when y_{l+i} decreases, the rate of increase of $v_1(t_{l+i})$ decreases. The third stage starts when y_{l+i} decreases to such a value (still larger than y^*) that $v_1(t_{l+i})$ ceases increasing. Since $v_1(t_{l+i})$ is now larger than ν^* , y_{l+i} must keep decreasing. Clearly, this must also reduce $v_1(t_{l+i})$. Consequently, y_{l+i} and $v_1(t_{l+i})$ approach y^* and ν^* , respectively.

The above discussion is for Case 1. It can easily be extended to Case 2. Thus, when the system reaches equilibrium, the time interval between each two consecutive zero crossings of $s_1(t)$ will be exactly T seconds.

REFERENCES

1. Gardner, F. M., *Phase-lock Techniques*, New York: John Wiley and Sons, Inc., 1966, p. 65.
2. Gersho, A., and Karafin, B. J., "Mutual Synchronization of Geographically Separated Oscillators," *B.S.T.J.*, 45, No. 10 (December 1966), pp. 1689-1704.
3. Karnaug, M., "A Model for the Organic Synchronization of Communication Systems," *B.S.T.J.*, 45, No. 10 (December 1966), pp. 1705-1736.
4. Brilliant, M. B., "The Determination of Frequency in Systems of Mutually Synchronized Oscillators," *B.S.T.J.*, Vol. 45, No. 10 (December 1966), pp. 1737-1748.

5. Williard, M. W., "Analysis of a System of Mutually Synchronized Oscillators," IEEE Trans. Commun. Tech., *COM-18*, No. 5 (October 1970), pp. 467-483.
6. Williard, M. W., and Dean, H. R., "Dynamic Behavior of a System of Mutually Synchronized Oscillators," IEEE Trans. Commun. Tech., *COM-19*, No. 4 (August 1971), pp. 373-395.
7. Saito, T., "Dynamic Characteristics of Mutually Synchronized Systems," Elec. Commun. Japan, *51-A*, No. 4, 1968, pp. 19-27.
8. West, N., "Synchronous Digital Switching in Highly Interconnected Communication Networks," Proc. Inst. Elec. Eng., *114*, No. 10 (October 1967*e*), pp. 1378-1384.
9. Chang, R. W., Mazo, J. E., and Salz, J., "Stability Analysis of a Digital Rate-Lock Loop," unpublished work.
10. Chang, R. W., unpublished work.

Contributors to This Issue

J. E. DAVID BATSON, JR., Sc.B., 1961, Brown University; M.S., 1971, Northeastern University; U. S. Navy, 1961–1967; Western Electric Company, 1967—. With the Navy, Mr. Batson was a carrier-based jet fighter pilot and flight instructor. At Western Electric, he worked on digital multiplex test development and many phases of D2 Channel Bank development and manufacturing. He is presently Project Coordinator for the Millimeter Waveguide project.

STEPHEN D. BLOOM, RCA Institutes, Inc., 1960; Bell Laboratories, 1960—. Since joining Bell Laboratories, Mr. Bloom has been engaged in the design and development of electronic power systems. Recently he has concentrated in the design of pulse-width-modulated-type converters and integrated-circuit-type series regulators.

ROBERT W. CHANG, B.S.E.E., 1955, National Taiwan University; M.S.E.E., 1960, North Carolina State University; Ph.D., 1965, Purdue University; Bendix Corporation, 1960–1963; Bell Laboratories, 1965—. Mr. Chang has worked on a variety of problems in data transmission and communication system theory. Member, Phi Kappa Phi, Eta Kappa Nu, Sigma Xi, IEEE.

ANTHONY J. CIRILLO, B.E.E., 1962, Manhattan College; S.M., 1963, Massachusetts Institute of Technology; Bell Laboratories, 1962—. Mr. Cirillo's initial work was on high-speed solid-state coders for pulse-code-modulation transmission systems. He also worked on timing and framing circuitry for the D2 Channel Bank. At present, he is involved in circuit and software development for the Voiceband Interface Frame which is part of the No. 4 ESS Toll Switching project. Member, Eta Kappa Nu, Tau Beta Pi.

CARL L. DAMMANN, B.S.E.E., 1961, University of Maryland; M.E.E., 1963, New York University; Bell Laboratories 1961—. Mr. Dammann has been engaged in development of PCM Terminal equipment. He is now Supervisor of the Coding Group. Member, IEEE.

M. EISENBERG, B.S. (E.E.), 1964, M.S. (E.E.), 1964, Ph.D. (E.E.), 1967, Massachusetts Institute of Technology; Bell Laboratories, 1967—. Mr. Eisenberg has worked on several problems in the fields of queuing theory, network management, and network design. He has also done research on certain aspects of communications theory. Member, Operations Research Society of America, Association for Computing Machinery, Sigma Xi, Tau Beta Pi, Eta Kappa Nu.

D. GLOGE, Dipl. Ing., 1961, Dr. Ing., 1964, Technical University of Braunschweig, Germany; Bell Laboratories, 1965—. Mr. Gloge's work has included the design and field testing of various optical transmission media and the application of ultra-fast measuring techniques to optical component studies. He is presently engaged in transmission research related to optical fiber communication systems.

KENNETH A. GLUCKOW, Electronics Eng. Cert., 1957, RCA Institutes; Newark College of Engineering, 1957-62; Bell Laboratories 1957—. Mr. Gluckow initially worked in the Bell System Repair Specification organization, preparing repair requirements and repair studies for transmission products. Since 1962, he has been involved with reliability studies for transmission products. He has recently assumed responsibility for coordinating physical design efforts for the Voiceband Interface Frame portion of the No. 4 ESS project.

JAMES W. GORMAN, B.S.M.E., 1958, University of Maine; Western Electric Company, 1958—. Mr. Gorman has worked in Machine Design, Tantalum and Mica Capacitor Engineering, D2 Channel Bank Development and Product Engineering, and is now Lead Engineer for D1B Channel Bank Network and Channel Unit Product Engineering.

HANSJUERGEN H. HENNING, B.E.E., 1955, Polytechnic Institute of Brooklyn; M.E.E., 1961, New York University; Bell Laboratories, 1955—. Mr. Henning has been engaged in the design of PCM transmission systems, including the D1 and D2 Channel Banks, and experimental high-speed transmission systems. He also was engaged in circuit design for the *Telstar*[®] experimental satellite. Since 1970, he has been a member of the Ocean Systems Technology Laboratory where he was responsible for a group concerned with the design of hardware

for digital processing and display. He is presently engaged in the design of digital transmission systems for underwater applications. Member, Sigma Xi.

ROBERT H. KRAMBECK, B.E., 1965, City College of New York; M.S.E.E., 1966, and Ph.D., 1969, Carnegie-Mellon University; Bell Laboratories, 1968—. Mr. Krambeck has been engaged in the analysis and development of new types of memory elements. Member, IEEE.

COLLIER LEE MADDOX, B.S.E.E., 1959, Illinois Institute of Technology; M.E.E., 1961, New York University; Bell Laboratories 1959—. Since joining Bell Laboratories, Mr. Maddox has been concerned with the development of terminals for PCM systems. Member IEEE, Tau Beta Pi, Eta Kappa Nu.

DIETRICH MARCUSE, Diplom Vorpruefung, 1952, Dipl. Phys., 1954, Berlin Free University; D.E.E., 1962, Technische Hochschule, Karlsruhe, Germany; Siemens and Halske (Germany), 1954-57; Bell Laboratories, 1957—. At Siemens and Halske, Mr. Marcuse was engaged in transmission research, studying coaxial cable and circular waveguide transmission. At Bell Laboratories, he has been engaged in studies of circular electric waveguides and work on gaseous masers. He spent one year (1966-1967) on leave of absence from Bell Laboratories at the University of Utah. He is presently working on the transmission aspect of a light communications system. Mr. Marcuse is the author of two books. Member, IEEE, Optical Society of America.

L. D. McDANIEL, B.S.E.E., 1964, Southern Methodist University; M.S.E.E., 1967, Polytechnic Institute of Brooklyn; Mobil Field Research Laboratory, 1960-64; Bell Laboratories 1965—. Before joining Bell Laboratories, Mr. McDaniel was engaged in designing and fabricating electronic equipment for usage in sonic and nuclear well logging tools. Since joining Bell Laboratories, he has been concerned with design and development of analog-to-digital and digital-to-analog converters for PCM Channel Banks. Member, Sigma Tau, Eta Kappa Nu.

JOHN W. PAN, B.S. (E.E.), 1955, University of Cincinnati; Sc.D. (E.E.), 1962, Massachusetts Institute of Technology; Bell Laboratories 1955—. Mr. Pan received a fellowship for study at MIT, 1958–1962. His early work at Bell Laboratories was concerned with the process of pulse code modulation and digital transmission. Since 1964, he has been responsible for a group that is engaged in design and analysis of digital transmission systems. Member, Eta Kappa Nu, IEEE, Sigma Xi, Tau Beta Pi.

KENNETH A. PICKAR, B.S. (Physics), 1961, Queens College of C.U.N.Y.; M.S., 1963, and Ph.D. (Physics), 1966, University of Pennsylvania; Bell Laboratories 1966—. Mr. Pickar has worked in radiation damage effects and device technology using ion implantation. More recently he has been studying electron beam lithography. During the academic year 1972–73, he is on leave from Bell Laboratories to teach at the Technion–Israel Institute of Technology.

SUKETU R. SHAH, B.Sc., 1961, Gujrat University, Ahmedabad, India; M.E. (E.E.), 1971, Stevens Institute of Technology; Bell Laboratories, 1967—. Mr. Shah is a member of the Radio Research Laboratory. Since joining Bell Laboratories, he has worked on thin-film solid-state devices and thin-film circuits at microwave and millimeter-wave frequencies. He is presently engaged in work on microwave integrated circuits for use in short-haul terrestrial communication systems.

DAVID J. VANSLOOTEN, B.S.E.E., 1949, Newark College of Engineering; Bell Laboratories, 1937—. Mr. VanSlooten's early work was in drafting, technical writing, and trial installation areas. His numerous physical design projects include various military projects, the *Telstar*[®] satellite, Sub-cable Test Set, and D2. His present assignment is the Voiceband Interface Frame of the No. 4 ESS. Member, Tau Beta Pi.

G. FOLKE SWANSON, I.M.E., 1932, Pratt Institute of Technology; attended Brooklyn Polytechnic Institute for E.E.; Bell Laboratories, 1937—. Mr. Swanson has been involved in the physical design of military radars, Nike-Zeus ABM, Unicom, and No. 1 ESS ADF Data

Switching System. Since 1966, he has been a member of the Physical Design Department of the Electronic Power Systems Laboratory. He has worked on physical design of carrier, microwave, military submarine cable power plant, and key telephone power supplies.

DENNIS K. THOVSON, B.S., 1960, M.S., 1961, Iowa State University; Bell Laboratories, 1961—. Mr. Thovson worked on the final development of the command circuitry for the *Telstar*[®] project, and did exploratory development work on logic circuitry for an experimental 224-Mb/s PCM system. More recently he has been engaged in the design of the digital circuitry and channel units for the D2 Channel Bank. Member, IEEE, Phi Kappa Phi.

R. H. WALDEN, B.E.S., 1962, M.E.E., 1963, and Ph.D., 1966, New York University; Bell Laboratories, 1966—. Mr. Walden's initial activities in the Semiconductor Device Laboratory were concerned with switching properties of VO_2 , followed by work on the conduction properties of Al_2O_3 films. For the last two years he has been involved in a study of the properties of charge transfer devices. Member, IEEE.

B. S. T. J. BRIEFS

Extension of Multidimensional Polynomial Algebra to Domain Circuits with Multiple Propagation Velocities

By S. V. AHAMED

(Manuscript received May 26, 1972)

I. INTRODUCTION

Multidimensional Polynomial Algebra¹ is a new technique for the analysis of circuits in which a finite and distinct time interval is necessary for the propagation of binary bits of data from one location to the next. These conditions exist specifically in magnetic domain circuits.^{2,3} The algebra expands the basic concepts of the coding theorists,^{4,5} and includes the representation of both time and space in the algebraic representation of data streams. A set of algebraic transformations¹ has been developed to correspond to the subfunctions in the circuit, and the overall function is thus modeled by a series of algebraic transformations. Such an analysis predicts the location and the value of all the binary positions at any prechosen instant of time, thus leading to the algebraic verification of the operation of a proposed circuit. In the foregoing technique for analysis, one velocity of propagation was assumed. However, the circuit designer may depend on more than one velocity for the successful operation of the circuit, and the technique suggested here accounts for different propagation velocities. Further, the analysis proposed determines the relationships between such velocities.

In discrete circuits, multiple velocities are generally derived from clock sources driven at different rates. The movement from one location to the adjoining location is finite, but the duration for the movement is derived from different clocks.

II. REPRESENTATION OF MULTIPLE CLOCKS

In Ref. 1 it was proposed that the number of clock cycles between a prechosen origin of time and a given instant of time be represented as the exponent of X ; X being defined as the carrier of the time dimension. In multiple clock systems it is proposed that X be subscripted to denote the various clocks available in the system. Hence, if one adopts the

notation of Ref. 1 to denote the binary values of data bit positions (i.e., a_0, a_1, \dots, a_{n-1} for the first, second, \dots , n th, data bit positions), and to denote their spatial locations (i.e., $Y_k^{l_0}, Y_k^{l_1}, Y_k^{l_2}, \dots, Y_k^{l_{n-1}}$ for the first, second, \dots , n th locations in the k th element of a circuit), then a stream of data, n bits long after j_0 clock cycles at the first clock (denoted by the subscript 0 and X), may be represented as

$$u = X_0^{j_0} \sum_{i=0}^{i=n-1} a_i Y_k^{l_i}, \quad (1)$$

where l_i indicates the location number of the i th data position.

If this stream traverses for m_0 clock cycles in the forward direction, and is propagated at the first clock rate, then the final condition is

$$u = X_0^{j_0+m_0} \sum_{i=0}^{i=n-1} a_i Y_k^{l_i+m_0}. \quad (2)$$

Now if the data stream is propagated at a second clock rate (denoted by the subscript 1 for X) for m_1 clock cycles, the binary data is then represented as

$$u = X_0^{j_0+m_0} X_1^{m_1} \sum_{i=0}^{i=n-1} a_i Y_k^{l_i+m_0+m_1}. \quad (3a)$$

In general, the stream (1) after $m_1, m_2, \dots, m_j, \dots, m_\delta$ clock cycles at $X_0, X_1, X_2, \dots, X_j, \dots, X_\delta$ clocks, having been propagated from an element k to an element t , with their intersection located at Y_k^z and Y_t^z , may be represented as

$$u = X_0^{j_0} \prod_{i=1}^{i=\delta} X_i^{m_i} \sum_{i=0}^{i=n-1} a_i Y_t^{l_i + \sum_{j=0}^{j=i} m_j}^{-z}. \quad (3b)$$

Similar expressions for streams after looping, duplicating, logical gating, etc., may also be written (Ref. 1), and it is thus possible to model a series of functions in the circuit with multiple clocks by a series of algebraic equations, as was indicated for circuits with a single clock.

III. RELATION BETWEEN CLOCK RATES

The functional constraints on the circuit demand that bubble positions or streams be physically present at certain predefined locations and at preselected intervals of time. For example, it may be necessary

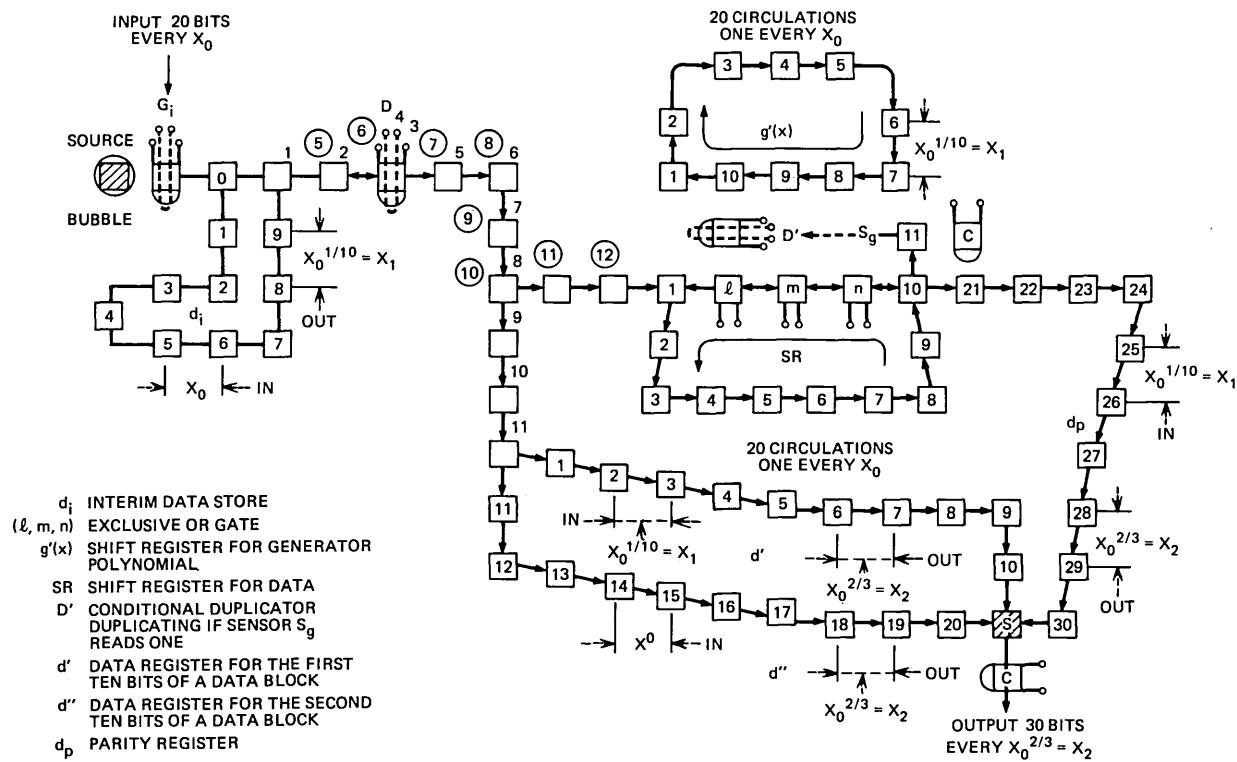


Fig. 1—Magnetic domain encoder for the (30, 20) code with conductor propagation.

for a binary stream to have completely circulated a loop (p periods) once and advanced one additional period* before the arrival of the next data bit.[†] If the incoming data is assumed to arrive every X_0 clock cycle, and the data in the p -period loop is being propagated one location every X_1 clock cycle, then one cycle at the rate X_0 should correspond to $(p + 1)$ cycles at the rate X_1 . This leads to the conclusion that

$$X_1^{p+1} = X_0, \quad (4)$$

and this equation should be construed to imply that the propagating clock at X_1 runs $(p + 1)$ times faster than the clock at X_0 .

IV. APPLICATION OF ALGEBRA TO THE DESIGN OF A (30, 20) DOMAIN ENCODER

The application of the algebra to the design of a (30, 20) domain encoder⁶ shown in Fig. 1 has yielded all the design parameters, the relation between different clocks, and the instants of synchronization of the various clocks to obtain a satisfactory operation of the encoder. It is foreseen that such an algebraic analysis of discrete circuits will help engineers to check the operation of the devices before they are constructed and reduce the debugging time once they have been made.

REFERENCES

1. Ahamed, S. V., "Multidimensional Polynomial Algebra for Bubble Circuits," B.S.T.J., 51, No. 7 (September 1972), pp. 1535-1558.
2. Bobeck, A. H., Fischer, R. F., and Perneski, A. J., "A New Approach to Memory and Logic Cylindrical Domain Wall Devices," Proc. FJCC (1959), pp. 489-498.
3. Bonyhard, P. I., Danylchuck, I., Kish, D. E., and Smith, J. L., "Application of Bubble Devices," IEEE Trans. Magnetics, MAG 6, No. 3 (September 1970), pp. 447-451.
4. Peterson, W. W., *Error Correcting Codes*, Cambridge, Mass.: MIT Press (1961).
5. Berlekamp, E. R., *Algebraic Coding Theory*, New York: McGraw-Hill Book Company, 1968.
6. Ahamed, S. V., "The Design and Embodiments of Magnetic Domain Encoders and Decoders for Block Cyclic Codes," B.S.T.J., 51, No. 2 (February 1972), pp. 461-485.

* A period is defined as that unit of physical distance by which a binary position is propagated in one clock cycle at any rate X_0 , or X_1 , or X_2 , etc.

[†] Such a requirement is placed in the magnetic domain encoders and decoders presented in Ref. 6.

Blooming Suppression in Charge Coupled Area Imaging Devices

By C. H. SÉQUIN

(Manuscript received June 23, 1972)

An intense spot of light projected onto the photo-sensitive surface of an imaging device can cause this device to saturate locally. Excess carriers generated by the light source can diffuse into the neighboring area which may also be driven into saturation. In the display the light source will then appear as a white area that can be considerably larger than its image in the true geometrical proportions. This effect, known as blooming, is present in most TV camera tubes, and demands special care by the operator to avoid bright objects in the scene being imaged. For the *Picturephone*[®] camera, which often has to operate in less than ideal conditions, the design of a camera with limited blooming is thus more than desirable.

The camera tube presently used in the *Picturephone* station set has a silicon diode array target scanned by an electron beam. Blooming is produced by the diffusion of carriers in the bulk silicon, leading to a circular spreading of the saturated area. In solid-state imaging devices, blooming can take on even more objectionable forms. The complicated potential distribution at the silicon surface can cause excess carriers to move along a preferred axis, generating quite irregular blooming patterns in the display.

In a recently demonstrated 128×106 -element charge coupled array¹ the excess charge spills preferentially in the vertical direction. In this n-channel frame transfer² device the isolation between adjacent CCD channels is achieved by a p-type channel stopping diffusion which keeps the potential at the Si-SiO₂ interface close to zero. Due to the negative flatband voltage of this particular MIS-system, the potential underneath a grounded transfer electrode, separating two adjacent potential wells in the vertical direction, is a few volts positive, and this barrier is thus distinctly lower than the one produced by the channel stopping diffusion. No negative voltage can be applied to the transfer electrodes because they are connected to diffused crossunders or to protection diodes. In this device, blooming appears in a very objectionable form. Bright light spots bloom out into a vertical line that quickly extends

the length of the picture and then, with increasing light intensity, starts to widen.

The basic idea behind blooming protection consists in providing an overflow drain for excess carriers. This drain can consist of a reverse biased diffused junction of the same polarity as the output diode. In area imaging devices with frame transfer organization these overflow channels can be placed between the vertical transfer channels and interconnected and accessed at the top of the device (see Fig. 1).

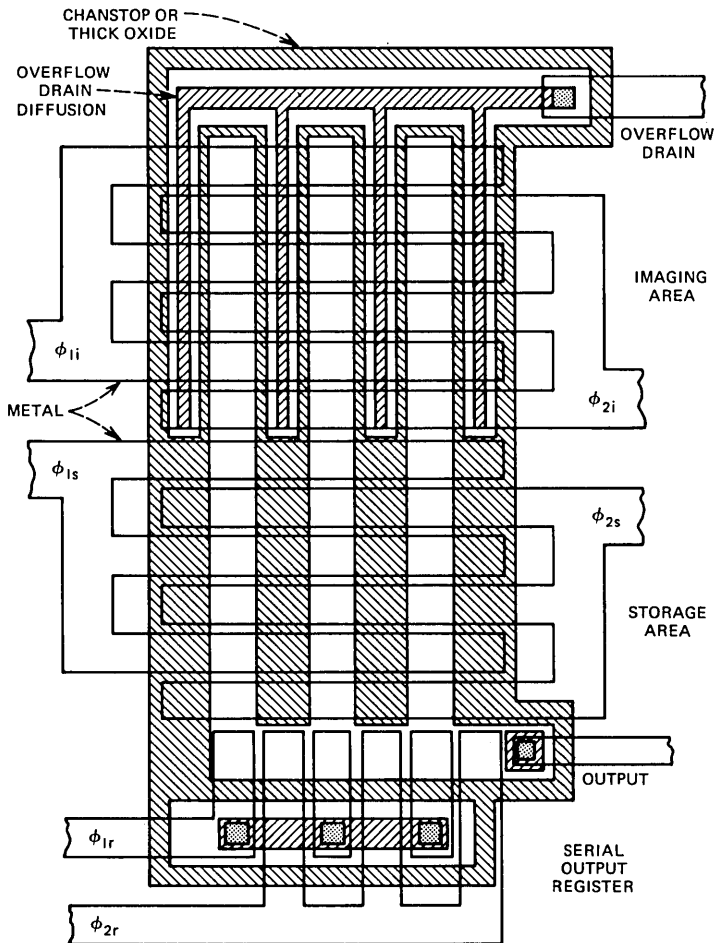


Fig. 1—Geometrical arrangement of the overflow channels in a charge transfer imaging device of frame transfer organization.

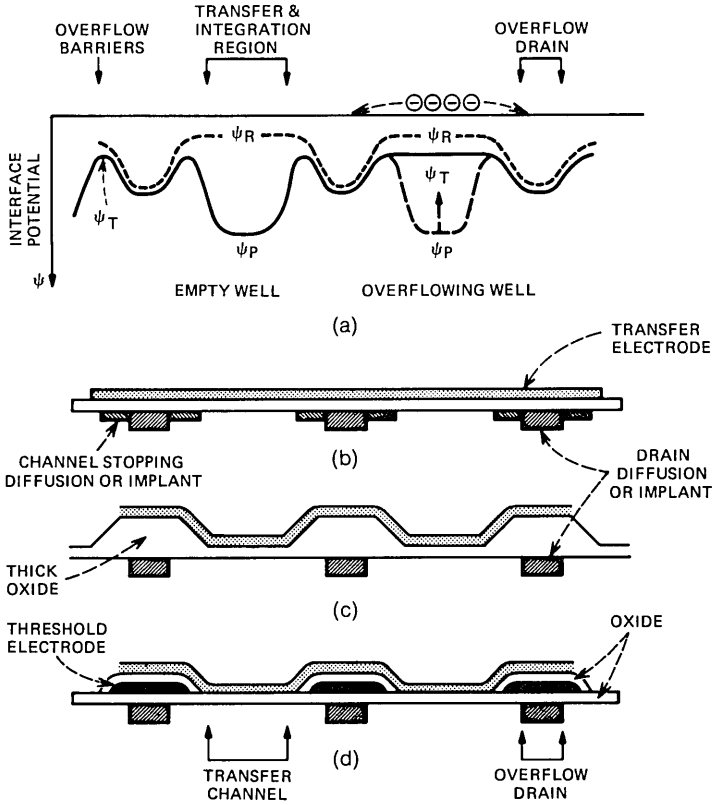


Fig. 2—Cross section through imaging area of a frame transfer image sensor. (a) Desired potential profile at Si-SiO₂ interface, underneath isolating electrodes (---) and underneath integrating electrodes (—). (b, c, d) Realization using a channel stopping diffusion or implant (b), using a thick-thin oxide structure (c), or a special threshold electrode (d).

Excess carriers can reach the overflow drain by passing over a potential barrier (Fig. 2a). This overflow threshold is established either by a light channel stopping diffusion or implant (Fig. 2b), by a thick oxide region (Fig. 2c), or by a special threshold electrode (Fig. 2d).

The centers of charge collection are the deep potential wells underneath the electrodes biased at V_P , generating an initial interface potential ψ_P (Fig. 2a). The resolution elements are isolated in the vertical direction by electrodes kept at V_R , producing interface potential ψ_R . This potential barrier has to be higher than the overflow threshold ψ_T . The potential well is filled when its interface potential

has reached ψ_T and additional carriers will then escape laterally into the overflow drain. Ideally, an overflowing cell should thus not affect its neighbors.

The overflow drain will directly collect a certain percentage of all generated carriers, thereby reducing the light sensitivity of the device. To a certain extent, the sensitivity can even be modulated by changing the potential in the overflow channel.³ On the other hand, the horizontal resolution might be somewhat improved since the overflow drain reduces the overlap of the sensitivity functions of adjacent resolution elements.

It is not anticipated that the introduction of a large area of p-n junction will lead to a problem with increased dark current. Excess current generated in the junction area should not influence the operation of the device but be swept away through the overflow drain.

A preliminary experiment performed on the described 128×106 -element area CCD showed that the migration of excess carriers can be controlled. One of the protection diodes in the imaging section has been burnt out so that one of the sets of electrodes can be biased negatively. This drives the Si-SiO₂ interface underneath into accumulation and generates a potential barrier equal in height to the barrier produced by the channel stopping diffusion. As expected, when a negative potential of a few volts was applied to the unprotected electrode set, blooming in the display changed from a vertical line into a circle of about 200 μm radius, corresponding to the diffusion length of the minority carriers.

The described protection scheme could also be applied to an electron-beam-scanned target if in biasing the overflow drain special care is taken not to increase the target-to-ground capacitance, which reduces the signal-to-noise ratio. In a charge transfer device the substrate is normally grounded and the additional capacitance is of no concern.

REFERENCES

1. Séquin, C. H., Sealer, D. A., Bertram, W. J., Tompsett, M. F., Buckley, R. R., Shankoff, T. A., and McNamara, W. J., "A Charge Coupled Area Image Sensor and Frame Store," unpublished work.
2. Tompsett, M. F., Amelio, G. F., Bertram, W. J., Buckley, R. R., McNamara, W. J., Mikkelsen, J. C., and Sealer, D. A., "Charge Coupled Imaging Devices: Experimental Results," *IEEE Trans. Electron Devices*, *ED-18* (1971), pp. 992-996.
3. Stupp, E. H., Kostelec, J., Steneck, W., Labak, J., Kidder, M., and Crowell, M. H., "Silicon Diode Array Target with Gating Capabilities," presented at the International Electron Device Meeting, Washington, D. C., October 1971.

Data Transmission Performance in the Presence of Carrier Phase Jitter and Gaussian Noise

By E. Y. HO and D. A. SPAULDING

(Manuscript received June 23, 1972)

I. INTRODUCTION

In the operation of data transmission systems over voice-grade telephone channels, phase jitter¹ is a commonly observed transmission impairment. It appears in the form of low-index angle modulation of the received data signals. It is believed that phase jitter is a very important parameter in determining system performance. Therefore, many complicated methods^{2,3} have been developed to recover the jittered carrier. However, recent field measurements¹ show that the phase jitter in Bell System carrier systems has improved significantly over the past few years. As a result of this improvement, the following question naturally arises: How much phase jitter recovery is required for two-level and four-level systems?

In this B.S.T.J. Brief, we analyze the system performance degradation caused by phase jitter. The results suggest that for two-level systems, jitter need not be recovered and that for four-level systems, a coarse jitter recovery system would provide acceptable performance.

II. GENERAL CONSIDERATIONS

A simplified block diagram of a general VSB-AM data system is depicted in Fig. 1. A random message sequence $\{a_n\}$ is used to modulate an identically shaped pulse train, which is then transmitted over a voice-grade telephone line. The received random pulse train is corrupted by additive Gaussian noise and intersymbol interference;* the latter is slowly time-varying and is caused by the phase jitter in the carrier system which causes crosstalk between the in-phase and quadrature channels. The received pulse train is processed and sampled to provide the estimates of the transmitted message sequence $\{\hat{a}_n\}$. A useful measure of the performance of such a data system is the error probability, $P_r\{\hat{a}_n \neq a_n\}$.

III. THE PROBABILITY BOUND

If the peak-to-peak phase jitter is small, then the pulse train presented

* Channel amplitude and delay distortion are assumed to be removed by an equalizer and other impairments such as nonlinear distortion and impulse noise are neglected.

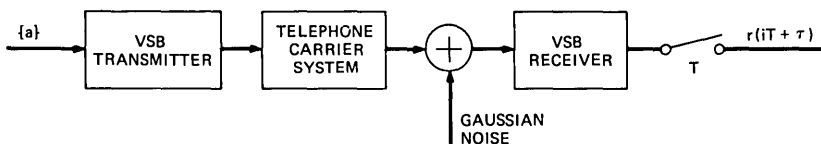


Fig. 1—Block diagram of a general VSB-AM system.

to the sampler at the receiver in the presence of additive Gaussian noise may be approximated by

$$r(t) = \sum_K a_K f(t - KT) + \phi(t) \sum_K a_K \bar{f}(t - KT) + n(t), \quad (1)$$

where the a_K are the transmitted symbols, $f(t)$ and $\bar{f}(t)$ are the basic in-phase and quadrature pulse response of the system, $1/T$ is the baud, $n(t)$ is a zero-mean Gaussian noise process and $\phi(t)$ is the phase jitter which is the sum of slowly varying sinusoids, $\phi(t) = \sum_i P_i \cos(\omega_i t + \theta_i)$. We assume the a_K are identically and independently distributed random variables with probabilities

$$P\{a_K = 2j + 1\} = \frac{1}{2M}, \quad j = -M, \dots, -1, 0, 1, \dots, (M - 1). \quad (2)$$

The probability density of $n(t)$ is

$$P(n(t)) = (2\pi\sigma_n^2)^{-1/2} \exp(-n^2(t)/2\sigma_n^2). \quad (3)$$

The m th transmitted symbol is determined by sampling $r(t)$ at $t = t_0 + mT$, i.e.,

$$\begin{aligned} r(t_0 + mT) &= a_m f(t_0) + \sum_{K \neq m} a_{K-m} f(t_0 - KT) \\ &\quad + \phi(t_0 + mT) \sum_K a_{K-m} \bar{f}(t_0 - KT) + n(t_0 + mT). \end{aligned} \quad (4)$$

The error probability is defined to be

$$\begin{aligned} P_e &= P_r\{\hat{a}_m \neq a_m\} \\ &= \frac{2M - 1}{M} P_r\left\{ \sum_{K \neq m} a_{K-m} f(t_0 - KT) + \phi(t_0 + mT) \right. \\ &\quad \left. \cdot \sum_K a_{K-m} \bar{f}(t_0 - KT) + n(t_0 + mT) > f(t_0) \right\}. \end{aligned} \quad (5)$$

Assume the sampling instant is perfect and that there is no channel amplitude or delay distortion, i.e., $f(t_0 - KT) \equiv 0$ for all $K \neq 0$, and $f(t_0) = 1$. This is a reasonable assumption if an adaptive equalizer is

incorporated in the system. Thus, (5) can be rewritten as

$$P_e = \frac{2M - 1}{M} P_r \{n(t_0 + mT) + \phi(t_0 + mT) \cdot \sum_K a_{K-m} \bar{f}(t_0 - KT) > f(t_0)\}. \quad (6)$$

Applying the Chernoff bound⁴

$$\Pr \{z > y\} \leq \exp \{-\lambda y\} \langle [\exp (\lambda z)] \rangle, \quad \text{all } \lambda > 0 \quad (7a)$$

and the following inequality

$$\langle \exp \{a_K \cdot x\} \rangle \leq \exp (x^2 \cdot \sigma_a^2 / 2) = \exp \{x^2 \cdot (2M - 1)(2M + 1) / 6\}, \quad (7b)$$

to (6), we obtain an upper bound on the conditional error probability

$$P_e |_{\phi(t_0+mT)} \leq \frac{2M - 1}{M} \cdot \exp \left\{ - \frac{f^2(t_0)}{2(\sigma_n^2 + \phi^2(t_0 + mT)) \cdot \frac{(2M + 1)(2M - 1)}{3} \cdot \sum_K \bar{f}^2(t_0 - KT)} \right\}. \quad (8)$$

From (8) it can be seen that an upper bound of the error probability is

$$P_e \leq \frac{2M - 1}{M} \cdot \exp \left\{ - \frac{f^2(t_0)}{2(\sigma_n^2 + \phi_P^2) \cdot \frac{(2M + 1)(2M - 1)}{3} \cdot \sum_K \bar{f}^2(t_0 - KT)} \right\}, \quad (9)$$

where ϕ_P is the maximum phase.

IV. EXAMPLE

A multilevel single-sideband AM system is used as a vehicle to determine the performance degradation caused by the phase jitter. The system signal-to-Gaussian noise ratio is assumed to be 24 dB and is defined as

$$S/N = \frac{\langle a_K^2 \rangle \cdot f^2(t_0)}{\sigma_n^2}. \quad (10)$$

The power of the signal is normalized to unity, i.e., $f^2(t_0) = 1$ and $\sum_K \bar{f}^2(t_0 - KT) = 1$. Figures 2 and 3 are plots of the probabilities of error bound versus peak-to-peak phase jitter for two- and four-level

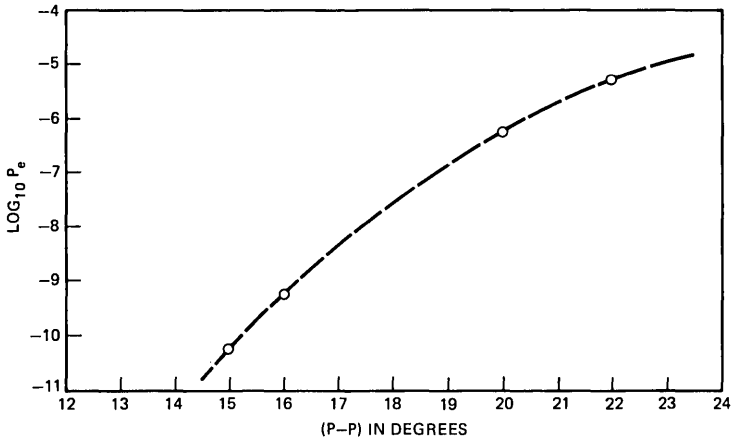


Fig. 2—The probability of error bound versus peak-to-peak phase jitter $f^2(t_0) \cdot \sigma_a^2 / \sigma^2 = 24$ dB, two-level system.

signaling. It can be seen from these curves that if the peak phase jitter for two- and four-level systems is limited to less than 20 and 6 degrees respectively, a probability of error less than 10^{-6} is achieved.

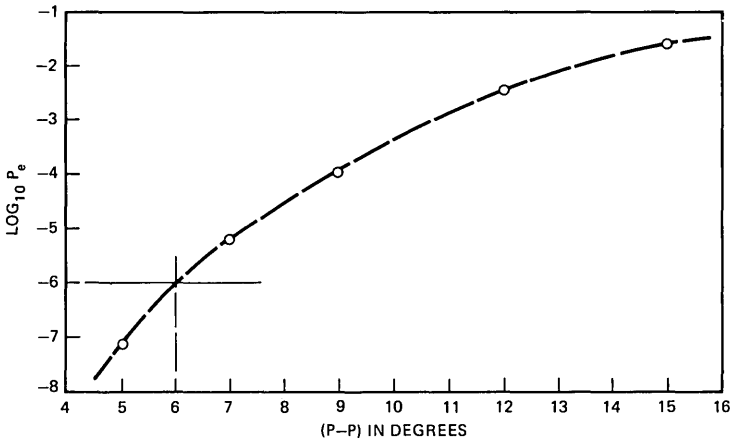


Fig. 3—The probability of error bound versus peak-to-peak phase jitter $f^2(t_0) \cdot \sigma_a^2 / \sigma^2 = 24$ dB, four-level system.

V. CONCLUSIONS

An upper bound of the error probability of a VSB-AM data system operated in the presence of additive Gaussian noise and phase jitter is presented in this correspondence. By restricting our attentions to these two parameters alone, it has been possible to calculate curves which can be used to estimate the accuracy required of a phase jitter recovery system.

REFERENCES

1. Duffy, F. P., and Thatcher, Jr., T. W., "Analog Transmission Performance on the Switched Telecommunications Network," *B.S.T.J.*, 50, No. 4 (April, 1971), pp. 1311-1347.
2. Bennett, W. R., and Davey, J. R., *Data Transmission*, New York: McGraw-Hill Book Company, 1965.
3. Lucky, R. W., Salz, J., and Weldon, Jr., E. J., "Principles of Data Communication," New York: McGraw-Hill Book Company, 1968.
4. Saltzberg, B. R., "Intersymbol Interference Error Bounds with Application to Ideal Bandlimited Signaling," *IEEE Trans. Inform. Theory*, *IT-14*, No. 4 (July 1, 1968), pp. 563-568.

CONTENTS

(Continued from front cover)

B.S.T.J. Briefs:

- | | | |
|---|-------------------------------------|-------------|
| Extension of Multidimensional Polynomial Algebra to
Domain Circuits With Multiple Propagation Velocities | S. V. Ahamed | 1919 |
| Blooming Suppression in Charge Coupled Area Imaging
Devices | C. H. Séquin | 1923 |
| Data Transmission Performance in the Presence of
Carrier Phase Jitter and Gaussian Noise | E. Y. Ho and D. A. Spaulding | 1927 |



Bell System