

The accordion imager, a new solid-state image sensor

A. J. P. Theuwissen and C. H. L. Weijtens

No image sensor can compete with the human eye. Even if a charge distribution could be produced in a solid-state image sensor with the same accuracy as an image is formed on the retina, the information could not be transferred in the same way. Each photosensitive element in the eye has its own channel for transferring the image information to the brain. A simulation of this system would require far too many connections. The information originating from all the individual image elements in solid-state image sensors is ultimately transferred through a single channel. The idea for the transfer of the image information in the 'accordion imager' described below was first put forward in the early eighties at Philips Research Laboratories.

Introduction

A solid-state image sensor seems to be an attractive alternative to the conventional television camera tube. In a camera tube an electron beam scans the charge distribution produced on a photoconductor by incident light. In most solid-state sensors the charge distribution is not scanned; the charge is transferred directly.

Major applications of such image sensors are to be found in monitoring and surveillance equipment, electronic cameras and even in toys. Because of the interesting applications in the consumer market it is important to keep the price of the imager low. Since the price is directly related to the size of the sensor, the sensor should be as small as possible without loss of resolution.

A solid-state image sensor has a number of advantages compared with a camera tube: it can be made in IC technology, its weight is low, it requires little power and it is small and robust. Image-lag or 'comet' effects can be avoided and the image area is not damaged if the light beam is too bright. The image quality, however, is not that of a camera tube. For a comparable

image quality, a solid-state sensor should have a certain minimum number of picture elements (pixels). But this must not require any change in the critical dimensions (and hence a new technology) or make the chip too expensive. In the current state of the technology the chip will only compare with other image sensors if its area is less than 40 mm².

A solid-state image sensor is simply a silicon chip. It consists of two parts, an image section and a storage section. In the image section incident light generates electrons, which are collected in potential wells at defined positions on the surface (the pixels). This results in a distribution of charge packets that corresponds to an image. The size of a charge packet corresponds to the quantity of light that arrives at a pixel. After charge has been accumulated during a specific period (the 'integration period') this charge distribution is transferred in its entirety^[1] to the storage section, which is shielded from light. The transfer has to be very fast to ensure that the charge packets in transit are not significantly affected by incident light. Nor, of

^[1] The charge distribution corresponding to an image is transferred in its entirety to the storage section situated below the image section. Sensors of this type are therefore called frame-transfer (FT) devices.

course, must there be any interaction between the charge packets. The potential wells in the silicon are produced by applying voltages to a number of electrodes positioned on the surface, perpendicular to linear channels of n-type silicon.

From the storage section the charge packets are transferred row by row (line by line) to the output register. Further processing for television pictures can then take place in the usual way.

The resolution of the sensor is determined by the number of pixels per unit area. The dimensions of a pixel are fixed by the width of one n-channel and by the number and dimensions of the electrodes that are necessary to produce the potential wells. In a four-phase image sensor a pixel comprises four electrodes (see fig. 1).

Fig. 2 shows a cross-section through part of a four-phase image sensor, which is essentially a shift register of the CCD type (charge-coupled device) [2]. A substrate of p-type silicon contains narrow channels of n-type material. Above this there is an insulating layer of silicon oxide, and on top of that layer there are the electrodes, consisting of polycrystalline silicon. Light

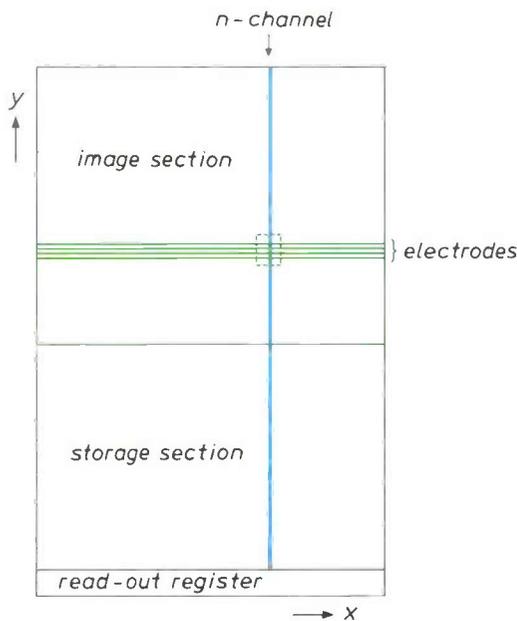


Fig. 1. Diagram representing a solid-state image sensor. From top to bottom, the sensor consists of three sections: an image section, a storage section and a read-out register. In the image and storage sections the n-channels are vertical (y-direction). For clarity only one n-channel is indicated. The electrodes run horizontally across the channels (x-direction). The figure shows only four electrodes for the image section. The dotted rectangle indicates a picture element (a pixel) in a four-phase image sensor. Each pixel contains a charge packet at the end of an integration period. After an integration period has been completed, all the charge packets are transferred to the storage section simultaneously. The storage section contains as many pixels as the image section. From the storage section the charge packets are moved towards the read-out register a row at a time. Each row of charge packets contains information that corresponds to a line in a television picture.

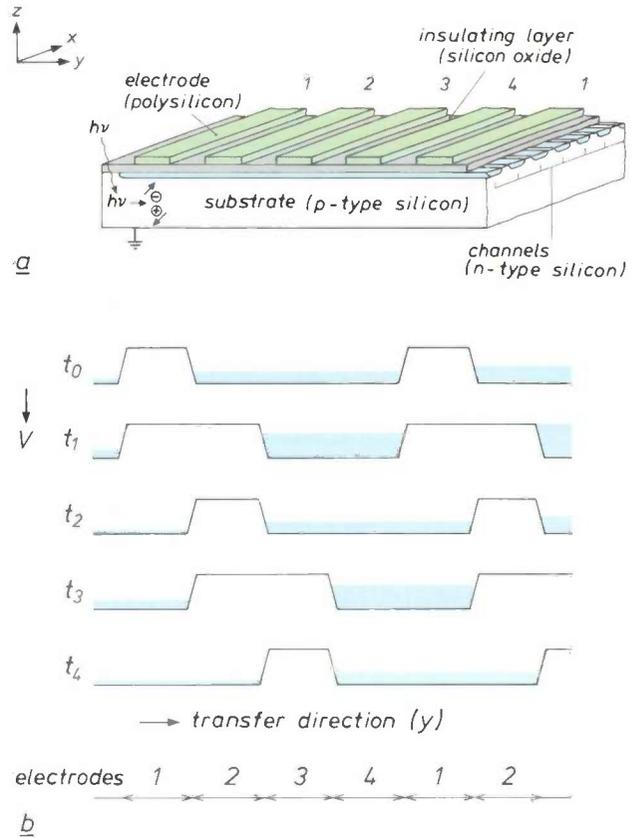


Fig. 2. a) Diagram of part of a four-phase solid-state image sensor. A p-type silicon substrate contains narrow channels of n-type silicon (in the y-direction). The surface of the substrate is covered with a layer of silicon oxide, and on top of that layer are the polysilicon electrodes (in the x-direction). The silicon oxide is necessary to prevent conduction between the substrate and the electrodes and between the electrodes themselves. Light passing through the transparent electrodes and the silicon oxide releases electrons and holes in the silicon ($h\nu \rightarrow +$ and $-$). The holes are conducted to earth and the electrons are collected in potential wells in the n-channels which are produced by applying voltages to the electrodes. b) Potential (V) in the silicon at consecutive times. The first signal represents the situation during an integration period. The electrons are collected in the potential well beneath the electrodes 2, 3 and 4. From the moment t_1 the changes in the voltages on the electrodes cause a peristaltic transfer of the charge packets. The blue colouring indicates the contents of a potential well. All the charge packets are transferred simultaneously.

passing through the electrodes and the insulating layer generates electrons in the silicon. In addition to electrons positive charge carriers are also generated, and these are conducted to earth. The potential wells in which the electrons are collected are created by applying a positive voltage to electrodes 2, 3 and 4 and a negative voltage to electrodes 1. The first signal (marked t_0) represents this situation, which remains unchanged for the full integration period of 20 ms. Then the charge distribution is transferred to the storage section. In the image sensor shown in fig. 2 this transfer is effected by generating a peristaltic 'potential movement' controlled by a clock signal. This propels the charge packets towards the storage section.

The peristaltic 'potential movement' can in principle be produced with only three electrodes. We use four electrodes to obtain the usual interlacing for television pictures. During successive integration periods the voltage patterns on the electrodes are shifted symmetrically with respect to each other by an integer number of electrodes. In the field duration following the transfer in fig. 2, electrode 3 acts as a barrier electrode and the voltage on electrodes 1, 2 and 4 produces the potential wells.

For the *collection* of the charge packets during the integration period two electrodes per pixel are sufficient. The peristaltic charge *transfer* can only take place if a pixel consists of more than two electrodes. In fig. 2 there are four. Because of these extra electrodes a row of pixels occupies a relatively large area on the chip. This has detrimental consequences for the resolution in the vertical direction (y), which is determined by the number of pixels per unit length. For this reason we wanted to make a sensor with smaller pixels. The obvious way of doing this is to make the electrodes narrower. However, this would require an entirely new technology for producing the sensors. Another way of making the pixels smaller is to reduce the number of electrodes in each row. This does not greatly affect the production process, but it does of course change the transfer mechanism. The way in which this problem has been solved in our new 'accordion imager' is treated in this article, which describes the operation of the sensor, its design and its characteristics.

The 'accordion' operation

The new 'accordion imager' [3] is a solid-state image sensor with only two electrodes per pixel, which serve alternately as integration and barrier electrodes in successive integration periods. The storage section also consists of elements with only two electrodes. Charge transfer takes place as follows. After an integration period the potential wells and barriers are doubled in width one at a time. This process starts at the interface between the image and storage sections of the sensor. The charge distribution is stretched like an accordion. When the lower edge of the storage section is reached the potential wells and barriers are reduced to a width of one electrode, one at a time. The charge distribution is now 'squeezed together' like an accordion. The result of this stretching and squeezing is that the entire charge distribution is moved from the image section to the storage section. The transfer from the storage section to the output register can be made in the same way.

It is clear that in this method of transfer the voltages on the electrodes will change in a more complex way than in the four-phase sensor. The benefits, however,

are considerable: with the same technology as used for producing a four-phase sensor an equally large sensor can be made that has twice as many pixels, or a smaller sensor with the same number of pixels. The accordion imager has the same number of pixels as the four-phase sensor but occupies a smaller chip area.

The detailed action for the charge transfer in the accordion imager is illustrated in fig. 3 and fig. 4. At time t_0 (fig. 3) the sensor is in one of the two states in which charge is accumulated during a period of 20 ms, the integration period. At time t_1 the first charge packet (a) is stretched. At time t_2 it is pushed further towards the output. The other charge packets (b, c, \dots) stay in position. At time t_3 the second charge packet (b) starts to move, and so on. This change from a static two-phase system to a dynamic four-phase system continues until every charge packet has been spread over two electrodes. The separation between each two charge packets has a width of two elec-

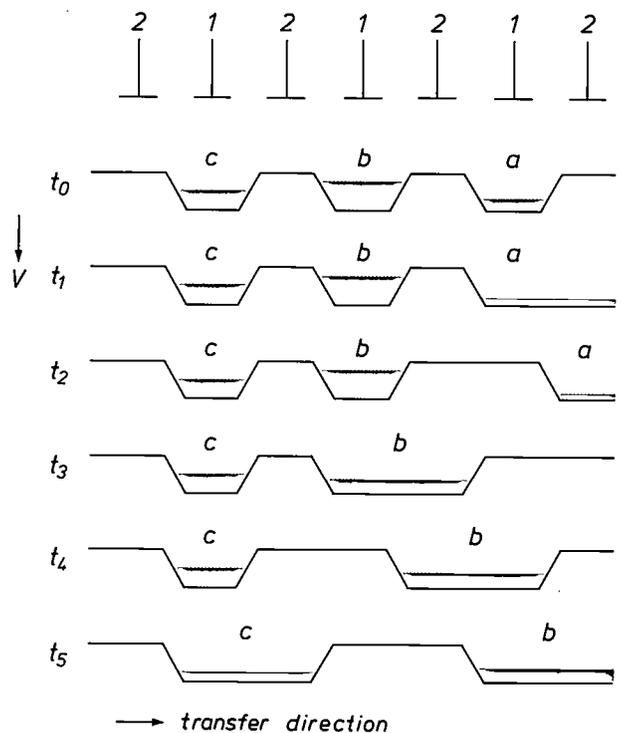


Fig. 3. Principle of the 'accordion' mechanism. The starting point is a two-phase system, in which charge is collected beneath the electrodes 1 during the integration period and the electrodes 2 form a barrier between the different charge packets. A four-phase system is built from this, a step at a time. The state at time t_0 occurs at the end of an integration period, t_1 to t_5 indicate the successive stages in 'stretching the accordion'. It can clearly be seen that charge packet a is transferred first, followed by charge packet b , and so on to the output of the sensor.

[2] F. L. J. Sangster and K. Teer, Bucket-brigade electronics — new possibilities for delay, time-axis conversion, and scanning, IEEE J. SC-4, 131-136, 1969.

[3] The ideas on which this work is based were put forward by A. J. J. Boudewijns, now with the Consumer Electronics Division, Philips NPB, formerly with Philips Research Laboratories, and M. G. Collet and L. J. M. Esser of these laboratories.

trodes. At this moment the first packet has arrived at the first part of the storage section and the accordion starts to close up. The potential wells and barriers are gradually made one electrode wide.

The entire process of stretching and squeezing the 'accordion' is illustrated in *fig. 4* for a sensor in which the image and storage sections each have eight electrodes. Note that at the end of the cycle the sensor is in the appropriate state for collecting charge beneath the electrodes that were barrier electrodes in the previous integration period; this provides the interlacing.

In the accordion imager the charge transfer takes place over two times 588 electrodes (image and storage). This takes 0.5 ms, which is sufficiently short compared with the integration period of 20 ms. The application of the voltages to the electrodes is obviously more complicated than in the four-phase sensor. The voltage change is no longer the same for electrodes with the same number, as it was in the four-phase sensor. In the following section we shall show how the required voltage pattern on the electrodes is produced.

The accordion imager

Principle

So far we have considered the consequences of the changing voltages on the electrodes for the charge transfer. Let us now see how these voltage changes are produced.

The image sensor is controlled by two shift registers, one for the image section and one for the storage section. These shift registers consist of cells, each with a clock input. The output of each cell is connected to an electrode of the sensor and to the input of the next cell. At the input of the first cell the signal can be either *IM* (for the image section of the sensor) or *ST* (for the storage section of the sensor). The clock signals ϕ_1 , which controls the clock inputs of the odd cells in each shift register, and ϕ_2 , which controls the clock inputs of the even cells, make this input signal shift step by step through the register, with signal inversion after each cell. This process produces the required voltage pattern on the electrodes. The signal at the input of the first cell of the shift register (*IM* or *ST*) causes the stretching and squeezing of the accordion. The way in which this is done is illustrated in *fig. 5* for a small part of the sensor (8 electrodes of the image section and 8 electrodes of the storage section with the associated shift registers).

At the moment when the clock associated with a cell generates a pulse, the output of the cell takes on a value opposite to the value at the input. The input signals *IM* and *ST*, as shown in *fig. 5a*, produce the

potential profile shown in *fig. 5b* on the electrodes A_1, B_1, \dots . As long as the input signal *IM* is changing, the potential wells and barriers are two electrodes wide and the charge packets are transferred. The variation of the potential wells and barriers with time can be seen in *fig. 5* from the values of the outputs A_1, B_1, \dots at consecutive times. We see that the required potential pattern is produced at the numbered times, and that it corresponds to the state of the sensor as illustrated in *fig. 4*.

As soon as the input signal *IM* (or *ST*) becomes constant, the voltages on the electrodes of the image section (or storage section) also become constant, after a short delay. Gradually the accordion is squeezed shut, and the process continues until the potential distribution consists of wells and barriers, each with a width of only one electrode. This potential distribution remains unchanged as long as the signal *IM* remains constant. By keeping the input signal *IM* high in one integration period and low in the other, the electrodes function alternately as barrier and integration electrodes (see *fig. 4* and *fig. 5*).

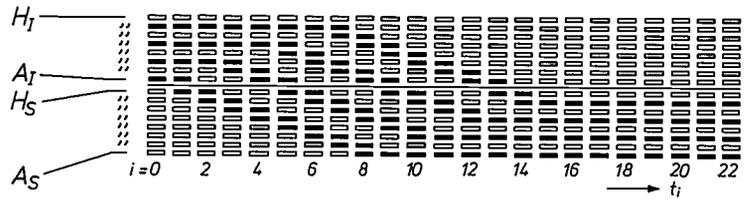
Correct operation of the sensor requires a number of synchronization signals — in addition to the clock signals ϕ_1 and ϕ_2 and the control signals *IM* and *ST*. These synchronization signals indicate the start and finish of the stretching of the two 'accordions' (the image and storage sections). They can be determined by counting the number of clock pulses, for example. Some of them can also be derived from the state of the shift register [4]. The first synchronization signal indicates the start of the stretching of the charge distribution in the image section and coincides with the end of the integration period. This signal must of course be supplied from outside the chip [5]. The end of the stretching process is reached when the potential of the final electrode of the image section, electrode H_1 in the example given in figures 4 and 5, changes for the first time. The arrows in *fig. 5* indicate the states from which the synchronization signals are derived. At this time, t_8 , the accordion of the storage section starts to close up. There must then be no further change in the signal *ST*.

When the charge distribution in the storage section has completely closed up, signal *IM* is kept constant. The potential distribution in the image section also closes up. The time at which this must start to happen is reached when there is no further change in the potentials of the last two electrodes in the storage section (t_{15}).

[4] A. J. P. Theuwissen, C. H. L. Weijtens and J. N. G. Cox, The accordion imager: more than just a CCD-sensor, Proc. Electronic Imaging 85, Boston 1985, pp. 87-90.

[5] J. N. G. Cox contributed to the design and construction of the special electronic units required here.

Fig. 4. A complete cycle of the stretching and squeezing of the 'accordion' of an image sensor consisting of 16 electrodes. A_1 to H_1 indicate the electrodes of the image section, A_s to H_s the electrodes of the storage section. A blue rectangle corresponds to a row of charge packets beneath it. A red rectangle corresponds to a barrier electrode, and a green rectangle to an electrode that has no charge beneath it yet, but is at a positive potential. The first column gives the potential distribution at t_0 , i.e. at the time when a complete integration period has just finished. Beneath the electrodes A_1 , C_1 , E_1 and G_1 there are rows of charge packets (a , b , c and d). At t_1 the first row of charge packets starts to move, at t_2 the second, and so on. At time t_8 the first row of charge packets has arrived at the lower edge of the storage section, and the 'accordion' of the storage section can then be squeezed shut. At t_9 the first row of charge packets is again one electrode wide. We note that from t_9 the electrodes of the image section that no longer take part in the transfer



of the current charge distribution take up a new potential distribution, which will collect charge during the next integration period. After the complete charge distribution has arrived in the storage section, a clock signal must be sent to provide the potential distribution in the image section with the required one-electrode-wide wells and barriers. From t_{21} charge can be collected beneath the electrodes (B_1 , D_1 , F_1 and H_1). After completion of this integration period A_1 , C_1 , E_1 and G_1 are again used to collect the electrons, so that the complete cycle of stretching and squeezing starts again.

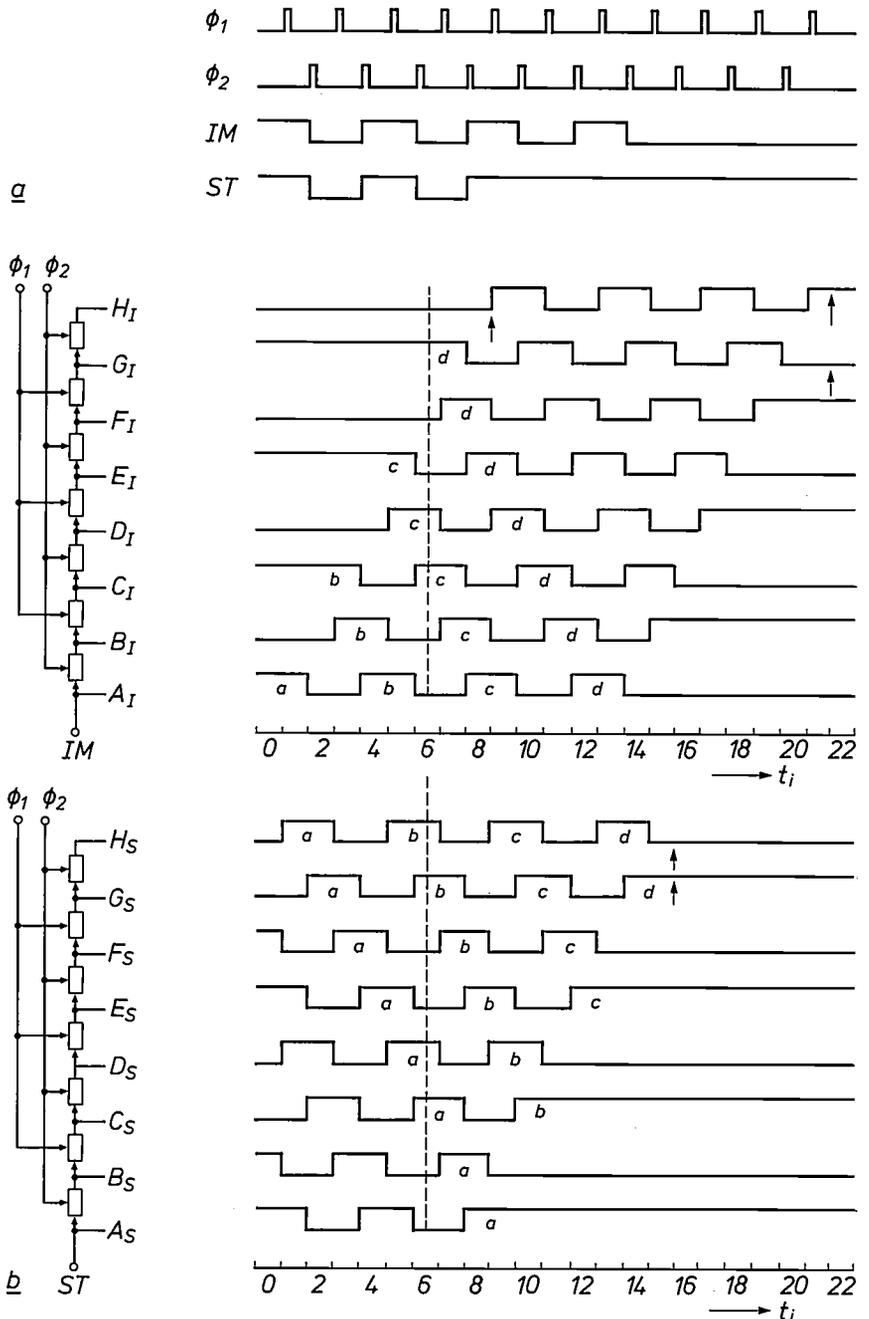


Fig. 5. a) Control signals for the transfer of the charge distribution from image section to storage section in the 16-electrode image sensor shown in fig. 4. IM and ST are the input signals of the 7-cell shift registers (shown as rectangles). These registers produce the potential profile shown in b on the electrodes A_s to H_1 . The clock signals ϕ_1 and ϕ_2 are necessary for 'clocking' the input signal through the shift register. This potential profile causes the row of charge packets (whose position is indicated by a , b , c , ...) to be transferred from the image section to the storage section. The numbers below the time axis and the designations of the electrodes correspond to those given in fig. 4. At time t_6 , for example (see dashed line), the potential pattern on the electrodes corresponds to the pattern in fig. 4 for t_6 . The arrows indicate the points in the potential pattern from which the synchronization signals can be determined (see text). After the transfer the functions (separation and integration) of the electrodes in the image section are exchanged.

The final synchronization signal, which indicates that the accordion of the image section is squeezed fully shut, is derived from the states of the two final electrodes of the image section. For the sensor to function correctly it is not necessary to 'know' when the

Practical design

The two shift registers for the control of the electrodes of the image and storage sections of the accordion imager have to be fabricated on the same chip as the CCD. The same design rules apply to the produc-

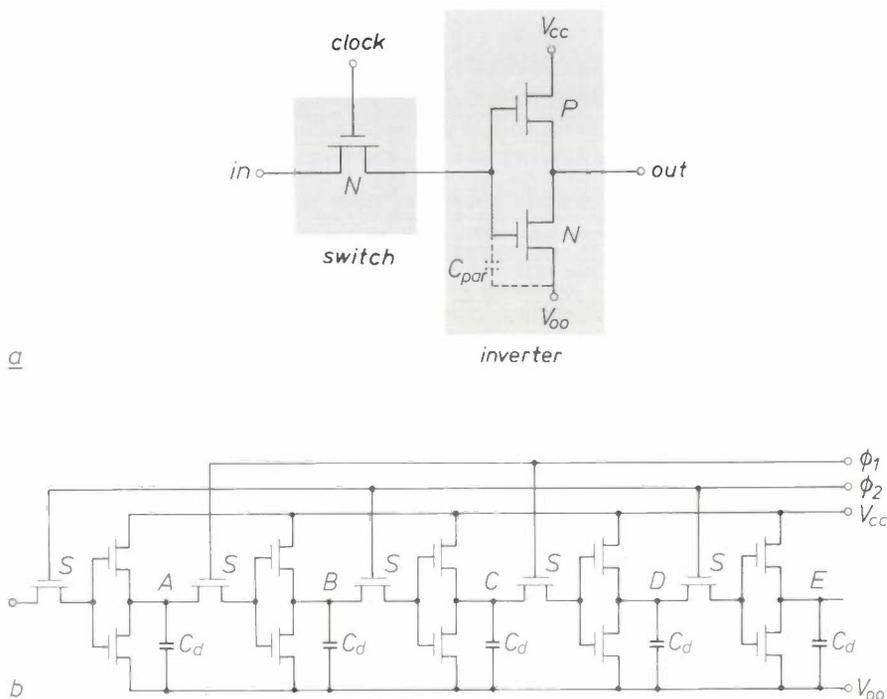


Fig. 6. a) Schematic diagram of a cell of the digital shift register that effects the transfer of the charge distribution. A cell consists of a combination of MOS transistors. The first NMOS transistor feeds the input signal along the register 'in phase with' the clock, which therefore acts as a switch. The output of this switch is connected to the gates of an NMOS transistor and a PMOS transistor. The source of the NMOS transistor is connected to a low voltage (V_{00}) and the source of the PMOS transistor is connected to a high (V_{cc}) voltage. This means that the drain electrodes, which are interconnected, are at a high voltage when the output of the switch — the value of the gate voltage — is low, or at a low voltage when the output of the switch is high. The final part of the cell acts as an inverter. The dashed capacitor C_{par} represents the parasitic capacitance of the connections. b) A row of these cells forms the shift register. The switches S are connected alternately to the clock signals ϕ_1 and ϕ_2 . The output of each cell (A_1, B_1, \dots) is connected to the input of the next one and to a CCD electrode (shown here as a capacitance C_d). The parasitic capacitances are not shown.

'accordion' is squeezed shut again. However, to keep the dissipation low, the clocks are stopped at that time (t_{22}).

During the next integration period the charge distribution from the storage section is read out line by line (not shown in fig. 5). This is also done with the 'accordion' action. We should note that after this transfer the 'accordion' of the storage section is squeezed shut. This implies that at time t_0 the potential distribution on the electrodes of the storage section consists of wells and barriers two electrodes wide. At the end of this integration period a new charge distribution has been produced in the image section. At this moment the clocks are started again, and then the signals IM and ST are triggered by external pulses.

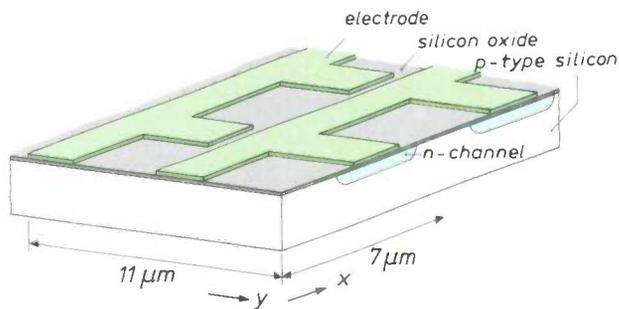


Fig. 7. Diagram showing the arrangement of the CCD electrodes and n-channels in the image section of the accordion imager. The electrodes are shaped in such a way that they do not cover the entire surface but leave some of the p-type material and a part of the n-channels accessible to incident light. This increases the blue sensitivity of the sensor. The polysilicon electrodes absorb light strongly in the blue wavelength range.

tion of the CCD and the shift registers. This means that the shift registers have to be fabricated in the 3.5- μm technology normally used for the CCDs. Moreover, the cells of the shift register should require as little energy as possible and occupy the smallest possible chip area. The simplest arrangement that delivers the required signal and also meets the above requirements is a combination of three MOS transistors per cell: two n-type and one p-type (see *fig. 6a*). The output signal from such a cell becomes the inverse of the input as soon as the clock gives a pulse. While the clock remains inactive (although the input may change) the output remains constant because of the stray capacitance of the inverter gate.

A shift register consists of 587 cells (the number of CCD electrodes less one). The inputs of the switches of these cells are connected alternately to one of the two clock signals. These two clock signals must not overlap in time. Otherwise, the shift register would behave like a row of unsynchronized inverters while both clocks were active.

Each CCD electrode is controlled by a cell of the shift register. The shift register obviously ought to be placed next to the electrodes on the chip, but this is not so easy.

The minimum dimensions of a photosensitive element of the image sensor are 7 μm in the horizontal direction and 11 μm in the vertical direction^[6] (see *fig. 7*). The electrodes, which are made of polysilicon, are shaped in such a way that some of the p-type material and a part of the n-channels are accessible to incident light. This improves the sensitivity of the sensor to blue light (polysilicon is almost opaque to blue light). In the storage section of the sensor it is not necessary to leave a part of the silicon surface uncovered by the electrodes. The area of the part of the channel beneath the electrode must be the same as in the image section if it is to have the same charge storage capacity. This means that the dimension in the y -direction of a pixel in the storage section can be smaller than in the image section. This vertical dimension is 9 μm . In *fig. 8* it can indeed be seen that the storage section of the sensor is smaller than the image section.

It is obvious that a shift-register cell cannot be located next to every horizontal electrode: a cell consists of a number of p-n junctions and therefore has a vertical dimension several times the minimum width of 3.5 μm . For this reason a few of these cells (four) are placed side by side and on opposite sides of the CCD. The distribution of the available chip area among the different components can be seen in *Table I*.

The synchronization signals that mark the start and finish of the stretching and closing of the 'accordions'

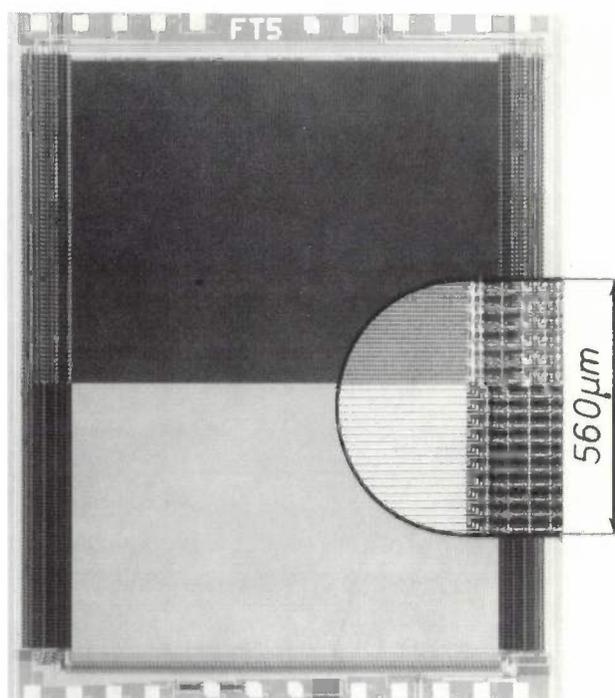


Fig. 8. The accordion imager. The storage section (the light region) is smaller than the image section (the dark region). The storage section is shielded from light by an aluminium layer. The n-channels are vertical, the CCD electrodes horizontal. On the left and right are the shift registers, which generate the voltages on the CCD electrodes. *Insert:* enlarged view of the transition from image section to storage section.

Table I. Allocation of the available chip space to the different sections of the sensor.

Section	Relative area
Image section	35.8%
Storage section	29.3%
CCD output register	2.9%
Digital shift registers	16.8%
Peripheral connections	15.2%

(see previous subsection) can be determined by means of a small on-chip electronic circuit consisting of about 20 gates. This circuit determines the logical state of the final two electrodes of the image and storage sections. The clock signals and input signals are started or stopped depending on this state.

[6] The minimum horizontal dimension follows from the 3.5- μm technology employed. The transfer in the horizontal direction takes place in a three-phase shift register. The length of a cell in this register, which really consists of three shift registers one above the other, is 21 μm . This is why the minimum dimension of an electrode of the shift register, and hence the minimum dimension of a channel, is 7 μm . The pixel matrix must have the standard TV format. This means that the ratio of the horizontal to the vertical dimensions of the image section must be 4:3. Since there are 604 pixels per line and 294 lines, a pixel must have a vertical dimension of 11 μm .

Characteristics of the sensor

To conclude, we shall recapitulate the characteristics of the accordion imager, and compare them with those of its predecessor, the four-phase image sensor. *Table II* gives the main characteristics of both types of sensor.

The area occupied by the accordion imager is only 56% of that of the four-phase device. This makes the accordion imager the smallest solid-state image sensor described in the literature with such a large number of pixels. This reduction of chip area has not entailed any sacrifice of resolution. One of the advantages of the smaller number of electrodes is that it reduces the dissipated power. A lower dissipated power per pixel results in a lower dark current, and therefore better

image quality. Dark-current fluctuations, which are local, are one of the main limitations in the use of a solid-state image sensor. These fluctuations can cause bright spots to appear at various points in the picture.

Another advantage is that the electronic control circuits not included on the chip are simpler and more compact. Less polysilicon is required on the chip, and this improves the photosensitivity of the sensor, particularly for blue light. The most important aspect of the accordion imager, however, is that a number of characteristics connected with the chip size and the number of pixels (especially the cost of the chip) are better than for other image sensors, yet it has not been necessary to design a completely new method of production or to develop an entirely new production process.

Table II. Some characteristics of the accordion imager compared with those of the four-phase sensor.

	Four-phase sensor	Accordion imager
Type ^[1]	frame transfer	frame transfer
Chip area ($y \times x$)	9.41×7.01 = 66 mm^2	7.01×5.45 = 38.2 mm^2
Number of pixels ($y \times x$)	588×604	588×604
Pixel dimensions ($y \times x$)	$15.6 \times 10 \text{ } \mu\text{m}^2$	$11 \times 7 \text{ } \mu\text{m}^2$
Read-out rate	11.5 MHz	11.5 MHz
Layout rules	3.5 μm	3.5 μm
Number of electrodes per pixel	4	2

Summary. Solid-state image sensors are silicon chips in which a charge distribution is generated by incident light. The use of such devices in video cameras for the consumer market will depend greatly on the price of the sensor. As with all chips, this is closely dependent on the chip area. With a new read-out mechanism, in which the charge distribution is stretched out and squeezed shut like an accordion, it is possible to build an image sensor with two electrodes instead of the usual four for each line of the frame. The new sensor can be produced with the same technology as used for the four-phase sensor. The electronic control circuits can be fabricated on the chip with the sensor in a single process.

Philips Technical Review 50 years ago

DEMONSTRATION MODEL ILLUSTRATING SUPERHETERODYNE RECEPTION

MARCH 1936

At the World Exhibition held in Brussels last year, a demonstration model was exhibited which illustrated in the most elementary manner the method of operation of modern radio receivers. It is of course generally known that in radio transmission a carrier wave of high frequency is employed and the much lower audio-frequencies are transmitted as a modulation of the amplitude of this carrier wave. The function of the rectifying stage in the receiver is to separate these low audio-frequencies from the high frequency of the carrier wave. Perhaps less well known is the sequence of operations which actually take place in a superheterodyne receiver in which an "oscillator", a "converter valve" and an "intermediate frequency" are used. The main object of the demonstration model described in the present article is therefore to illustrate the principles underlying superheterodyne reception.

The Superheterodyne Principle

In superheterodyne receivers, not only is the incoming high-frequency carrier wave appropriately dealt with, but provision is also made for the simultaneous generation of an additional oscillation by means of an oscillator.

To generate the auxiliary oscillation and to modulate it on the incoming oscillation a special converter valve, the octode, is used in Philips superheterodyne receivers. This valve may be regarded as a triode and a pentode, connected in series, in which the triode serves for the generation of the auxiliary oscillation, while the pentode

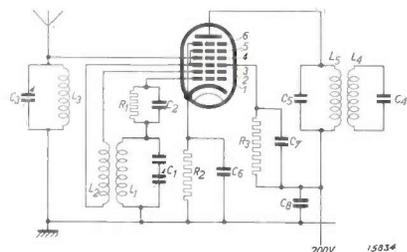


Fig. 1. Circuit diagram of the octode. The circuit L_1, C_1 is tuned to the auxiliary oscillation. L_2, C_2 to the incoming signals. L_3, C_3 and L_4, C_4 to the intermediate frequency.

- 1 = Control grid
- 2 = Screen-grid
- 3 = Auxiliary anode
- 4 = Control grid
- 5 = Screen-grid
- 6 = Interceptor grid

further handles the alternating current generated in the triode.

Mechanical representation

Electrical oscillations are usually represented diagrammatically by means of a sinusoidal or similar type of wavy line. It appeared therefore that these oscillations, which are usually drawn with chalk on a blackboard, could be usefully reproduced mechanically for general exhibition by means of a sand figure on a slowly-moving belt. This would enable a practical demonstration of the principles of reception and the properties and uses of the carrier wave. In the demonstration model which was evolved for this purpose, the utilisation

of the actual carrier wave is also demonstrated by means of a number of cathode ray tubes, which produce a visible trace of the electrical oscillations on their fluorescent screens.

In *fig. 2* the apparatus is shown which was exhibited on the Philips stand at Brussels. The graphing of the various oscillations occurring in radio receivers in the form of sand figures is performed on a moving belt which moves in a horizontal direction from right to left. The belt is only just visible in this picture, but is more clearly shown in *fig. 3*. The feed tubes for the sand and the pendulums fitted with funnels which swing to and fro back to front and thus produce the sand figures, can be picked out in *fig. 2*.

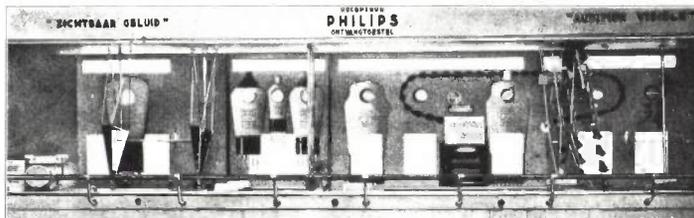


Fig. 2. Front view of model demonstrating superheterodyne reception, as shown at the Brussels 1935 World Exhibition.

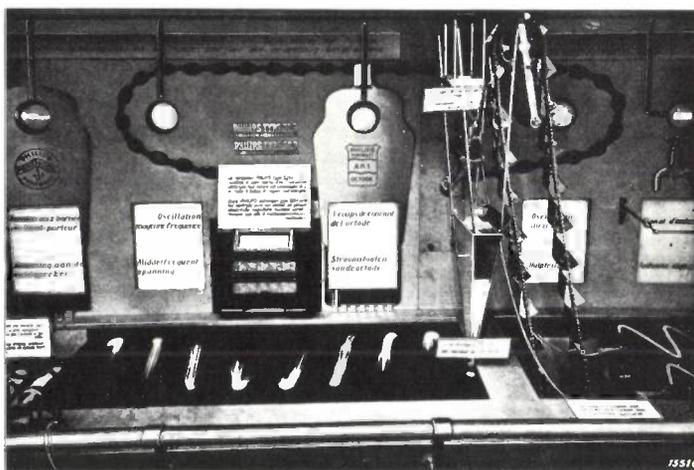


Fig. 3. Sand figures demonstrating the anode current of the octode.

PRACTICAL APPLICATIONS OF X-RAYS FOR THE EXAMINATION OF MATERIALS

By W. G. BURGERS.

Differentiation between Natural and Cultured Pearls

A natural pearl is built up of concentric layers of calcium carbonate (mother-of-pearl) which have been deposited by the oyster round any available nucleus (e.g. a grain of sand). Each layer consists of crystallites with a sixfold symmetry, the axis of each crystallite being perpendicular to the layer. This structure is shown schematically in *fig. 2a*.

A "cultured" (or Japanese) pearl is obtained by inserting in the oyster a bead of mother-of-pearl cut from an oyster shell. This shell is built up of plane layers of mother-of-pearl with the sixfold axis of symmetry of the individual crystallites again arranged perpendicular to the layer. The nucleus of a cultured pearl therefore has a structure as shown schematically by the horizontally shaded

section in *fig. 2b*. Round this "nucleus" the oyster deposits several thin concentric layers so that in external appearance the finished cultured pearl cannot be distinguished from a natural pearl.

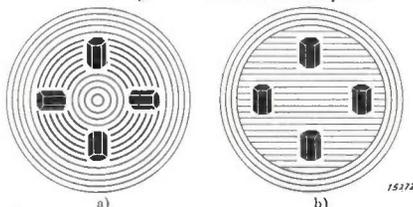


Fig. 2. Diagrammatic representation of the structure of a natural and a cultured pearl. a) The natural pearl is made up of concentric layers of mother-of-pearl (crystallites of calcium carbonate). b) The cultured pearl consists of a nucleus of plane layers of mother-of-pearl, on the outside of which several thin concentric layers have become secreted.

From the connection between the symmetry of the crystals and the X-ray pattern (which is produced by reflection of the rays at the crystal lattice planes) referred to at the beginning, it follows that the patterns obtained with a natural pearl must

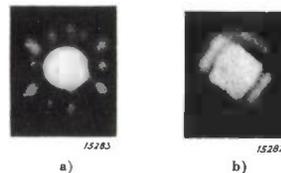


Fig. 3. a) The diffraction pattern of a natural pearl always exhibits a sixfold symmetry. b) A cultured pearl gives another type of pattern.

always exhibit an arrangement of dots with a sixfold symmetry, which latter cannot be obtained with a cultured pearl. *Figs. 3a and 3b* confirm this conclusion, and show that it is indeed possible to distinguish between natural and cultured pearls by means of X-rays.

FEBRUARY 1936

PHILAN, a local-area network based on a fibre-optic ring

J. R. Brandsma

Information is one of the main resources of modern life. The availability of up-to-date, reliable and sufficient information at the right place and at the right time is very important for the quality of public, social and economic services. This requires good means of communication, geared to the situation in which they are used. A specific situation is found in factories and offices, where large numbers of people work closely together. Within such a geographically defined area there is a need for the easy exchange of large amounts of information of all kinds and in accordance with an often varying communication pattern. For communications within this business environment ever-increasing use is being made of the 'local-area network' or LAN. At the Philips Research Laboratories Project Centre in Geldrop, near Eindhoven, a novel and versatile concept for a LAN based on a fibre-optic ring has been developed and given shape in the form of a trial system for demonstration purposes. The main aspects of this type of network, which is called PHILAN, are described in the article below.

Introduction

Few people would seriously consider buying a car if the appropriate infrastructure were not available in the form of roads and rules of the road, etc. Without those provisions there would be absolutely no way of making use of the many possibilities offered by that product of modern engineering.

A comparable situation is to be found in the field of information technology. Instead of the car we find one of the many kinds of data-processing equipment, and the infrastructure required consists of a communication network with the appropriate rules or code of practice ('protocols'). In this situation there is of course not much point in storing or rapidly processing large amounts of data at one particular place if it cannot subsequently be sent to the destination with equal rapidity and with a high degree of reliability. In situations where many people have to work closely together or use the same information, as in offices, factories, hospitals and universities, there is a great need for special communication facilities. The requirements these have to meet are the following:

- *Versatility.* Many different kinds of information have to be sent out in parallel (telephone conversa-

tions, telex, electronic mail, communication with computers and between computers, monitoring signals such as 'slow-scan TV', radiography, etc., etc.).

- *Universality.* Connections must be as universal as possible, which means to say that at every connection point it must be possible to receive and transmit a wide variety of information (provided the appropriate terminals are present).

- *Flexibility.* Connected terminals must be readily transportable without having to make radical changes such as recabling.

- *Adaptability.* It must be easy to change the kinds of information that can be transmitted and to extend the number of connections.

In the past *ad hoc* solutions to these problems have been sought, and these have tended to result in a rather confused combination of various kinds of communication networks with different types of cables, network structures, wall sockets, etc. The best solution is one that meets all the above requirements with a single communication network, called a 'local-area network (LAN)' ^[1] or, when special emphasis is placed on the great diversity of the information that can be trans-

Ir J. R. Brandsma is with Philips Research Laboratories, Eindhoven. Until recently the author led the team working on the PHILAN project at the Project Centre in Geldrop, near Eindhoven.

^[1] W. A. M. Sniijders, *Bedrijfscommunicatie en glasvezel, een symbiose, I & I (Inf. & Informatiebeleid)*, No. 7, 20-30, 1984; D. D. Clark, K. T. Pograd and D. P. Reed, *An introduction to local area networks*, Proc. IEEE 66, 1497-1517, 1978.

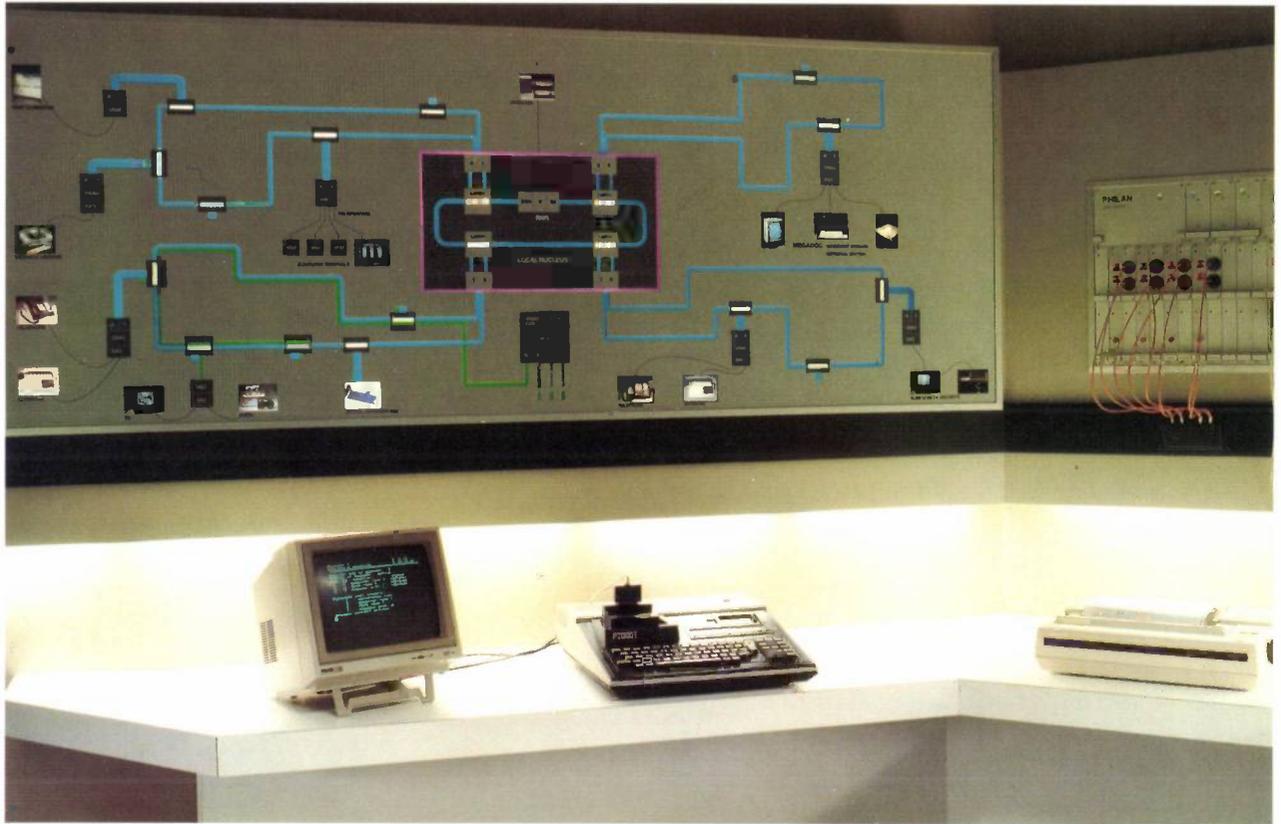


Fig. 1. At the Philips Research Laboratories Project Centre a local-area network based on the PHILAN concept has been built as a trial system for demonstration purposes. The equipment shown in the upper photograph forms the central unit, from which the whole network is controlled. Above it, on a large panel, is a diagram of the complete trial system showing the exact configuration of the

network. The lower photograph shows a large number of possible applications of the system. From left to right can be seen an advanced telephone set, a printer, a display screen specially suited to the display of documents, a word processor and an intercom set. The same transmission medium, a single optical fibre, carries the input and output signals of all these devices.

mitted, an 'integrated local-area network (ILAN)'. Owing to the local nature of these LANs, it is not in principle necessary to take account of existing communication networks — usually trunk (long-distance) networks — at the design stage. Designers consequently have a large measure of freedom in the choice of the network structure, the transmission medium and the type of signal.

In this article we shall describe the PHILAN network, a highly sophisticated ILAN network [2], recently developed at Philips, which has been set up as a trial system for demonstration purposes (fig. 1). The network has a 'meander-type' ring structure that uses fibre-optic cable. The information (e.g. speech, certain video signals and all kinds of data) is transmitted in digital form at a rate of about 20 Mbit/s, and the error probability is smaller than 10^{-10} .

First of all we shall examine the main considerations that have played a part in the design of PHILAN, such as the various requirements applicable to the transfer of various kinds of information and the possible alternative network structures. We shall then describe the composition ('frame structure') of the digital signals and consider the handling of the traffic between the users of the network in accordance with certain protocols. Finally we shall go somewhat deeper into a number of the building blocks of PHILAN, such as the ring-management system, the ring-access units and a small ingenious optical relay.

Local-area communication

The various kinds of information to be transmitted in a local-area network can be divided into two main groups: continuous information and intermittent information. The first kind is found for example in telephone and video communication and in communication for security and surveillance purposes. The other kind is found in data processing, electronic archiving and in comparable cases of communication between man and machine (computer).

Continuous information requires a connection which is permanently available and whose only side-effect is a small and constant delay. Traditionally this is realized by means of *circuit switching*, in which there is — or at least seems to be — a continuous connection between information source and destination. Intermittent information can best be sent via a data link established by means of *packet switching*, where a channel is occupied only during the transmission of a specific, clearly defined 'packet' of information at a time. Each packet is separately addressed, transmitted and received. Long messages are split into several packets. In a modern local-area network the most effi-

cient operation is achieved by combining both types of information transmission [3].

This can be done for example by first converting all information into digital form, i.e. into bit streams. These separate bit streams can be combined by means of time-division multiplex (TDM) to produce one total bit stream with a much higher bit rate, which can be transmitted via a common transmission medium (e.g. an optical fibre) to which all users are connected. The combined bit stream thus represents a considerable number of communication processes that are taking place simultaneously; the circuit-switched links each have a piece of the combined bit stream that recurs at strictly periodic intervals (e.g. a string of 8 successive bits 8000 times a second) and the packet-switched traffic receives parts of the remaining bit positions as required. We shall return to this point later.

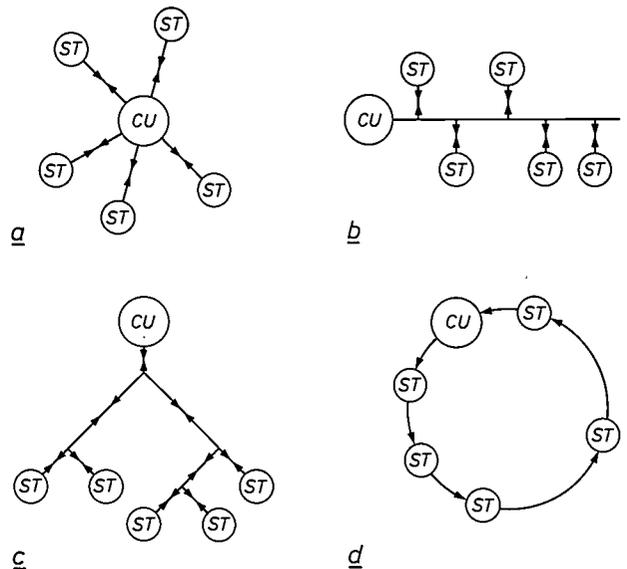


Fig. 2. Four elementary network structures. a) Star. b) Bus. c) Tree. d) Ring. CU central unit. ST station.

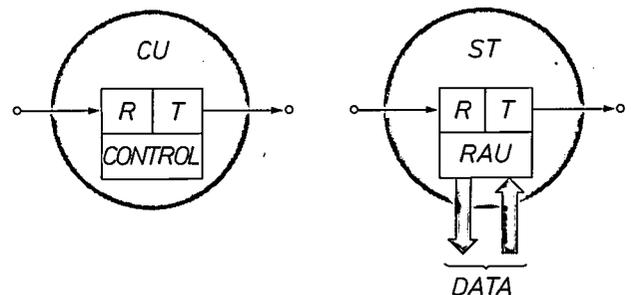


Fig. 3. Basic block diagrams of the central unit CU and a station ST of a network with ring structure. In both diagrams the receiver unit R and the transmitter unit T are shown. In CU a number of monitoring and management operations are carried out in the part marked CONTROL. The station contains a ring-access unit RAU, which handles the exchange of data between the ring and the outside world.

An important external characteristic of any communication network is the structure or topology of the network [4][5]. Fig. 2 shows four basic structures: star, bus, tree and ring. In principle a distinction can always be made between the central unit, which controls the operation of the system, the satellite stations, which function as data source or destination (data sink), and the actual transmission paths.

Every network structure has its own way of handling traffic between stations, and in every structure the role of the central unit is different. In fact with a bus structure it is common practice to decentralize the control in such a way that no separate central unit is needed. Both with the star structure and the tree structure all the data goes from the source via the central unit to the sink. In the ring structure the data circulates from station to station. If a station does not itself have any information to transmit, the incoming information is passed on to the next station. When all stations do that, the result is the formation — via the central unit — of a 'circulating memory', whose capacity depends on the total delay that occurs in the stations and the transmission paths and, of course, on the bit rate.

Another important external characteristic of a communication network is the transmission medium employed. For a LAN three kinds of media enter into consideration: twisted-wire pairs, coaxial cable and fibre-optic cable. Twisted pairs are cheap and easy to install. Coaxial cable has a higher transmission capacity and is popular because of the large number of other existing applications (such as cable TV); this accounts for the large variety of accessories available. Fibre-optic cable offers the largest transmission capacity, little signal attenuation and the unique features of complete electrical insulation and insensitivity to electromagnetic interference.

In the rest of this article we shall mainly be concerned with the ring structure, in which the signals in digital form (briefly: data) are transmitted by means of time-division multiplex. Fig. 3 shows basic block diagrams of the central unit and a station in such a ring. Both comprise a receiver and a transmitter. The minimum operation that takes place is the unaltered retransmission of the received data. The station also possesses, however, a 'ring-access unit' (RAU), with which data can be read from the ring or written to it. The central unit controls and monitors the operation of the ring; it determines the bit rate and, among its other functions, plays a part in the allocation of transmission capacity to the stations.

From the foregoing it may be deduced that the ring structure has two potential disadvantages. In the first place a ring structure is relatively vulnerable because

Table I. List of abbreviations used for component parts of PHILAN.

	Meaning
CRAU	RAU for Circuit-switched traffic
CU	Central Unit
EBS	Electronic Bypass Switch
LN	Local Nucleus
MG	Meander Gate
PRAU	RAU for Packet-switched traffic
RAU	Ring Access Unit
R/T	Receiver/Transmitter
ST	STation
US	USer
WR	Wall Receptacle
WS	Wall Socket

one defective station or one broken cable can put the whole ring out of action. In the second place the repeated reception and retransmission of data in the form of a bit stream gives rise to a gradually increasing inaccuracy in the bit positions ('jitter'), which may ultimately result in transmission errors. Timely steps therefore have to be taken in the central unit to restore the original accurate time division. This limits the maximum number of stations that can be included in the ring. Both disadvantages can be overcome by segmentation of the ring. This can be done in PHILAN and has resulted in a modified ring structure which we shall refer to as a meander ring.

PHILAN

The network structure chosen for the PHILAN system corresponds essentially to a ring, but parts of the ring ('meanders') can be cut off whenever necessary. This is done in a special unit that we call a meander gate or MG. (A list of abbreviations frequently used in this article will be found in *Table I*.) Each meander

[2] J. R. Brandsma, PHILAN: a fiber-optic ring for integrated traffic, Proc. GLOBECOM '85, New Orleans, LA, 1985, pp. 468-471;

J. L. W. Kessels, PHILAN: a LAN providing a reliable message service for real-time applications, Proc. INRIA Conference 'Advanced seminar on real-time local area networks', Bandol, France, 1986;

J. R. Brandsma, A. M. L. Bruekers and J. L. W. Kessels, Method and system of transmitting digital information in a transmission ring, U.S. Patent No. 4 553 234 (12th November 1985).

[3] The term 'integrated local-area network' (ILAN) mentioned earlier is used for communication networks in which this combination is found.

[4] C. D. Tsao, A local area network architecture overview, IEEE Commun. Mag. 22, No. 8, 7-11, 1984.

[5] D. Hutchison, J. A. Mariani and W. D. Shepherd (eds), Local area networks: an advanced course (Proceedings, Glasgow 1983), Lect. Notes Comput. Sci., Vol. 184, Springer, Berlin 1985.

contains a limited number of stations which are interconnected in accordance with the usual ring structure (fig. 4). The transmission medium used in the meanders is fibre-optic cable. The central unit and the meander gates together with the links between them form the local nucleus of PHILAN. This is concentrated at one location and, because of the short distances *inside* the local nucleus, the links used are purely electrical connections.

In each meander gate an interference-free clock signal from the central unit is used to form the signal to be sent to the next meander gate. In this way any jitter that may have arisen in a meander is eliminated, so that in theory an unlimited number of meanders can be connected in cascade.

A more detailed block diagram of a meander is given in fig. 5, corresponding to the part A-A' in fig. 4. The total digital signal that appears at a given moment in the PHILAN system arrives (in electrical form) at the input of the meander gate. During normal operation an optical equivalent of this signal is formed in the transmitter, which sends it on to the meander. Each possible position in the meander where a station may have to be connected is provided with an optical connector box, called a 'wall receptacle'. If no station is connected, each wall receptacle forms an optical through-connection. If a station is included in the ring, it is connected in the usual way by a special optical plug (see fig. 3). Situated at the end of the meander is the optical receiver of the meander gate. Here the optical signal from the meander is converted into an electrical signal and processed for transmission to the next meander gate.

Meander length

The total length a meander can have depends on the number of wall receptacles it contains. Upon passing each wall receptacle to which an *active* station is connected the circulating digital PHILAN signal is regenerated (refreshed) before being sent on; any signal attenuation that may have occurred earlier is thus eliminated. A less favourable situation is found in each wall receptacle to which *no active* station is connected. The signal then passes straight through the wall receptacle, involving a nominal attenuation of 1.3 dB. The worst situation occurs when all stations are passive except for one of the two that lie closest to the meander gate. The total attenuation occurring in the meander is then at a maximum. At a given maximum attenuation there is a direct relation between the maximum meander length and the number of wall receptacles in the meander. This relation is shown in fig. 6 for a total 'power budget' of 32 dB, of which more than 22 dB is available for losses in wall receptacles and fibre-optic cable, given an attenuation of 5 dB/km in the optical fibre. With the very common number of 10 wall receptacles, the maximum meander length is about 2 km.

If for one reason or another a meander is not functioning properly, the electronic bypass switch (EBS,

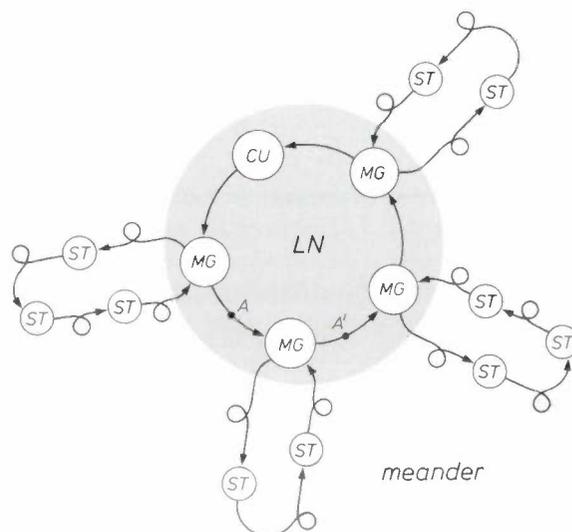


Fig. 4. The meander ring structure of PHILAN. In addition to the central unit *CU* and the stations *ST* found in every ring structure, the PHILAN ring contains a number of meander gates *MG*. These divide the ring into sections, called meanders, each containing a number of stations. The *CU*, the *MGs* and the links between them form the local nucleus *LN* (grey). The transmission medium used in the meanders is fibre-optic cable (indicated by the symbol $\text{---} \circ \text{---}$); electrical connections are used in the local nucleus.

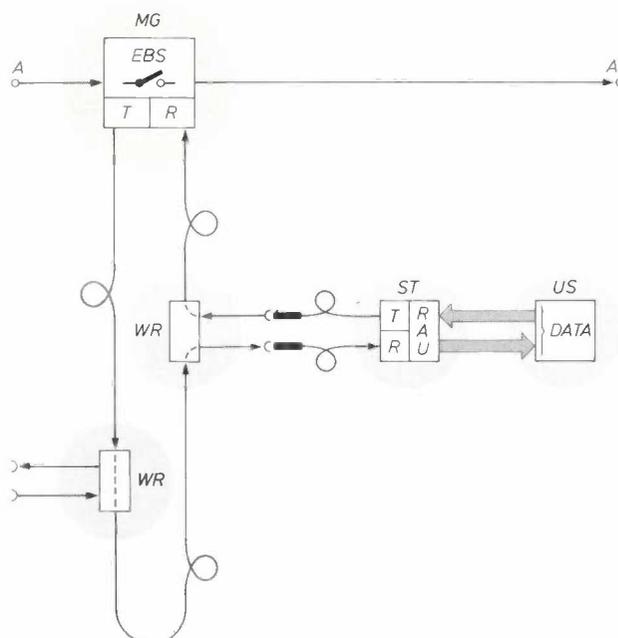


Fig. 5. Basic diagram of a PHILAN meander. The meander gate *MG* forms the link between the actual meander and the rest of the network. This link is established via the optical transmitter *T* and the optical receiver *R*. In addition, *MG* contains an electronic bypass switch *EBS*, which can 'cut off' the meander. At a number of fixed positions along the fibre-optical cable of the meander are wall receptacles *WR*. A user *US* can be connected to a *WR* via a station *ST* and a plug, *WR* forms an (optical) through-connection. *ST* contains the same optical receiver/transmitter combination *R/T* as *MG*. *US* is coupled to the ring-access unit *RAU* of *ST* via an electrical connection.

see fig. 5) is closed in the meander gate, thereby cutting out that meander. The rest of the PHILAN system is therefore not affected by the malfunctioning of this meander. This is not the only protective provision made; in each wall receptacle the station connected can be optically bypassed by means of a relay, in the same way as when there is no plug in the socket. This safeguard comes into operation automatically if it is found, for example, that the transmitting unit of a station is not sending out any optical signals. Yet another safety feature is the possibility of making a

direct electrical through-connection in a station — between the output of the receiver and the input of the transmitter. This is a feature that plays a particularly important part in the localization of errors ('error diagnosis').

During the design of PHILAN, reliability was a paramount consideration. The aim throughout was to provide *safe communication services*, not only as regards network structure and hardware but also as regards the signals and the protocols used. By this we mean that after a certain time interval the user must have an assurance that his message has arrived correctly at the right destination, without him having to do anything else to verify this. Before looking at this more closely we shall first examine the signal structure used in PHILAN.

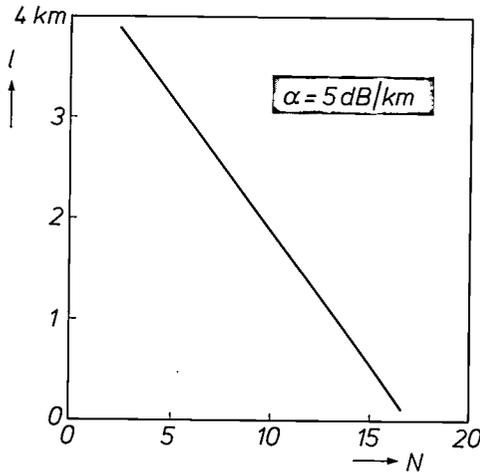


Fig. 6. For each connection point some attenuation of the signal circulating in the ring must be taken into account. In PHILAN the number of connection points N per meander and the maximum meander length l are related by the straight line shown (for a fibre-optic cable with an attenuation of 5 dB/km and a maximum permissible total attenuation in wall receptacles and glass fibre of just over 22 dB).

The signal structure: frames and subrings

Signal transfer in PHILAN takes place in the form of a continuous series of bits passing through all connected stations in succession at a rate of 20.48 Mbit/s. Each group of 2560 successive bits is identically organized and is called a 'frame'. A frame corresponds to 125 μ s. Measures are taken in the central unit to ensure that the signal delay corresponding to one complete circuit or *lap* around the ring is always equal to an integral number of frames. The frames therefore continue to appear at each point in the ring at a fixed rate.

The first 16 bits of a frame always form a special and unchanging synchronization word, called the 'preamble', so that all parts of the PHILAN system remain in synchronism with each other and 'know' exactly when a new frame starts. The remaining bits of the frame belong together in groups of 8 bits ('bytes'). The grouping of the bits of a particular byte from all successive frames results in a transmission channel ('subring') with a capacity of $8/(125 \times 10^{-6}) = 64$ kbit/s. Similarly a combination of N bytes produces a subring of $N \times 64$ kbit/s (fig. 7). The bits in one particular frame that belong to the same subring form a 'field'; the size of a field may vary from 1 to 64 bytes (i.e. from 8 to 512 bits).

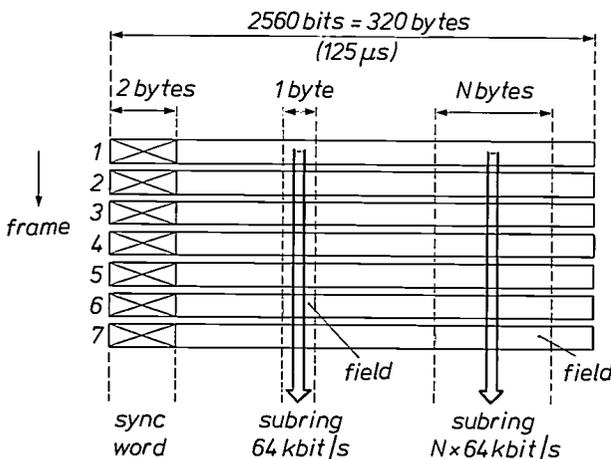


Fig. 7. Frame structure of the digital signal used in the PHILAN ring. Every 2560 successive bits form one frame consisting of 320 bytes of 8 bits. The first two bytes form an unchanging synchronization word, called the 'preamble'. By joining together the bits of one particular byte from each frame a transmission channel ('subring') is formed with a capacity of 64 kbit/s. By taking N bytes together a subring of $N \times 64$ kbit/s is obtained. The share of each subring in a particular frame is called a 'field'.

Between the preamble and the fields that form the subrings there are a further 6 bytes in each frame that have a special function. These are used for intercommunication between the units of the local nucleus (CU and MGs) or between the units of a particular meander (MG and STs). Details will not be given in this article.

Communication within PHILAN takes place by the bits of a particular subring being written by the transmitter of one station and read out by the receiver of another station. Writing and read-out can also take

place in a station simultaneously, in which case it is referred to as a 'swap' [6]. The bits coming from the ring and those destined for the ring are stored in the stations in buffers, which form a kind of separation between the fast (20 Mbit/s) transmission on the ring and the preceding or subsequent processing. Processing rate and transmission rate are thus made independent of each other. One advantage of this is that no processing operations need to be carried out at transmission rate.

All communication takes place through the mechanism of the subrings: this includes both the actual data transfer in packet-switched or circuit-switched traffic and all incidental messages concerning matters such as origin and destination (data source and data sink) and related to the operation of the whole network. Three distinct types may be distinguished: traffic subrings, memory subrings and insertion subrings. Each subring is completely independent; what happens on one particular subring has no influence whatever on the others.

Traffic subrings

As the name suggests, these are the subrings in which the actual traffic in PHILAN takes place, and they account for about 90% of the total ring capacity.

In circuit-switched traffic a subring with the required capacity of $N \times 64$ kbit/s is made available for the whole duration of a 'conversation', which in principle may be unlimited. Even in two-way ('full-duplex') traffic between two stations only one subring is required if both stations use the swap operation.

In packet-switched traffic a subring is made available for the transfer of one packet (a maximum of 8 kbytes). This subring covers a field of 1, 4, 16 or 64 bytes per frame; on this basis there are four categories of traffic subrings with a capacity of 64, 256, 1024 and 4096 kbit/s respectively. Owing to the limited dimensions of the fields, a packet will in general be distributed over more than one frame. Messages that are longer than the longest possible packet are in turn divided into a number of packets.

Traffic subrings will generally be referred to here from now on as 'channels'.

Memory subrings

Memory subrings function as a 'circulating memory' for various system data that may be of interest to all connected stations. They are used for example to allocate a channel to a particular station for a particular time for the transmission of a data packet. This takes place in principle in the following way. In each frame the channel controller (a part of the central unit) places the number of a free channel — if there is

one — in a field of the memory subring. By means of a swap, each connected station may substitute a zero for this number and thus 'seize' the channel concerned. The channel controller also continuously performs the swap operation on the memory subring and finds out, by receiving the number zero, when the channel is occupied. When a channel becomes free, this is communicated to the channel controller by returning the appropriate number in a similar way via the memory subring.

The channel controller monitors the free channels continuously. Each of the four categories of channels for packet-switched traffic has its own memory subring. Depending on the actual volume of messages the channel controller may transfer channels from one category to another by combination or splitting. The channel controller also ensures that no channel numbers can become permanently lost or are duplicated in circulation as a consequence of transmission errors.

Insertion subrings

In each frame two times 8 bytes are reserved for forming two 'insertion subrings'. They each therefore have a transmission capacity of 512 kbit/s and are used in packet-switched traffic for exchanging information ('protocol messages') between the various connected stations. One insertion subring is used for sending protocol messages that connect a transmitter and a receiver for the duration of the complete message ('call messages'). The other insertion subring is used for exchanging protocol messages during actual packet-switched traffic between partners that have already been connected. These messages prevent packets from being lost through overflow of storage

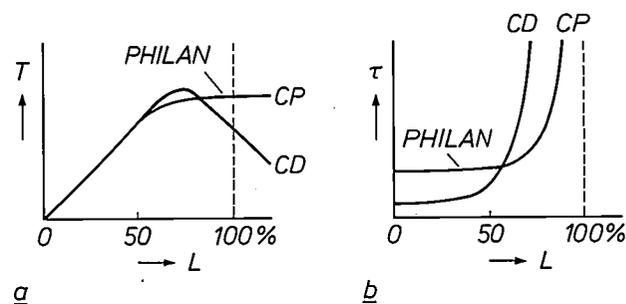


Fig. 8. Comparison of the network stability for packet-switched traffic in two types of LAN. The first type of network, of which PHILAN is an example, offers the special feature of collision prevention (CP) between the transmitted packets. In the other type of network mutual interference is permitted and subsequently restored ('collision detection' or CD). a) With increasing traffic load L , the effective throughput T in PHILAN increases monotonically to the maximum value; in the other system, T becomes smaller at large values of L . b) The price to be paid for this benefit is a somewhat longer average delay τ for the packets in the PHILAN network when the traffic load L is small. At high values of L , PHILAN again has the advantage. (If L is 100% or higher, the mean delay in both systems is infinitely large.)

buffers in the receiver (flow control) or as a consequence of transmission errors (error control).

The name of this type of subring derives from the manner in which it is used: each station can insert a protocol message of 8 bytes into a field in this subring by means of a swap operation. The present contents of the field are read into a receiving buffer and replaced by new contents. The present contents may however represent a valid message for a subsequent station on the ring and should not therefore be destroyed. In the next frame it is consequently re-inserted (again via the swap operation). In this way each protocol message goes through the whole PHILAN system until it finally comes back like a kind of echo at the station which originally sent it out. This station then finally removes the message from the ring.

The total time that elapses before an echo arrives depends on the number of stations in which swapping

An illustrative example

The function of the various subrings may be illustrated by the example given in *fig. 9*. This shows two users *US-A* and *US-B*, between which a packet-switched connection has to be built up via the stations *ST-A*, *ST-B* and the rest of the PHILAN system. *US-A* and *US-B* are connected by means of a standardized link known as a 'VMEbus' [7] with *ST-A* and *ST-B* respectively. These stations have a ring-access unit specially designed for packet-switched traffic and are therefore designated *PRAU*. Let us assume that *US-A* is the initiator of the connection and that *US-B* is not engaged; this means that *PRAU-B* constantly monitors the first insertion subring, which is used for sending call messages.

When *US-A* wants a connection, he sends a request for it, including the address of *US-B*, over the VMEbus to *PRAU-A*. The latter places a call for *US-B* on the first insertion subring by means of the swap operation. This call also contains the address of *US-A*. In response, *PRAU-B* performs two actions. He places a call for *US-A* on the same insertion subring to indicate that he is ready to receive information, and he informs *US-B* that a message is about to be received. When the call for *US-A* by *PRAU-A* has been received, *US-A* is requested, via the VMEbus, to fill the packet buffer of *ST-A* with a packet. *US-A* does this and at the same time indi-

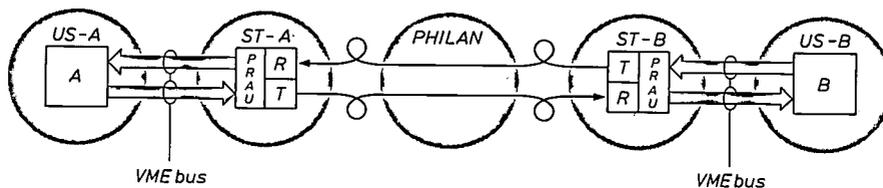


Fig. 9. Simplified representation of a packet-switched link between users *US-A* and *US-B*. Between each user and the corresponding station *ST* electrical signals are exchanged via a standardized connection, called a VMEbus. The stations are connected by fibre-optic cables with the rest of the PHILAN system (not further specified here). *PRAU* ring access unit for packet-switched traffic. *R* optical receiver. *T* optical transmitter.

has been carried out, but the maximum delay is known beforehand. The correct reception of an echo at the place of origin means that the protocol message in question has passed correctly through all operating stations. If, on the other hand, no correct echo is received within the maximum delay, the same protocol message is sent once again. This procedure forms the basis on which the correct transmission of information can be guaranteed in the PHILAN system ('safe-message service').

One of the main advantages of the signal structure and associated protocols adopted for PHILAN is network stability: an increasing volume of traffic does *not* — as in some other networks — lead to a lower effective transmission capacity, caused by the fact that an ever-decreasing proportion of the packets are transmitted correctly at the first attempt (*fig. 8*). In addition, special measures have been taken to ensure that PHILAN functions fairly, i.e. without any discrimination against transmitters or receivers in the allocation of transmission capacity.

cates how large the packet is and also whether it is the last packet of this message.

PRAU-A then occupies a free traffic subring (channel) by means of a swap on the memory subring in the category with the capacity selected by *PRAU-A*. Next, *PRAU-A* sends a start message to *PRAU-B* via the second insertion subring and at the same time gives information about the channel used, the size of the packet and whether it is the last packet or not. *PRAU-B* then switches to reception on the relevant channel, waits for the arrival of the packet and writes it into his packet buffer. *PRAU-B* then places the number of the used channel into the memory subring by means of a swap operation to indicate that the channel is free again. In addition an

[6] Read-out alone does not affect the signal circulating on the ring; this only happens when new information is written to certain bit positions (possibly as part of a swap operation).

[7] A number of manufacturers in the computer field (in particular Motorola, Mostek and Signetics/Philips) have defined a collective standard for communication links in applications of 8-bit, 16-bit and 32-bit microprocessors. It is called the 'VMEbus', but is now also referred to as the 'IEC 821 bus' or the 'IEEE P1014'. The abbreviation VME stands for 'Versa Module Europe', but the full name has now no more than historical interest. An organization called VITA (VMEbus International Trade Association) has been set up to disseminate information on the VMEbus. It has its own journal, 'VMEbus Systems', published by Intratech Communications, St Clair Shores, Michigan 48081, U.S.A.

error-detection procedure is carried out to check whether the packet has been correctly received and, if this is so, the packet is made available to *US-B* via the VMEbus. Finally, confirmation of this and a request for the next packet are sent by *PRAU-B* to *PRAU-A* via the second insertion subring. *PRAU-A* then requests *US-A* for the next packet, again occupies a free channel, and so on.

If *PRAU-B* discovers that a particular packet has not been correctly received, he immediately sends a message back to *PRAU-A* via the second insertion subring with a request for the relevant packet to be retransmitted. This is done as soon as a free channel has been found. This entire procedure is repeated for as long as is necessary for all packets of a particular message to be sent correctly from *US-A* to *US-B*. The connection between *US-A* and *US-B* is cancelled as soon as *PRAU-A* has received confirmation that *PRAU-B* has received the last packet of the message correctly.

The actual procedure differs in one particular respect from what has been suggested so far: each *PRAU* does not have only one packet buffer; it has two — equivalent — packet buffers of 64 kbits each. While *US-A* fills one packet buffer of *PRAU-A*, the contents of the other packet buffer of *PRAU-A* are sent to one of the packet buffers of *PRAU-B*, and at the same time the other packet buffer of *PRAU-B* is emptied by *US-B*. If the transmission procedure takes place at the maximum rate, packets of the same message are handled simultaneously at three places in the PHILAN system: on the actual ring and on the two VMEbuses of *A* and *B*. In this way messages of any given length can be transmitted from *A* to *B* at a maximum rate of about 4 Mbit/s.

The component units of PHILAN

The stations

As we have seen (fig. 5), a station can be divided into a receiver/transmitter combination *R/T* and a ring-access unit *RAU*. The two main tasks of *R/T* are to perform the opto-electrical and the electro-optical conversions of the digital signals circulating in the PHILAN ring. A more detailed block diagram of *R/T* is given in fig. 10. Whereas the electrical signals be-

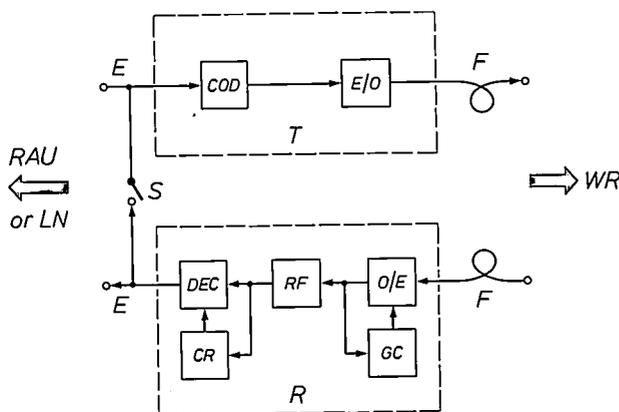


Fig. 10. Optical receiver/transmitter combination *R/T* used in each station (and also in each meander gate) of the PHILAN system. *E* electrical connection. *F* fibre-optic cable. *COD* biphasic encoder. *E/O* electro-optical converter. *O/E* opto-electrical converter for the signal coming from the fibre-optic cable. *GC* gain control for the electrical signal. *RF* receive filter. *DEC* biphasic decoder. *CR* clock recovery. *S* electrical bypass switch. *WR* wall receptacle. *RAU* ring-access unit. *LN* local nucleus.

tween *R/T* and *RAU* take the form of a simple bit stream (with no encoding), an elementary form of encoding (biphase coding^[8]) has been adopted, for various technical reasons connected with transmission^[9], for the digital signals to be transmitted by optical fibre (fig. 11). The transmitter has therefore been provided with an encoder and the receiver with a special receive filter, a clock-recovery circuit and a decoder.

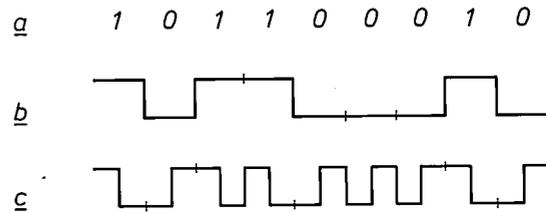


Fig. 11. *a*) Arbitrary bit stream consisting of ones and zeros. *b*) Corresponding 'uncoded' binary signal; a '1' is represented by a high level and a '0' by a low level. *c*) The same bit stream after biphasic encoding. A '1' is now represented by a sequence of a high level and a low level, an '0' by the reverse. This encoding gives a transmission signal with a better frequency spectrum, and the required clock signal is easier to recover at the receiver end.

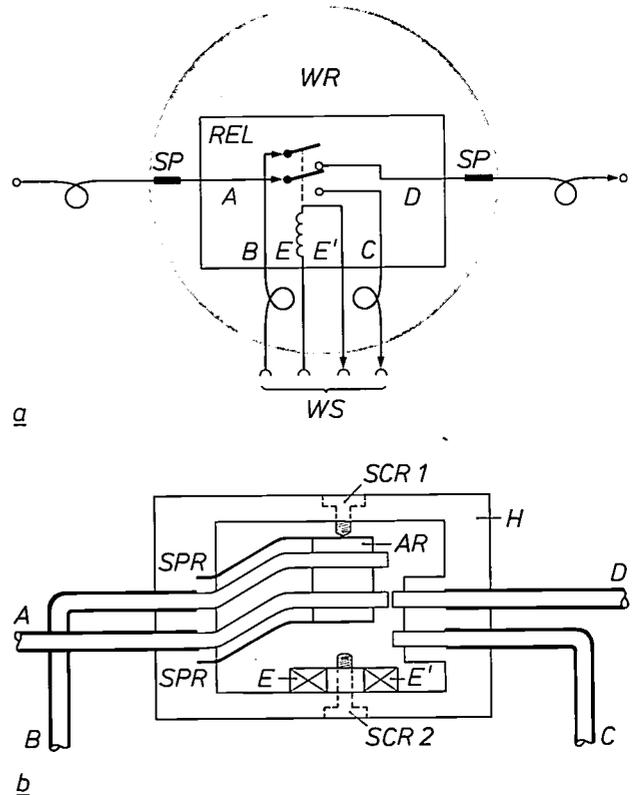


Fig. 12. *a*) Basic diagram of the wall receptacle *WR*. Two optical splices *SP* are used to include *WR* permanently in the meander ring. A station *ST* can be connected by means of a wall socket *WS*, which contains two optical contacts and two electrical contacts. The heart of *WR* is a relay *REL* which controls two optical switch contacts. *b*) Cross-section through the relay *REL*. *A, B, C, D* optical fibres. *SPR* leaf springs. *E, E'* connections for the relay coil. *SCR1, SCR2* set screws. *AR* armature. *H* housing^[11].

There are two versions of the ring-access unit, each built around a standard microprocessor (type '68000'): a CRAU for circuit-switched traffic and a PRAU for packet-switched traffic. At the user end a CRAU has a connection (in both directions) of 2.048 Mbit/s. This gives the user access to 32 channels on the PHILAN ring, each with a capacity of 64 kbit/s, or to one or more combinations of a number of these channels.

must be an economic proposition. The answer to this problem has been found in the use of a specially developed optical relay capable of switching fibre-optic cables. The basic diagram of the wall receptacle is shown in *fig. 12a*, a cross-section of the relay used in it is given in *fig. 12b*, and a photograph of the relay is shown in *fig. 13*. When the coil of the relay is not energized (e.g. when there is no plug in the socket) the

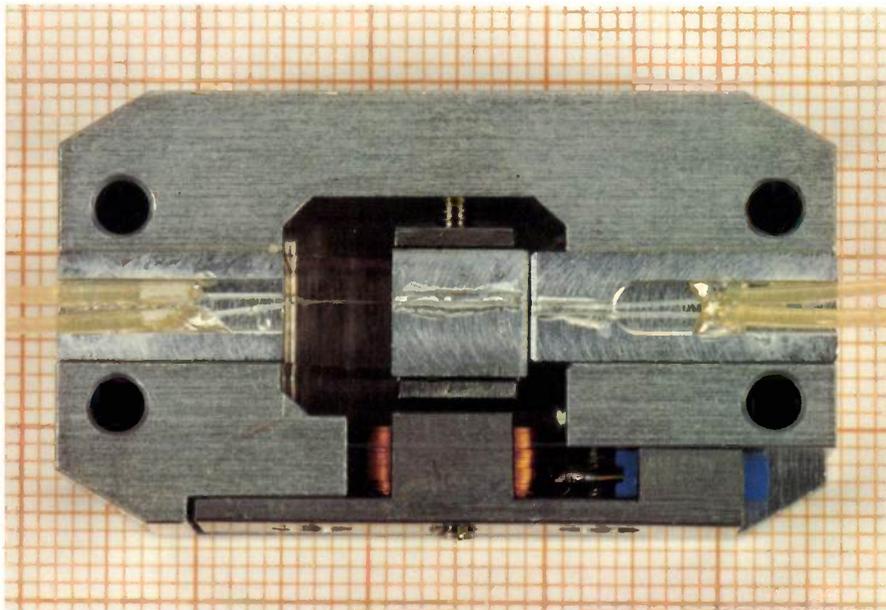


Fig. 13. Photograph of the interior of the relay *REL* in *fig. 12*. The optical fibres are bonded by adhesive into V-grooves, which are made in a single operation in both the housing and the armature during manufacture. This means that there is no need for extremely tight tolerances. Set-screws on either side determine the two extreme positions of the armature. The two leaf springs that attach the armature to the housing are hardly visible in this photograph.

Signalling information between a user and a CRAU is exchanged via a VMEbus^[7].

The VMEbus is also used with the PRAU for exchanging the actual packets and the related messages between user and station. As already mentioned, packets can be sent over the PHILAN ring at different rates, from 64 to 4096 kbit/s. The maximum packet length is always 64 kbits (8 kbytes); longer messages have to be sent in more than one packet. Each PRAU possesses hardware for error detection after the reception of each packet. This is done by means of the Cyclic Redundancy Code (CRC)^[10].

The wall receptacles

In PHILAN the optical wall receptacles occupy a place of their own. Their first task is to make it possible to modify the number of PHILAN connections in a simple and reliable way. In addition they must be capable of forming an optical 'bypass' and, in spite of the critical mechanical requirements found here, they

relay forms an optical through-connection; otherwise the optical signals go via the connected plug (*fig. 14*). The nominal attenuation in the relay is only 0.7 dB.

The meander gates

The three main functions of the meander gates were touched upon earlier in this article (see *fig. 5*):

- Using the same type of receiver/transmitter combination *R/T* found in each station, an opto-electrical connection and an electro-optical connection are

[8] See for example p. 348 in the special issue: F. W. de Vrijer, Modulation, Philips Tech. Rev. 36, 305-362, 1976; see also: P. J. van Gerwen and W. A. M. Sniijders, Communication system for bi-phase transmission of data and having sinusoidal low-pass frequency response, U.S. Patent No. 4573 169 (25th February 1986).

[9] Such as measures to deal with the large spread (by a factor of 1000) in the strength of the received optical signal, to facilitate clock generation in the receiver and to optimize the signal-to-noise ratio after reception.

[10] See p. 129 in: A. S. Tanenbaum, Computer networks, Prentice-Hall, Englewood Cliffs, NJ, 1981.

[11] W. A. M. Sniijders and J. P. Klomp, Optical switch, European patent application No. 0181657-A1 (21st May 1986).

established between the local nucleus and the individual meanders.

- In appropriate cases the *EBS* switch can be used to 'cut off' a particular meander.
 - Any jitter produced in the meander is eliminated.
- In addition the synchronization pattern (preamble) at the beginning of each frame is regenerated.

To enable the electronic bypass *EBS* to be switched on and off without interfering with the rest of the



Fig. 14. Wall socket and plug for the PHILAN meander ring. This provides a reliable and easily detachable double optical and double electrical connection (see fig. 12a) [12].

PHILAN system (e.g. with the synchronization), the signal delay between input and output of the meander gate must be the same before and after switching. Special measures are taken to ensure that this delay always amounts to 256 bits.

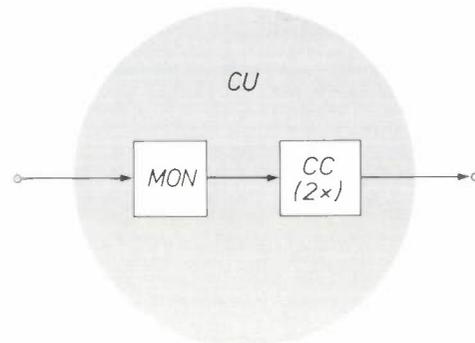
In addition, when *EBS* comes into operation, the meander gate has to take a number of extra measures to guarantee the 'safe-message service' of PHILAN. This service is based on operation with echoes (page 17) and in certain circumstances these echoes can give rise to erroneous conclusions, since a correct echo will *not* have passed through an operational station in any cut-off meander.

Finally, each meander gate also monitors the synchronization of all the stations in the corresponding meander and takes special measures to restore the synchronization if this is lost.

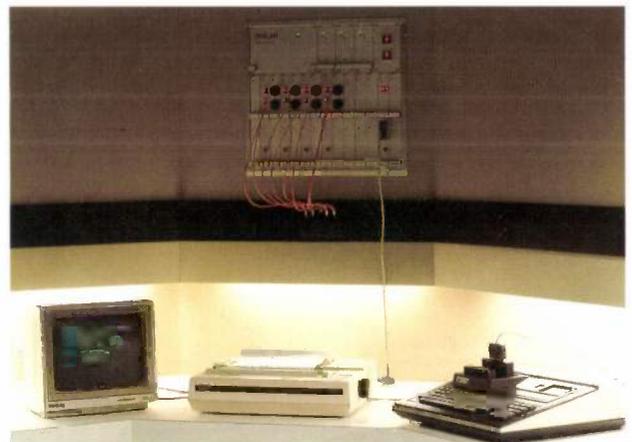
The central unit

The central unit *CU* of PHILAN (fig. 3) functions as a 'ring management system'. In its simplest version it contains a monitor and two channel controllers for circuit-switched and packet-switched traffic, respectively (fig. 15).

The main tasks of the monitor are the generation and monitoring of the correct frame structure of the PHILAN signal, and to make up the delay experienced



a



b

Fig. 15. a) In its simplest form the central unit *CU* in PHILAN consists of a monitor *MON*, which monitors the signal structure, and two channel controllers *CC*. If required, other equipment can be added for error diagnosis, management and direct control by an operator. b) Photograph of the extended central unit used in the demonstration system at the Project Centre. The specific PHILAN circuits are contained in a wall rack (top centre). Below it can be seen, from left to right, a display screen, a printer and a P2000 computer, with a keyboard and a cassette deck. These are used as the input and output peripherals for the operator. (See also fig. 1.)

by the signal in one lap of the meander ring to an integral number of frame periods.

The channel controllers are responsible for recording free channels and making them available as required.

[12] J. W. Faber and J. P. Klomp, Optical connector device, European patent application No. 0189609-A1 (6th August 1986).

Possible extensions

So far we have confined the description of PHILAN to the basic configuration; there are, however, a number of important extensions that can be made. Two of these will be mentioned briefly to conclude this article. In the first place it is possible in principle to connect the central units of several PHILAN systems with each other so as to form larger transmission networks. In the second place, specially developed hardware can be added to the central unit, the meander gates and the stations to increase the diagnostic and control capabilities of the system, enabling the system to make regular fully automatic checks on the functioning of all meander gates and stations. One such check is the measurement of the optical transmitting power of each station. All information about the state of the network is recorded in the central unit and can be displayed on a screen or communicated to an operator by means of alarm signals. Conversely, the operator using the keyboard can give instructions to the central unit to carry out specific tasks.

By making the appropriate additions to the PHILAN basic configuration the optimum local-area network can thus be created for practically any situation.

Summary. PHILAN is an integrated local-area network in which digital information is transmitted at a bit rate of about 20 Mbit/s in time-division multiplex. The system has a ring network structure, with sections called meanders that can be cut off in the event of local malfunctioning. The transmission medium used in the meanders is fibre-optic cable. Both circuit-switched and packet-switched traffic can be handled simultaneously in the PHILAN system. For packet-switched traffic the use of specific protocols gives practically 100% reliability ('safe-message service'). Connections to the meanders are made via wall receptacles, in which a double optical plug connection is established. Each wall receptacle contains a specially developed relay for electromechanical switching of two fibre-optic cables. In addition to a relatively simple basic configuration, there are more elaborate versions of the PHILAN system with sophisticated provisions for error diagnosis and control. At the Philips Research Laboratories Project Centre a PHILAN trial system has been built for demonstration purposes.

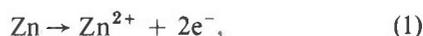
Investigation of a new type of rechargeable battery, the nickel-hydride cell

J. J. G. Willems

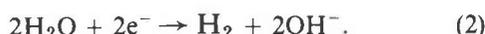
The author describes the research work on a new type of rechargeable battery, which has resulted in a hermetically sealed nickel-hydride cell. The new cell has stable electrode material and an effective control of the hydrogen and oxygen flow in the system. These features have enabled a first experimental version to operate well for more than a thousand cycles of high-rate charge and discharge.

Brief historical background

The first galvanic cell — or more accurately 'battery' of galvanic cells — was Alessandro Volta's pile (the 'voltaic pile') of 1800, which consisted of a stack of sheets of zinc and silver (or copper), separated from each other by pieces of cardboard soaked in a saline solution (*fig. 1*)^[1]. As became clear much later, the electrical energy from each individual cell is due to the oxidation of the zinc to form zinc ions, with the release of electrons:



and to the reduction of water to gaseous H_2 , with the uptake of electrons:

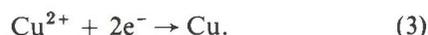


Why these particular materials fulfil these functions can be explained by their relative position in the electrochemical series of the elements (*fig. 2*). This thermodynamically determined position in the series is not the only significant factor, however: it is also necessary to take the kinetics of the process into account.

The silver (or copper) is not used up in the voltaic cell, but acts as a catalyst in the reduction of the water; its only other function is that of a conductor of current (*fig. 3*). Zinc is still used as the material for negative electrodes, for example in dry batteries, because of its negative electrode potential, but the use of

water (zero electrode potential in *fig. 2*) as the oxidizer at the positive electrode has remained confined to this first galvanic cell: the electromotive force (e.m.f.) available from this element is only 0.5 V.

In 1836 John Daniell developed a galvanic cell in which copper ions from a copper-sulphate solution are reduced at a positive copper electrode to metallic copper for the required electron transfer:



The negative electrode consisted of a rod of amalgamated zinc, which was immersed in dilute sulphuric acid. The amalgamation prevents an excessive side

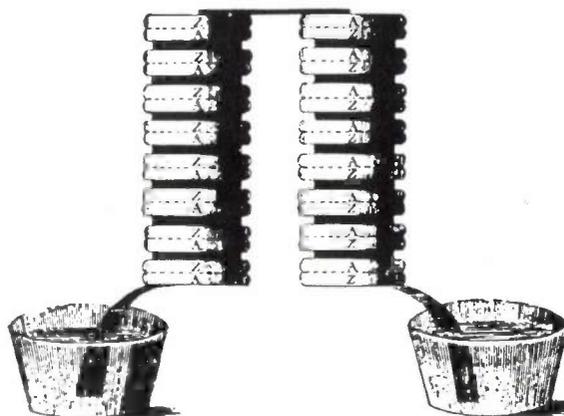


Fig. 1. The voltaic pile, as illustrated in the first publication (in *Phil. Trans. R. Soc.* 90, 403, 1800)^[1]. A stands for silver (argent) and Z for zinc.

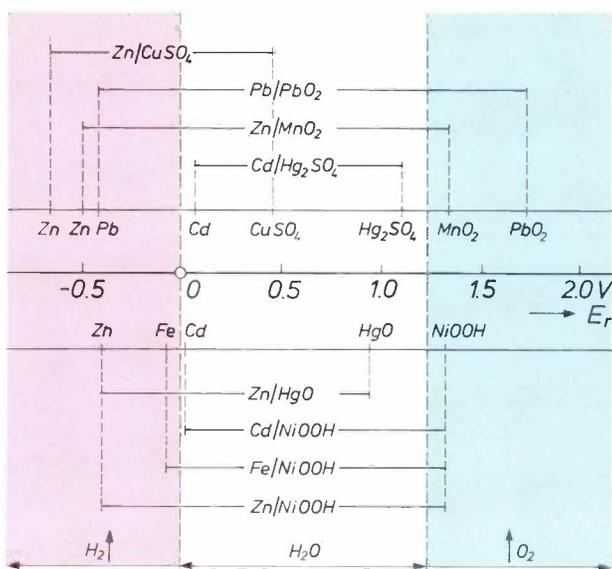


Fig. 2. The electrode potential E_r of some of the more familiar electrode materials (combined with an electrolyte), determined with respect to the potential of a reversible hydrogen electrode in the same electrolyte. Oxidation potential increases from left to right (from non-noble to noble in the electrochemical series of the elements). If the electrode potential is lower than 0 V (the electrode potential of the voltaic couple water/hydrogen), the electrode material tends to reduce H_2O to H_2 (the red region); if the electrode potential is higher than +1.23 V (the electrode potential of the couple water/oxygen), the material then tends to oxidize H_2O to O_2 (the blue region). Only in the Weston cell (Cd/Hg_2SO_4) do both electrode couples show no tendency to decompose water. The electrode potentials of materials used in combination with an acid or neutral electrolyte are shown above the horizontal axis; the potential values of materials used in an alkaline medium are shown below the horizontal axis.

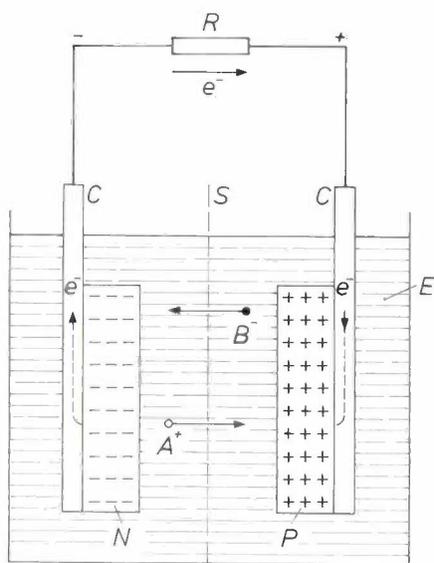


Fig. 3. Diagram of a current-producing galvanic cell. P active positive electrode material. N active negative electrode material. e electrons. E electrolyte. A^+ positive ions. B^- negative ions. C current collectors. S separator. R external load.

reaction between the zinc and the sulphuric acid. Even with the cell not working, hydrogen would otherwise have been liberated and the zinc would have been used up. (This process is equally detrimental to the operation of the voltaic pile.) To prevent copper ions from being deposited on the zinc — this happens because of their position in the electrochemical series — the two electrodes are contained in separate compartments, each with its own electrolyte. A porous diaphragm, originally of cow's gullet, permitted the necessary transfer of ions between the electrodes.

The first rechargeable battery, Gaston Planté's lead-acid cell, dates from 1860. It consisted of two interwound lead sheets separated by flannel strips and immersed in sulphuric acid. The electrodes were formed by repeated charging from primary cells. (Current generators like the dynamo had already been invented, but they were not yet widely used.)

The development of dry cells also dates from the same period, about 1865. In the dry cell the oxidizing substance is a solid (a metal oxide) and the electrolyte is effectively 'immobilized', so that the battery can be made 'leakproof'. Georges Leclanché used a mixture of manganese dioxide and carbon for the positive electrode, and a carbon rod inserted into the mixture was the current conductor. The electrolyte consisted originally of a saturated solution of ammonium chloride; in a later development this was thickened (immobilized) by means of starch or lime, or both. Amalgamated zinc, which was then coming into use for the manufacture of batteries, again served for the negative electrode. Today, more than a century later, the Leclanché cell is still the most widely used dry battery.

Apart from their use as a current source, batteries were soon required to meet the need for a constant-voltage reference source. These 'standard cells' operate in phase equilibrium; when they are being used it is only the position of the equilibrium that changes, though only very slightly. The nature of the substances reacting in the cell does not change. In the Weston cell (1892) the positive mercury/mercurous-sulphate electrode and the negative cadmium amalgam/cadmium-sulphate electrode are immersed in a saturated solution of crystals of these salts. The e.m.f.

[1] Further historical data will be found in:
 G. W. Vinal, Primary batteries, Wiley, New York 1950, see especially the first chapter, pp. 1-24;
 G. W. Heise and N. C. Cahoon, The primary battery, Vol. 1, Wiley, New York 1971, first chapter, pp. 1-58;
 S. U. Falk and A. J. Salkind, Alkaline storage batteries, Wiley, New York 1969, first chapter, pp. 1-41.
 For the general principles of the battery discussed here, see for example:
 V. S. Bagotskii and A. M. Skundin, Chemical power sources, Academic Press, New York 1980;
 J. O'M. Bockriss *et al.* (eds), Comprehensive treatise of electrochemistry, Vol. 3, Plenum, New York 1981.

of the Weston cell is 1.01463 V at 25 °C and its variation with temperature is only -4.05×10^{-5} V/°C.

A major objective in the development of batteries was, and still is, to make the most efficient use of the relatively expensive electrode materials, and this is most easily done with rechargeable batteries. To start with such batteries will be required to have no spontaneous side-reactions that give an unproductive loss of active material, in other words minimum 'self-discharge'. This requirement led to the development of alkaline accumulators at about the turn of the century. These included the nickel-cadmium cell and the nickel-iron cell (Junger and Edison). Alkaline reagents are in general less reactive with metals and metal oxides than acids.

It is difficult to avoid such non-productive side reactions completely, since most electrode materials are not thermodynamically stable in combination with aqueous electrolytes and are thus always inclined to react with water with the formation of hydrogen or oxygen, depending on the value of the electrode potential. As can be seen in fig. 2, it is only for the pair of electrodes in the Weston cell that the position of the electrode potentials is such that no spontaneous decomposition of water can take place. A possible way around this difficulty would be to use non-aqueous electrolytes, but the usual practice has been to manage with electrodes that, although not thermodynamically stable, are sufficiently stable kinetically. In other words, the electrode materials are chosen or treated in such a way (for example by amalgamation) that their rate of reaction with the electrolyte is very slow.

The extent to which unwanted amounts of hydrogen or oxygen are nevertheless produced then often depends on the state of charge of the battery and on the charging rate. This is because the battery can only be fully charged if the charging voltage is *greater* than the cell e.m.f., which implies moving further into the unstable region (fig. 2). The charging voltage will have to be higher the faster the battery is to be charged. Overcharging the battery, which is usually unavoidable because it is impossible to know exactly when the battery is fully charged, also causes the formation of hydrogen and oxygen. The same applies to overdischarge, which may be the consequence of differences in capacity between the series-connected cells of a battery. Here a cell that has discharged first can be forced by adjacent cells to discharge itself still further, a process inevitably accompanied by the formation of hydrogen and oxygen.

This problem of unwanted gas generation becomes particularly serious when it is intended to produce a completely sealed rechargeable battery, one that can be used in all positions and without requiring mainte-

nance, preferably for many years. The dry battery may be seen as an early answer to this problem, without however achieving the goal of rechargeability. Moreover, as pre-war users will have found to their displeasure, the 'leakproof' dry battery was by no means the same thing as the hermetically sealed battery now in view.

Good results in this respect have now been achieved with nickel-cadmium cells. Since 1950 these cells have been provided with an excess of cadmium hydroxide, $\text{Cd}(\text{OH})_2$, at both electrodes (fig. 4). In overcharging and overdischarging the current tries to find a path, oxidizing or reducing materials with the appropriate characteristics, and this can give a high gas pressure in a sealed cell. The excess of $\text{Cd}(\text{OH})_2$ 'channels' any such unwanted current so that it causes no damage, or very little.

If on overcharging the nickel electrode starts to form oxygen, the cadmium electrode will still not be

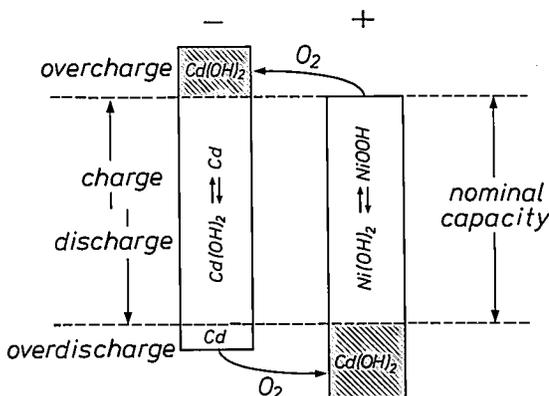


Fig. 4. Diagram illustrating the protection from overcharging and overdischarging in a sealed cadmium battery, introduced in about 1950. Adding $\text{Cd}(\text{OH})_2$ to both electrodes introduces an oxygen cycle during both overcharging (blue region) and overdischarging (red region). The added $\text{Cd}(\text{OH})_2$ is indicated by hatching. An excess of cadmium (red) is also added at the cadmium electrode; this is explained on p. 32.

fully charged because of the excess of $\text{Cd}(\text{OH})_2$. It will therefore not generate hydrogen (fig. 2) and continue to convert $\text{Cd}(\text{OH})_2$ into metallic Cd. This conversion, however, is now no longer used for the process of energy storage but for reducing the oxygen formed at the nickel electrode and retaining it as $\text{Cd}(\text{OH})_2$. So from that stage onwards there is no change in the state of charge of the two electrodes (since the nickel electrode was already fully charged).

In overdischarging the electrodes change roles as compared with the roles they took in overcharging. Because of the addition of $\text{Cd}(\text{OH})_2$ to the nickel elec-

trode, this electrode in the empty state will take on the potential of the cadmium electrode, and so will again be able to counteract the evolution of hydrogen. The oxygen produced at the empty cadmium electrode in overdischarging is reduced in a similar reaction cycle at the nickel electrode. It could be said that there was a deliberately engineered 'chemical bypass', which channels both forms of an unwanted passage of current via an oxygen cycle.

As mentioned, because of the position of its electrode potential, the excess of $\text{Cd}(\text{OH})_2$ is in principle capable of *preventing the evolution of hydrogen gas* (both upon overcharging at the cadmium electrode and overdischarging at the nickel electrode). This is therefore important because the evolution of hydrogen is a much greater problem than that of oxygen, since hydrogen gas once formed is *not* converted at the nickel electrode owing to kinetic impediments, and consequently a channelling comparable with the oxygen cycle is not possible.

Even though $\text{Cd}(\text{OH})_2$ can be introduced into the battery as explained, the evolution of hydrogen can never be completely eliminated. Some hydrogen may form because of the ageing of $\text{Cd}(\text{OH})_2$, particularly during charging at high currents and on overdischarging. Hermetically sealed NiCd cells are therefore always provided with a safety valve that opens if too much hydrogen gas accumulates.

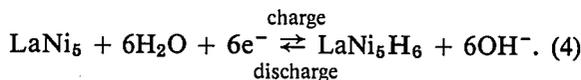
For some time now a similar expedient has been used in the lead-acid accumulator so that it can also be produced as a hermetically sealed battery. The greater reactivity of the acid electrolyte makes it even more difficult in this case to prevent the accumulation of gas, which requires at the very least a specially adapted cell construction, at the expense of durability.

The investigations at Philips Research Laboratories that will now be described should be viewed against this background of efforts to achieve a hermetically sealed rechargeable battery that can be used in all positions, has a very long service life, can be charged and discharged at a very high rate and can withstand considerable overcharging and overdischarging.

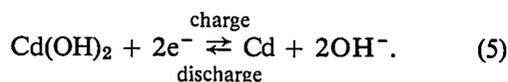
Initial study of LaNi_5 and LaNi_4Cu as electrode material

In this section we shall be considering the use of hydrogen as 'fuel' for a rechargeable battery, and not as the product of an unwanted side reaction, as it so often was in the previous section. This implies that it must be possible to oxidize hydrogen at the surface of an electrode while the battery is being discharged, and that hydrogen can be stored in this electrode while the battery is being charged. Our investigations at Philips

Research Laboratories of materials for such an electrode were mainly made with LaNi_5 and related compounds. Water from an aqueous electrolyte can be reduced electrochemically at these materials, the hydrogen evolved can be stored as a hydride and the stored hydrogen can then be reoxidized to form water:



The formation of hydrogen in this case is therefore, once again, not an unwanted side reaction but part of the chemical system responsible for the storage of energy. LaNi_5H_6 can be used here as the active material of the *negative* electrode, which is converted into LaNi_5 during the discharge; the LaNi_5H_6 is then reformed by the charging. The metal hydride may be compared in this respect with the cadmium at the negative pole of the NiCd cell:



If we compare the electrode reactions of LaNi_5 with those in Volta's cell — since water is converted in both — we see, however, that the reactions in the voltaic pile have exactly the opposite function. In the voltaic pile water is converted into hydrogen, and electrons are taken up, during discharge at the positive pole, whereas with LaNi_5 this occurs during charging at the negative pole (and then the conversion in the opposite direction on discharge results in the production of electrons).

Reactions of this kind are also used with the familiar reversible hydrogen electrode, with platinum as the electrode material. The difference compared with LaNi_5 in this case is that the hydrogen gas has to be supplied, since hydrogen is not absorbed by platinum. An electrode of this type cannot therefore be used for the storage of electrical energy.

About fifteen years ago it was discovered that intermetallic compounds such as LaNi_5 and TiFe are capable of absorbing large amounts of hydrogen gas and of desorbing it at pressures of the order of 1 atmosphere and at room temperature^[2]. Fig. 5 shows the absorption and desorption isotherms of the $\text{LaNi}_5\text{-H}_2$

[2] Publications on hydrogen-absorbing compounds discovered in the past 15 years include:
H. Zijlstra and F. F. Westendorp, *Solid State Commun.* 7, 857-859, 1969;
J. H. N. van Vucht, F. A. Kuijpers and H. C. A. M. Bruning, *Philips Res. Rep.* 25, 133-140, 1970;
J. J. Reilly and R. H. Wiswall Jr., *Inorg. Chem.* 13, 218-222, 1974;
F. A. Kuijpers, *Philips Res. Rep. Suppl.* 1973, No. 2;
H. H. van Mal, *Philips Res. Rep. Suppl.* 1976, No. 1;
K. H. J. Buschow, P. C. P. Bouten and A. R. Miedema, *Rep. Prog. Phys.* 45, 937-1039, 1982.

system at a number of temperatures [3]. As the figure illustrates, in a temperature range up to about 80 °C the system has a phase transition between a hydrogen-poor phase and a hydrogen-rich phase at a constant equilibrium pressure, called the plateau pressure. At room temperature LaNi_5 absorbs six hydrogen atoms at a plateau pressure of 2 atmospheres, and at a pressure of 1.6 atm the hydride formed is converted back to LaNi_5 , accompanied by the desorption of hydrogen gas. The amount of hydrogen that can be absorbed and desorbed in this way per unit volume is enormous: 40% greater than the hydrogen density in liquid hydrogen at 20 K. Even at a temperature of 80 °C large amounts of hydrogen gas can still be reversibly stored in LaNi_5 at a pressure of less than 20 atmospheres. The resultant absorption-desorption hysteresis, which is greater at higher plateau pressures and gives an energy loss, is of course an unwanted effect. Compared with TiFe , LaNi_5 has lower plateau pressures at the same temperatures and a much lower absorption-desorption hysteresis.

On repeated absorption and desorption of hydrogen gas, LaNi_5 disintegrates into small fragments (fig. 6), forming a fine powder with a specific surface area of 0.25 m^2/g . The powder has a high catalytic activity, which, combined with the high diffusion coefficient for H atoms in the metal lattice, results in high sorption rates. At 20 °C complete hydrogen desorption takes no more than a few minutes. In addition, LaNi_5 and LaNi_5H_6 are both good electrical conductors. All in all, LaNi_5 therefore seems to be a promising electrode material for a rechargeable battery.

How does LaNi_5 behave in the presence of the electrolyte, however? To find this out we carried out a number of experiments with this material and related compounds in a 'half-cell' arrangement. The materials, which are very hard and brittle, were pulverized by repeated hydrogen gas absorption/desorption cycles, and then mixed with powdered copper. This mixture was cold-pressed to form a porous pellet, which was found to satisfy all electrode requirements [3]. The electrolyte used in the cell was a 6 M KOH solution, a platinum plate was used as the counter-electrode and an $\text{Hg}/\text{HgO}/6 \text{ M KOH}$ electrode was used as the reference electrode (fig. 7) [3]. Fig. 8 shows a charge and discharge curve for an LaNi_5 electrode as measured in this cell. As can be seen, there is close correspondence with the absorption and desorption curves in fig. 5. Starting from the fact that the absorption of one hydrogen atom is associated with the transfer of one electron (eq. 4), it can be shown that the theoretical charge-storage capacity for the composition LaNi_5H_6 is 372 mAh/g, which is slightly

more than the theoretical storage capacity for a cadmium electrode (366 mAh/g). This is important because our aim was to be able to use the LaNi_5 elec-

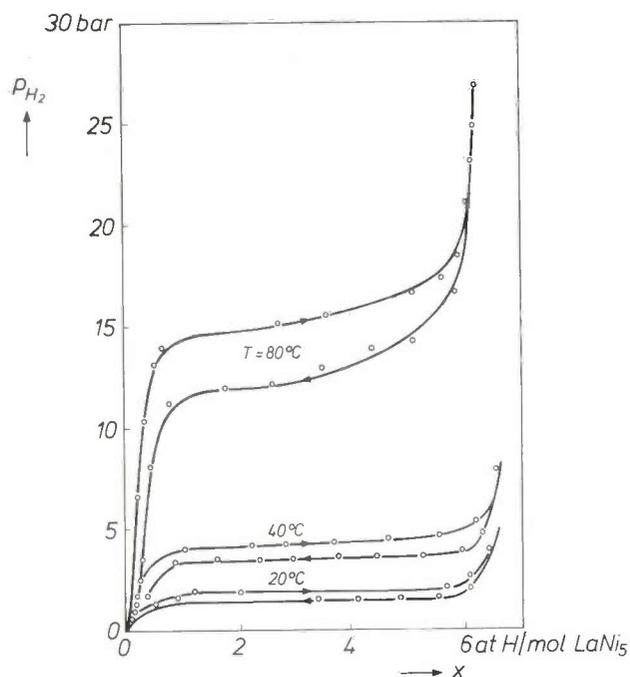


Fig. 5. Absorption and desorption isotherms (hydrogen pressure p_{H_2} as function of the concentration of hydrogen atoms x) of the $\text{LaNi}_5\text{-H}_2$ system at 20, 40 and 80 °C.

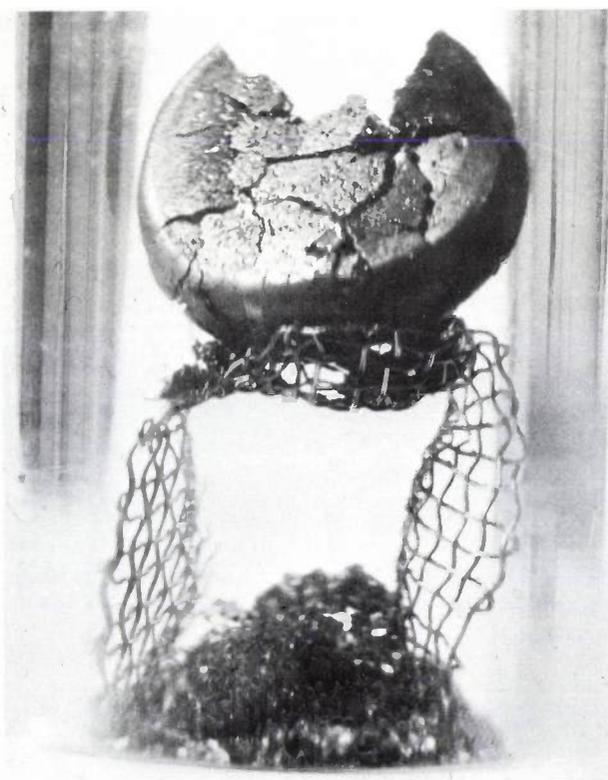


Fig. 6. If LaNi_5 is made to absorb and desorb hydrogen repeatedly, it gradually disintegrates into small fragments, seen here in a small heap under the piece of metal gauze acting as a sieve. (This is the effect of a single absorption of hydrogen by an LaNi_5 casting.)

trode as an alternative to the cadmium electrode in a nickel-cadmium battery, for reasons that will appear shortly.

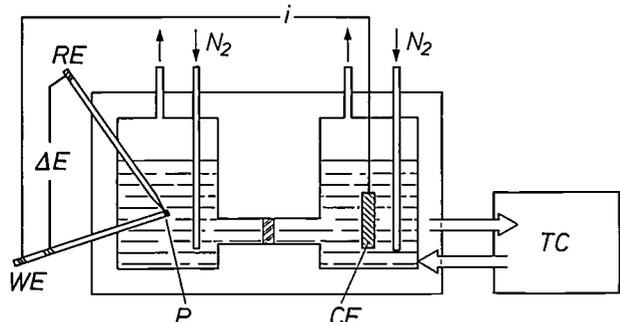


Fig. 7. Diagram of the experimental arrangement used for many of the 'half-cell' electrochemical measurements. *WE* working electrode; a porous pellet *P* of the active material being investigated, e.g. LaNi_5 , is attached to its tip. *RE* reference electrode, usually an $\text{Hg}/\text{HgO}/6\text{ M KOH}$ electrode. *CE* platinum counter-electrode. The electrolyte is kept free from oxygen by purging with nitrogen. The half-cells have a water-jacket and are kept at a temperature of 25°C by a thermostat *TC*.

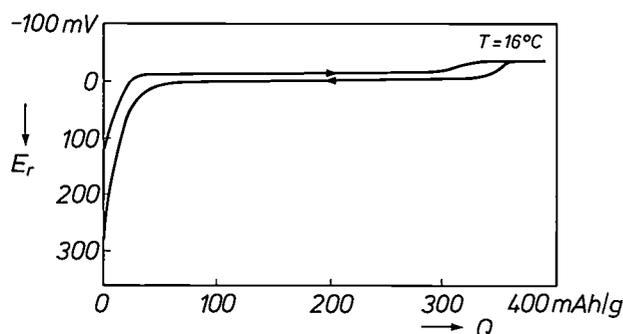


Fig. 8. Charge and discharge curve for an LaNi_5 electrode, with respect to the reversible hydrogen electrode, as a function of the charge Q , measured in 6 M KOH at 16°C .

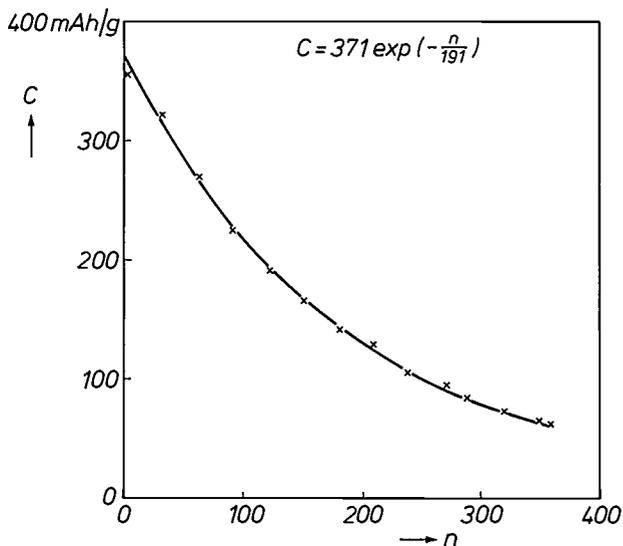


Fig. 9. Storage capacity C of an LaNi_5 electrode as a function of the number of charge/discharge cycles (n).

Although LaNi_5 has these desirable properties, it has one major disadvantage. The electrochemical experiments discussed soon revealed that the storage capacity of LaNi_5 electrodes decreases drastically after repeated charging and discharging. As can be seen in *fig. 9*, the storage capacity then decreases exponentially — by as much as 40% after only 100 cycles. This rules out LaNi_5 as an electrode material for rechargeable batteries — at least if it is to meet the strict requirements formulated in the previous section.

In our investigation we therefore set ourselves the primary goal of establishing the cause of this drastic decrease in storage capacity. We included in our investigations the compound LaNi_4Cu , which had been investigated previously at our Laboratories. This material also absorbs large amounts of hydrogen, but it does so at a lower plateau pressure more suitable for practical purposes, and gives the same exponential decrease in storage capacity on repeated charging and discharging.

This experimental investigation relied mainly on X-ray diffraction analyses and electron microscopy^[3], in addition to the electrochemical methods described. *Fig. 10* shows X-ray diffraction diagrams of powdered specimens of LaNi_5 , LaNi_4Cu and their hydrides, obtained after exposing the specimens to hydrogen gas at 150 atm for 24 hours. Diagrams *a* and *d* give the pattern of a phase consisting of hexagonal crystals (space group $P6/mmm$). Diagrams *c* and *e* show the same pattern as *a* and *d*, with the difference that all peaks have clearly been shifted to smaller diffraction angles, indicating the same hexagonal structure but with larger lattice spacings, resulting in a volume increase of 20% or more, due to the absorption of hydrogen.

We studied the effect of the periodic absorption and desorption of hydrogen by repeatedly charging and discharging electrodes of LaNi_5 and LaNi_4Cu in a solution of 6 M KOH , and then making diffraction diagrams of powders of the electrode material thus treated (*fig. 11*). After a number of charge/discharge cycles a number of small peaks and 'shoulders' start to appear in the diagrams — in addition to the diffraction pattern of the original material and large peaks at 21.7° and 25.3° , which are attributable to the copper powder. The small peaks are due to $\text{La}(\text{OH})_3$ and the

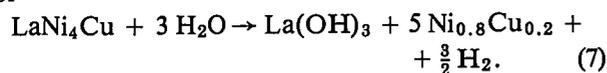
[3] A more extensive and detailed description of these investigations will be found in the thesis by J. J. G. Willems, Metal hydride electrodes; stability of LaNi_5 -related compounds, Philips J. Res. 39, Suppl. No. 1, 1984. Major contributions to this research were made by J. R. van Beek and H. C. Donkersloot; see for example: J. R. van Beek, H. C. Donkersloot and J. J. G. Willems, in B. W. Baxter (ed.), Power Sources, Vol. 10 (Proc. 14th Int. Power Sources Symp., Brighton 1984), Academic Press, London 1985, pp. 317-338.

'shoulders' to free Ni (or to $\text{Ni}_{0.8}\text{Cu}_{0.2}$). These diffraction peaks become distinctly larger as the number of cycles increases.

The conclusion is that the periodic charging and discharging of the electrode in 6 M KOH assists the conversion of LaNi_5 (or LaNi_4Cu) in the presence of water into $\text{La}(\text{OH})_3$ and Ni (or $\text{Ni}_{0.8}\text{Cu}_{0.2}$):



or



Scanning electron microscopy (SEM) gave the following supplementary picture of the degradation process in LaNi_5 on periodic electrochemical charging and discharging.

Fig. 12a shows the state of the material in an unused electrode. The cracks that can be seen were produced when the electrode was made, during the pulveriza-

tion in an H_2 atmosphere (fig. 6). After 12 charge/discharge cycles (fig. 12b) the first needles of $\text{La}(\text{OH})_3$ appear at the surface, and the number of cracks and tears has been increased. In fig. 12c, after 25 cycles, the number of needles is considerably larger, and the grain size has fallen on average by a factor of three. In fig. 12e the specimen is completely overgrown with $\text{La}(\text{OH})_3$, and in fig. 12f (after 359 cycles) the conversion of LaNi_5 into $\text{La}(\text{OH})_3$ is practically complete.

The 'driving force' for the decomposition of LaNi_5 in 6 M KOH is the strong affinity of LaNi_5 for water. The gain in Gibbs free energy for this oxidation reaction is 472 kJ/mol of LaNi_5 — nearly *four times* the heat of formation of LaNi_5 .

The actual oxidation reaction will occur at the interface between the metal and the electrolyte. As we have

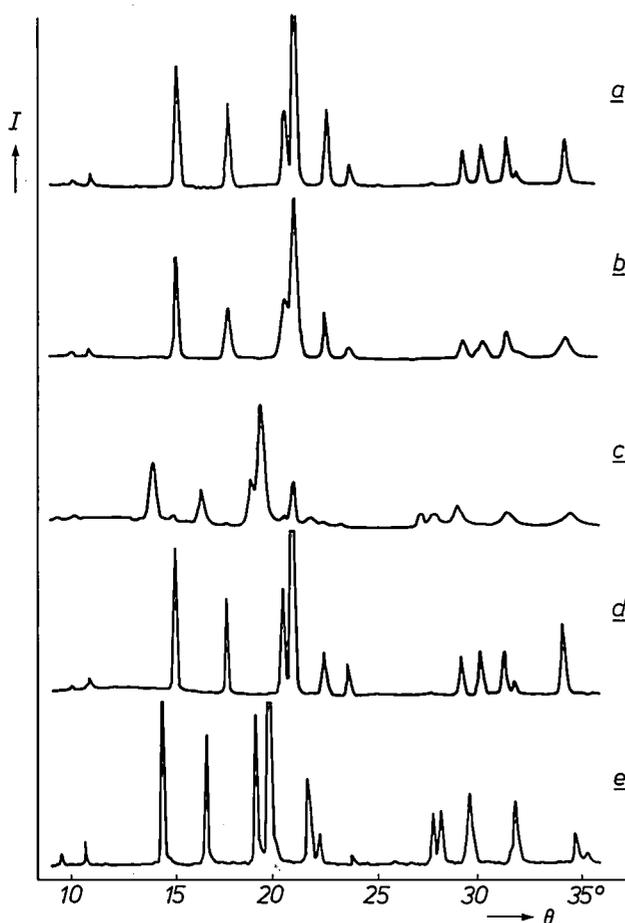


Fig. 10. X-ray diffraction diagrams of powdered specimens of LaNi_5 and LaNi_4Cu before and after hydrogen absorption (exposure to hydrogen gas at 150 atm for 24 hours). a) Ground LaNi_5 . b) LaNi_5 pulverized by 100 cycles of absorption and desorption of hydrogen. c) LaNi_5H_6 , formed by causing the material of (a) to take up hydrogen. d) Ground LaNi_4Cu . e) The same material as (d) after absorption of hydrogen.

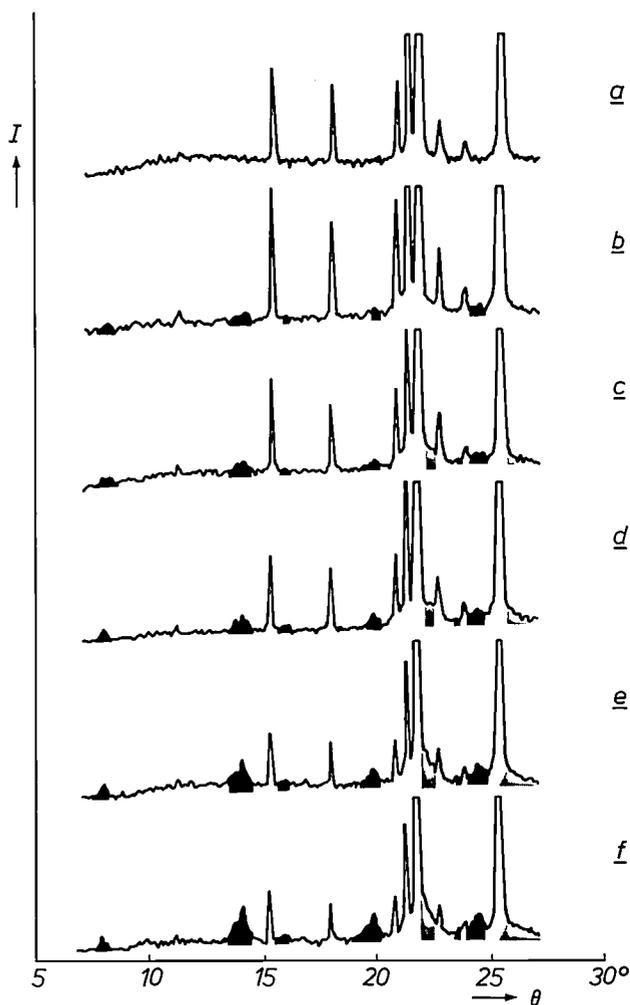


Fig. 11. X-ray diffraction diagrams of dismantled LaNi_4Cu electrodes that have been charged and discharged a number of times in 6 M KOH. The photographs show the effect of zero (a), 23 (b), 65 (c), 87 (d), 148 (e) and 188 (f) charge/discharge cycles. The positions of the reflection peaks due to $\text{La}(\text{OH})_3$ are shown in black, and those due to $\text{Ni}_{0.8}\text{Cu}_{0.2}$ are shown in grey. The large peaks at 21.7° and 25.3° can be attributed to the copper powder added to the electrode material. The other peaks are due to LaNi_4Cu .

seen, however, the conversion is not limited to the surface, but is virtually complete. This is remarkable, since the mobility of lanthanum at room temperature in the LaNi_5 lattice is virtually zero. Nor is the rate at which this lanthanum is converted into $\text{La}(\text{OH})_3$ proportional, as might be expected, to the area of the active surface, but rather to the *amount of active material remaining*, i.e. to the amount of stored hydrogen, as is evident from the exponential decay of the storage capacity. This suggests that the decomposition of LaNi_5 is connected with the absorption and desorption process.

This led us to the view that the 'respiration' of the lattice might produce a kind of peristaltic movement, which would drive the lanthanum present inside the lattice up to the surface.

The search for a more stable electrode material

From the work of H. H. van Mal, K. H. J. Buschow and F. A. Kuijpers [2] at our Laboratories it was already known that the partial substitution of Co for Ni in LaNi_5 had the effect of considerably reducing the increase of volume on hydriding. Altogether we made some twenty intermetallic AB_5 compounds in which La and Ni were partially replaced by other metals. As can be seen in *fig. 13*, the substitution of Co for 3.3 Ni atoms has the strongest effect in the desired direction. The substitution of 1 or 2 atoms has much less effect, and the substitution of 4 or 5 atoms gives rise to widely deviating effects, which are not yet properly understood [3]. Apart from a diminished decrease of capacity on repeated charging and discharging, we see that in all cases the initial capacity is also reduced.

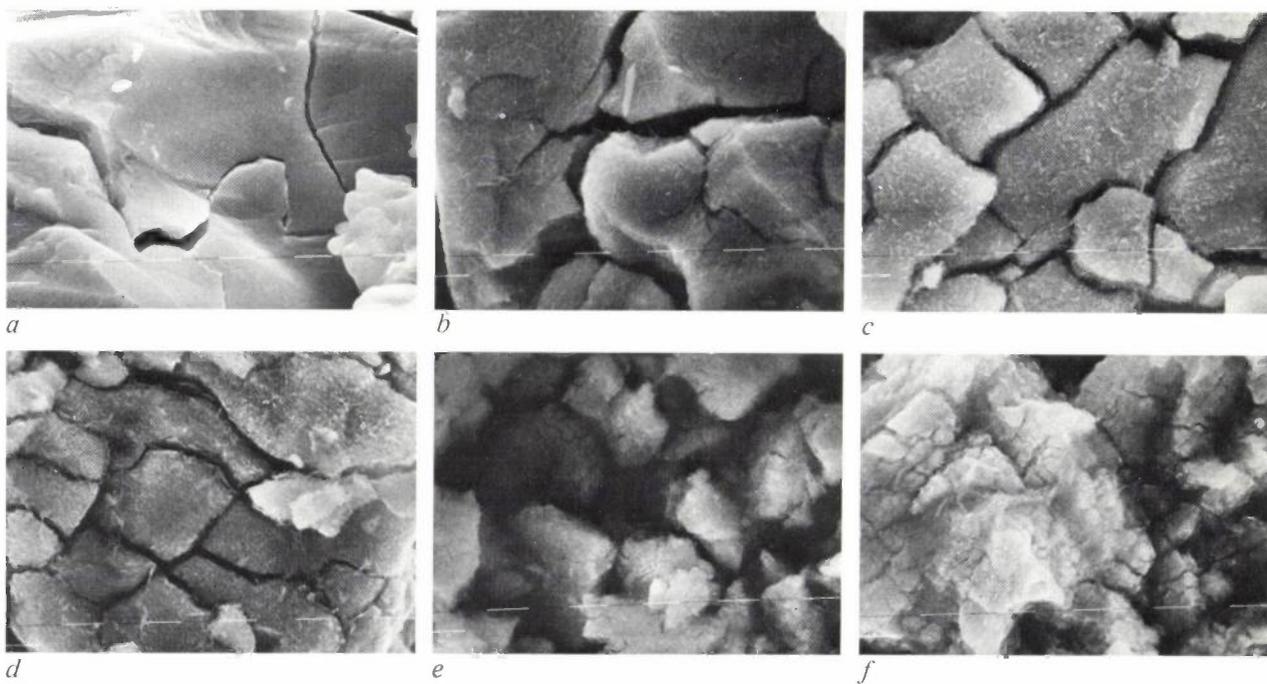


Fig. 12. SEM micrographs of LaNi_5 electrodes after repeated charging and discharging. The results are shown after zero (a), 12 (b), 25 (c), 63 (d), 176 (e) and 359 (f) charge/discharge cycles. One 'dash' represents 1 μm . $\text{La}(\text{OH})_3$ needles, which are visible from (b) onwards, have completely overgrown the original material in (f).

The mechanical stresses associated with the 'respiration' of the lattice will of course increase with the changes in volume that take place on the absorption and desorption of hydrogen. In the case of LaNi_5 the increase of volume due to the absorption of hydrogen may be as much as 24% of the initial volume.

We therefore looked for related hydrogen-absorbing intermetallic compounds that expand less during this absorption than LaNi_5 , and thus give less diffusion and oxidation of lanthanum.

The stability of some compounds, particularly those with 2 or 3 atoms of Co (*fig. 14*), can be significantly improved by adding a small amount of aluminium to the constituent components in the preparation of the compounds. The stability of LaNi_5 itself is not affected by this small addition. A similar improvement in stability is also obtained by adding small amounts of silicon. Both additives would be expected to impede the diffusion of lanthanum in the metal lattice or cause a protective skin to form on the surface.

No comparisons will be given of all the substitutions that we investigated or their effects. We found that the compound $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ answered our purpose extremely well. This electrode material has an initial capacity of 290 mAh/g and the capacity remaining after a thousand charge/discharge cycles is still 70% (fig. 15). The excellent properties of this compound are also illustrated by the SEM micrographs in fig. 16. These show that after 666 cycles there are only relatively few $\text{La}(\text{OH})_3$ needles present and that the fragments formed still have the jagged character of uncorroded crystals (see fig. 12).

The question remained as to whether the stability improvement achieved really was connected with a diminished expansion in volume on hydriding. To start with, we established that the decrease in the capacity of the electrodes only occurs during repeated charging and discharging and is not the consequence of mere contact with 6 M KOH over a long period. We also measured the volume increment ΔV due to H_2 absorption by X-ray diffraction methods and plotted the results against the fraction S_{400} of the capacity remaining after 400 cycles (fig. 17). The measured points are grouped along two straight lines in the figure, showing that a decreasing volume increment does in fact result in improved stability, and that this is most pronounced in compounds that contain small amounts of Al or Si.

We thus arrived at the following picture for explaining the corrosion of LaNi_5 -type compounds. The rate-determining step in the corrosion process is the transport of lanthanum to the interface between the electrode and the electrolyte. Under normal conditions the mobility of lanthanum in the metallic AB_5 lattice is virtually zero at room temperature. This is not the case at the phase boundary. As indicated in fig. 18, a phase boundary is always present during charging and discharging in the grains of the intermetallic compound. During charging and discharging continuous 'recrystallization' takes place from a hydrogen-poor phase to a hydrogen-rich phase, and vice versa. Since the hydrogen-rich phase has a larger unit cell (figs 10 and 17) than the hydrogen-poor phase, the lattice at the phase boundary is deformed, so that the mobility of the lanthanum in this zone can assume high values. During charging and discharging this phase boundary, and hence the zone of enhanced mobility, migrates through the grains, and therefore the lanthanum can be transported from the bulk of the grains towards the surface, even at room temperature. The lowest mobility is found when the volume expansion due to hydriding is smallest, of course.

Small amounts of aluminium or silicon lead to the formation of a closed skin at the surface of the grains,

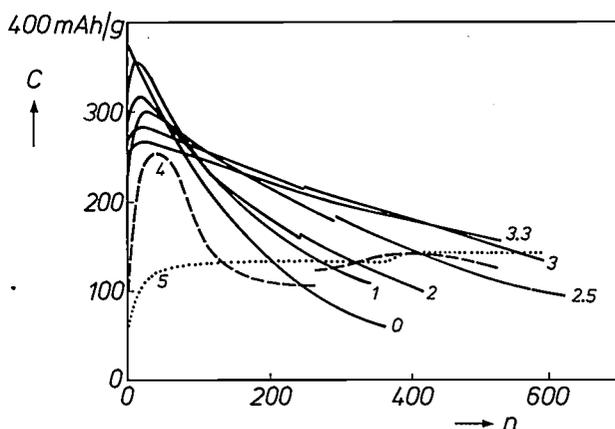


Fig. 13. The effect of substituting cobalt for nickel in LaNi_5 on the curve for the storage capacity C as a function of the number of charge/discharge cycles (n). The number of nickel atoms replaced by cobalt atoms is shown beside each curve. The optimum effect is found when 3.3 nickel atoms are replaced by cobalt.

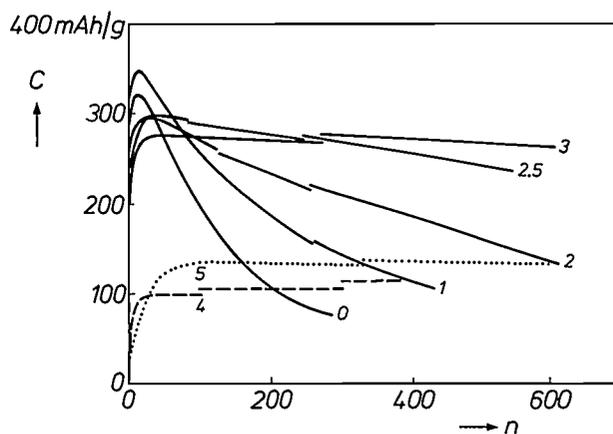


Fig. 14. As in fig. 13, but after the addition of 0.1 atomic fraction of Al. The addition of small amounts of Si has the same effect.

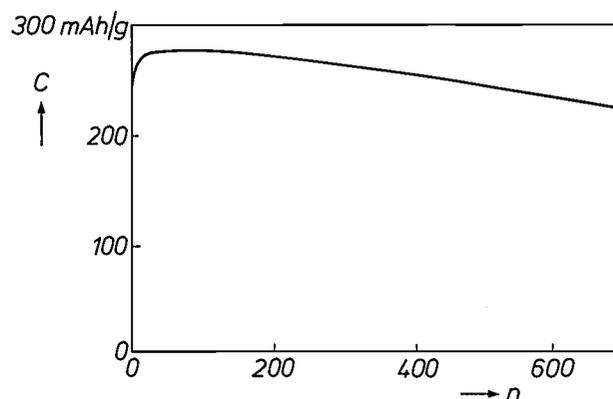


Fig. 15. The storage capacity C of an electrode of $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ as a function of the number of charge/discharge cycles n . The improvement in stability with respect to LaNi_5 is clear from a comparison with fig. 9.

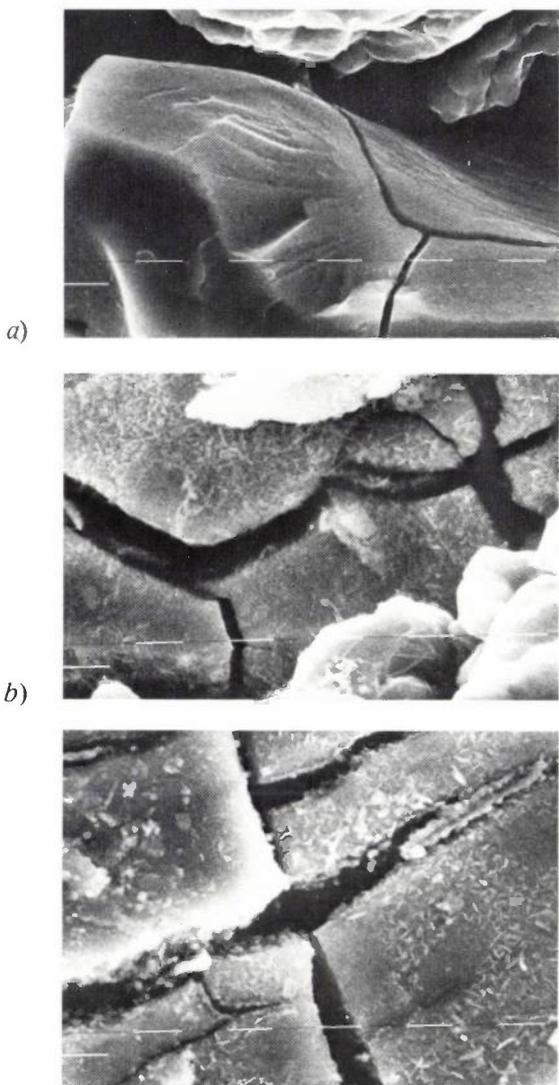


Fig. 16. SEM micrograph of $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ electrodes after zero (a), 41 (b) and 666 (c) charge/discharge cycles. One 'dash' represents 1 μm .

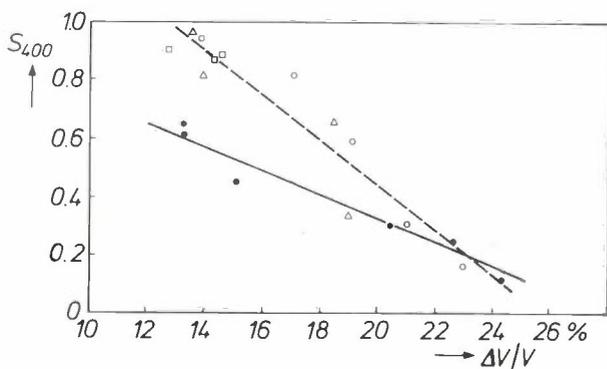


Fig. 17. Effect of the relative volume expansion $\Delta V/V$ on the stability of an electrode made from LaNi_5 and related compounds. The stability is expressed as the fraction S_{400} of the storage capacity remaining after 400 charge/discharge cycles. The filled circles relate to materials of the type $\text{LaNi}_{5-x}\text{Co}_x$, the open symbols to materials of the same type with small amounts of aluminium or silicon added. For small volume changes the materials with added aluminium or silicon are clearly the most stable.

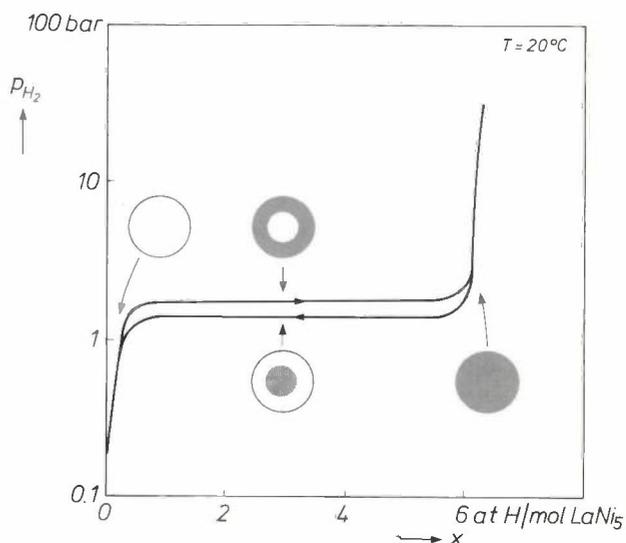
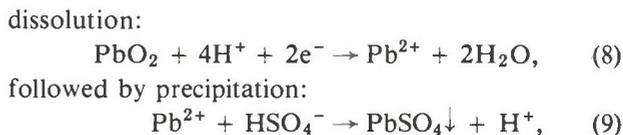


Fig. 18. Schematic representation of the phase boundary shift between the hydrogen-poor and hydrogen-rich phases in LaNi_5 grains on charging and discharging. The hydrogen-poor phase is white, the hydrogen-rich phase is shown grey.

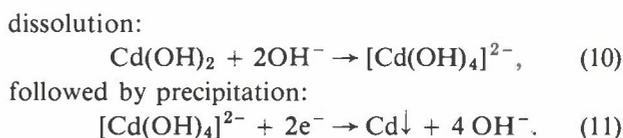
which acts as an additional restriction of the exposure of lanthanum to the electrolyte. This skin is thought to be composed of islands of cobalt and nickel with an aluminium- or silicon-containing oxide between them. Its inhibiting effect can only be maintained if the volume changes are not too great, if not the enclosing skin disintegrates; this explains why it is not found in LaNi_5 itself (fig. 14).

Is the hydride electrode a suitable alternative to the cadmium electrode of the NiCd cell?

Most electrodes of rechargeable batteries with an aqueous electrolyte are charged and discharged in a dissolution-precipitation mechanism. The discharge of the positive electrode in the lead-acid battery, for example, proceeds as follows:



and charging of the negative electrode of the NiCd cell:



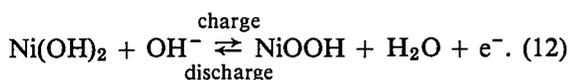
In the first case the soluble intermediate product is an ion, Pb^{2+} , and in the second case an ion complex, $[\text{Cd}(\text{OH})_4]^{2-}$. In view of the relatively low concentra-

tions of these ionic intermediates, large surface areas and short diffusion paths are required to obtain acceptable rates of charge and discharge. This makes it necessary to use electrodes with an appropriate porous structure. An undesirable feature of this charging and discharging mechanism is therefore that changes in shape occur during charging and discharging because of the redistribution of the active material within the porous electrode. The main processes are:

- nucleation and growth of the new solid phase at the surface of the active material still present;
- passivation, which consists of shielding the active material by a thin layer of non-active material formed around it, e.g. the sulphation of both electrodes of the lead-acid cell;
- recrystallization, leading to the formation of larger grains or more stable modifications, as found for instance during the ageing of the cadmium electrode in the NiCd battery, and finally
- dendrite formation, which may cause failure by internal short-circuiting.

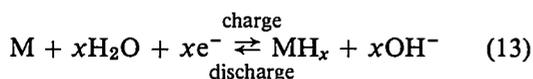
The worst effects of these processes are that they diminish the size of the active area and reduce its accessibility. These effects increase strongly with the rates of charge and discharge and with the 'depth of discharge'. Consequently, during average charge and discharge periods of 5 to 20 hours the percentage of optimally used active material is only between 25 and 45%. (This is why, for example, the capacity of the cadmium electrode in the nickel-cadmium cell has to be much larger than the nominal capacity of the nickel electrode.) On repeated discharging in less than an hour, which may be required in some applications, the percentage of active material optimally utilized is much lower than the stated percentage of 25%.

Electrodes without soluble intermediates, such as the nickel electrode, do not suffer from such limitations. Here charge and discharge take place via a solid-state transition, which may be written as:



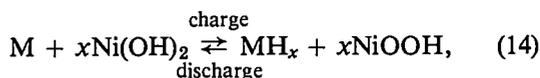
Owing to the absence of the degradation effects described, which are connected with the dissolution-precipitation mechanism, the nickel electrode has developed into one of the most valued products of present-day battery technology — it can be charged and discharged thousands of times at high current without noticeable deterioration.

As we have seen above, the hydride electrode is charged and discharged via a solid-state transition:



(where M could for example represent the compound $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.6}\text{Co}_{2.4}\text{Si}_{0.1}$).

When both types of electrode are combined as the positive and negative poles of a rechargeable nickel-hydride battery, the overall cell reaction for charging and discharging is:

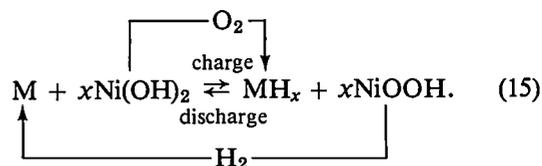


so that during operation of the cell there is no net consumption of electrolyte. This permits a simple and compact construction with little electrolyte^[9].

The resultant new type of rechargeable battery should have the following four desirable features:

- No cadmium or other toxic heavy metals.
- A high energy density, especially when used at high discharge rates. Both the NiCd cell and the nickel-hydride cell have a theoretical energy density of about 200 Wh/kg, but at the discharge rates in common use the actual energy density of the NiCd cell will not be more than 40 to 50 Wh/kg, partly because of the low degree of utilization of active material. In view of this considerable reduction in the theoretical energy density, the nickel-hydride cell, which has no dissolution-precipitation mechanism and consumes no electrolyte, seems to offer much better prospects.
- A high power density. Equation (14) shows that charging and discharging a nickel-hydride battery is rather like pumping hydrogen atoms, or protons, from one electrode to the other. Owing to the high mobility of the hydrogen atom and of the proton, the system should be capable of handling high current densities with low internal resistance, so that this electrode combination should result in a battery that can be charged and discharged at high rates.
- The possibility of effective protection from overcharging and overdischarging. In overcharging, the chemical bypass mechanism operates in the same way as in the NiCd cell, via an oxygen cycle. In overdischarging, however, the situation is different from that in a nickel-cadmium cell: the hydrogen gas evolved in the nickel-hydride cell is turned to good use, since it introduces a chemical bypass in the form of a hydrogen cycle:

overcharge



overdischarge

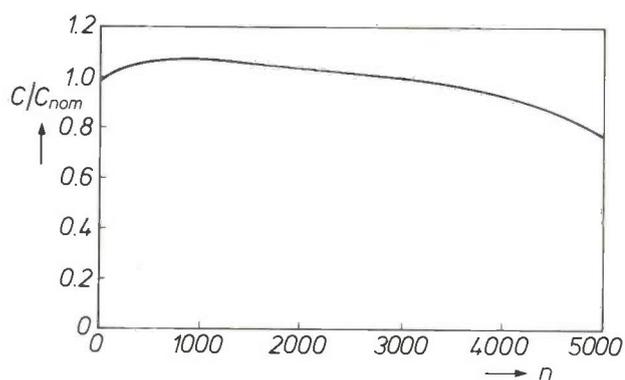


Fig. 19. The relative storage capacity C/C_{nom} of an experimental nickel-hydride cell as a function of the number of charge/discharge cycles n . The hydride electrode is oversized by a factor of 1.85 with respect to the nickel electrode. C_{nom} is 300 mAh.

Experimental nickel-hydride cell

In provisionally completing this experimental investigation, in which only half-cells have as yet been tested, we made a number of experimental nickel-hydride cells. At this point we were not trying to establish whether such a battery had all the good features that seemed likely in the previous section. This would have required the construction of an optimum battery that could provide the predicted high energy density and power density. Our aim was simply to show that predictions based on half-cell measurements were reasonably valid for the complete battery.

We were particularly interested in the validity of the following two hypotheses, which found support in the

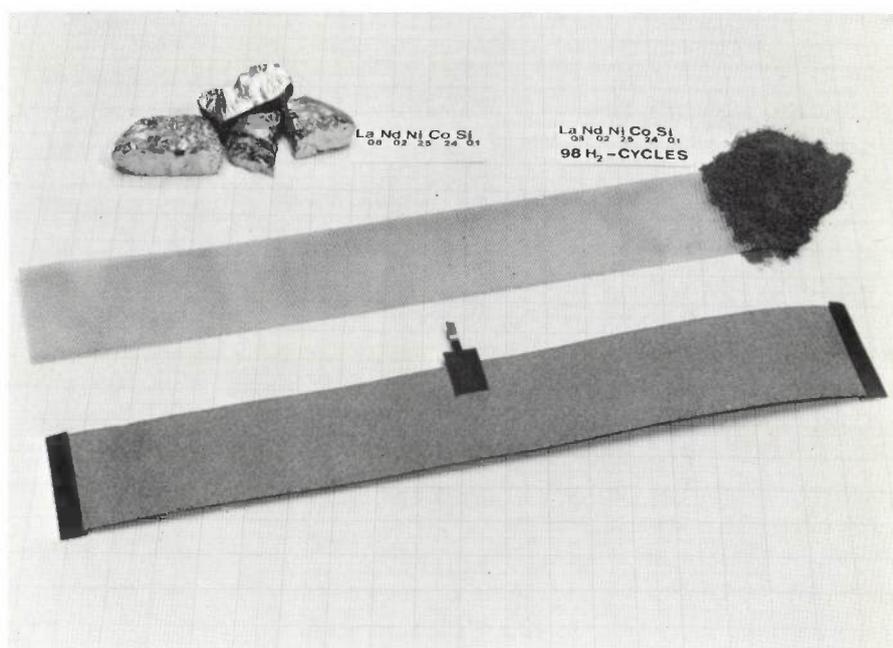


Fig. 20. One of our methods of manufacturing the hydride electrode. Metal castings of composition $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ (top left) are pulverized by repeated absorption and desorption of hydrogen gas (top right) and the resultant powder is applied as a paste to a strip of expanded nickel.

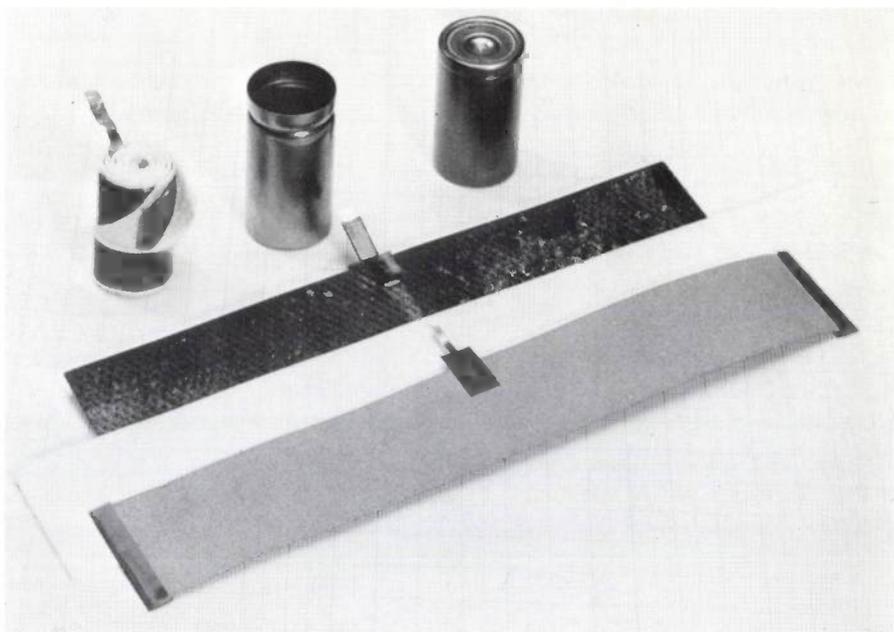


Fig. 21. Illustrating the second manufacturing step. A hydride electrode (grey; see the lower strip in fig. 20) is rolled up with a sintered nickel electrode (black) and a separator (white) to form a cylindrical battery. After addition of electrolyte the battery is hermetically sealed.



Fig. 22. A complete experimental nickel-hydride cell and pressure vessel.

investigation described here: firstly, that electrodes made from $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ would only give a 30% decay in storage capacity after 1000 charge/discharge cycles (fig. 15) and secondly, that no such decay of storage capacity would occur if the active material did not take part in the expansion and contraction process of hydrogen absorption and desorption.

To test these two hypotheses we made the hydride electrode, which we combined with a nickel electrode (overdimensioned by a factor of 1.85), so that we had active material that would expand and contract and material that would not do so.

Overdimensioning is also necessary in a practical battery (although usually by a much smaller factor), as part of the approach for protecting the cells of a hermetically sealed battery from overcharging and overdischarging. As soon as the overdimensioning ceases to exist because of the decay of storage capacity, a steep rise in gas accumulation will soon put an end to the life of the battery. (Battery life is defined here as the time for which the battery will operate before the gas pressure in the sealed unit rises above a specified value, e.g. 20 bars.)

From the measured decrease in the storage capacity of active material that participates 100% in hydrogen absorption and desorption (fig. 15) we can calculate the cell life for an overdimensioning of 1.85 times. For our experimental battery we calculated a life of 4000 charge/discharge cycles.

Fig. 19 shows the capacity decay of this battery as a function of the number of charge/discharge cycles. After 4000 cycles the storage capacity was still 92% of the nominal value, and even after 5000 cycles, that is to say after 60 weeks of continuous high-rate charging and discharging, the storage capacity was still better than 78% of the nominal value. After these 5000 cycles the pressure in the cell was about 10 bars.

The conclusion from these results is that the capacity decay of the new electrode material in a complete cell is at all events not greater than that found from the measurements on the half-cell.

Figs 20 and 21 illustrate one of the ways in which we made such experimental nickel-hydride cells. The intermetallic compound $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ was pulverized by repeated hydrogen-gas absorption/desorption cycles, and the powder was then mixed with fine nickel powder and applied as a paste to a strip of expanded nickel. The positive electrode was made by impregnating a porous sintered-nickel plate with $\text{Ni}(\text{OH})_2$. The electrolyte was a solution of 5 M KOH and 1 M LiOH. The two electrodes, separated by a polyamide strip, were rolled up to form a cylindrical package and placed in a nickel-plated steel can. After addition of the electrolyte the battery can be hermetically sealed or (as in our experiments) placed in a pressure vessel fitted with a pressure gauge (fig. 22).

Summary. LaNi_5 and related compounds have been investigated for their usefulness as electrode material in hydrogen-based rechargeable batteries. Although LaNi_5 and LaNi_4Cu give the required rapid absorption and desorption of large amounts of hydrogen, they are not suitable electrode materials because of the very large decay in their storage capacity with repeated charging and discharging. It was found that this capacity loss is due to the conversion of the active material into $\text{La}(\text{OH})_3$, which proceeds more rapidly as the volume changes in the active material caused by repeated charging and discharging become larger. Partial substitution of Co for Ni gives smaller volume changes and hence a more stable electrode material. The stability is further improved by adding small quantities of aluminium or silicon. The capacity of electrodes made from the newly developed $\text{La}_{0.8}\text{Nd}_{0.2}\text{Ni}_{2.5}\text{Co}_{2.4}\text{Si}_{0.1}$ was found to have decayed by only 30% after 1000 charge and discharge cycles. The experimental — hermetically sealed — nickel-hydride battery made with this material combines a long life with high energy density and power density, and effective protection from overcharging and overdischarging. It also contains no toxic heavy metals such as cadmium.

Dielectric resonators for microwave integrated oscillators

G. Lütteke and D. Hennings

Microwave oscillators used in radar and telecommunications must operate at constant frequency and generate little noise. Until recently this meant that oscillator circuits had to include a cavity resonator that functioned as a selective filter. The dimensions of the cavity resonator, however, were not compatible with the miniaturization of microwave integrated circuits in which waveguides and coaxial lines have been replaced by microstrip. Dielectric resonators are much more suitable for microwave integrated circuits, because they are small. Careful consideration must be given to the choice of the ceramic material for these resonators, to ensure high frequency stability and good noise characteristics.

Introduction

Now that direct television broadcasts via geostationary satellites have become a reality, there is a need for inexpensive circuits for microwave reception. Microwave integrated circuits (MICs) combine small dimensions and low weight with low cost. Modern MIC technology uses bipolar transistors and field-effect transistors with ever higher cut-off frequencies, and special active devices for the microwave region such as Gunn diodes and IMPATT diodes (IMPATT stands for 'IMPact Avalanche and Transit Time'). In MIC technology these components are interconnected by microstrip lines, which consist of a thin metal conductor (the strip) separated from a metal ground layer by a dielectric [1].

For telecommunication microwave links and for satellite television the frequency of the microwave carrier is required to be highly constant. A high signal-to-noise ratio is also essential. These requirements can only be met by the use of frequency-stabilizing selective band filters in the oscillator circuits. Until recently this meant using relatively large cavity resonators, which had to be made of Invar to ensure low temperature sensitivity. With these resonators a relative frequency stability of $10^{-6} \text{ }^\circ\text{C}^{-1}$ and Q-factors (quality factors) of more than 10 000 have been achieved.

Because of their size, cavity resonators are not very compatible with the structure of microwave integrated circuits. Efforts have therefore long been made to replace these resonators by ceramic types, which can be smaller because the permittivity of ceramic is higher than that of air. Many attempts to produce such resonators foundered because no ceramic materials could be found that combined a high permittivity and Q-factor with a low temperature dependence. It was not until the seventies that it was discovered that barium nonatitanate ($\text{Ba}_2\text{Ti}_9\text{O}_{20}$) was a suitable material for microwave dielectric resonators. This material was also found to be suitable for selective filters consisting of a number of coupled dielectric resonators. The dimensions of these filters can be kept within reasonable limits even at the lower frequencies. Its high permittivity also makes the material an excellent substrate material for microwave integrated circuits operated at relatively low frequencies. In addition, with specific microstrip dimensions low characteristic impedances can be achieved. (The characteristic impedance of a microstrip line is approximately inversely proportional to the square root of the permittivity of the dielectric.)

The temperature dependence of the frequency f_R , of electromagnetic standing waves in a dielectric reso-

Dr G. Lütteke is with Philips GmbH Apparatefabrik Krefeld, and was formerly with Philips GmbH Forschungslaboratorium Aachen (PFA), Aachen, West Germany. Dr D. Hennings is with PFA.

[1] J. H. C. van Heuven and A. G. van Nie, Microwave integrated circuits, Philips Tech. Rev. 32, 292-304, 1971.

nator is expressed by the temperature coefficient τ_{f_R} , which is defined as:

$$\tau_{f_R} = \frac{1}{f_R} \frac{df_R}{dt}, \quad (1)$$

where t is the temperature. This temperature coefficient is a function of the thermal expansion coefficient and of a similarly defined temperature coefficient for the permittivity. In oscillator circuits it is not desirable to make τ_{f_R} for the resonator exactly zero. This is because the overall temperature coefficient for the resonant frequency of the complete oscillator includes a small negative contribution from the other components in the circuit, and this has to be compensated by a positive contribution from the resonator.

To keep the dimensions of the resonator within reasonable limits the relative permittivity ϵ_r of the material must be high. (If we compare the dimensions of a cavity resonator and a dielectric resonator, we see that the cavity resonator at the same frequency and height/diameter ratio is about $\sqrt{\epsilon_r}$ times as large as the dielectric resonator [2].) In addition it must be possible to control the temperature coefficient of the resonant frequency in a small range by varying the composition of the material. The region of miscibility ratios of the two compounds BaTi_4O_9 and $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ in the phase diagram of the BaO-TiO_2 system is particularly suitable for controlling the temperature coefficient. With resonators of $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ containing about 30% of BaTi_4O_9 , oscillator circuits for 10 GHz have been made whose frequency stability is comparable with that of oscillators that have Invar cavity resonators. The dimensions of the dielectric resonators are very much smaller, however.

In this article we shall first look at the theoretical background of dielectric resonators in connection with measurements of material properties. We shall then deal with the composition of the ceramic materials for these resonators. After discussing various kinds of oscillator circuits, we shall conclude with a number of applications of stable oscillator circuits for microwave frequencies.

The dielectric resonator and the measurement of material properties

A dielectric resonator is usually a cylindrical body made of a material that has a high permittivity. Electromagnetic standing waves are set up in the resonator, and their wavelength is determined by the dimensions and the permittivity. Part of the high-frequency electromagnetic field extends beyond the resonator. This stray field can be used for coupling the resonator to the microwave circuit.

Ceramic materials with a high permittivity have been known for a long time. Barium titanate (BaTiO_3) can have ϵ_r -values as high as 10 000. Most ceramic materials, however, are unsuitable for use in resonators because their dielectric losses are excessively high at microwave frequencies. In ferroelectric materials like barium titanate these losses are due to the dielectric relaxation of the domain walls. The magnitude of the dielectric losses is determined by the loss factor $\tan \delta_m$, where δ_m is the loss angle of the material. The magnitude of the loss factor follows from the equation giving the complex permittivity ϵ :

$$\epsilon = \epsilon_r \epsilon_0 (1 - j \tan \delta_m), \quad (2)$$

where ϵ_0 is the permittivity of free space. In resonators the dielectric losses must be kept as small as possible. This is expressed in the Q-factor Q_m , which is defined as the reciprocal of the loss factor.

It is also important that the relative permittivity should not vary much as a function of temperature and that the thermal expansion coefficient of the material should be small. This can be seen from the expression for the temperature coefficient τ_{f_R} for the resonant frequency of dielectric resonators [3]:

$$\tau_{f_R} = -\frac{1}{2} (\tau_{\epsilon_r} + 2\tau_{\alpha}), \quad (3)$$

where τ_{ϵ_r} is the temperature coefficient of the permittivity and τ_{α} is the thermal expansion coefficient. A low value of τ_{f_R} can thus be obtained by using a material in which the quantities τ_{ϵ_r} and τ_{α} largely compensate one another. A problem is that for some materials there are no published values for these quantities, so that they have to be measured.

For these measurements we used the arrangement shown schematically in *fig. 1a* [4]. In a cylindrical sample microwave resonances are generated whose geometry is comparable with the geometry of resonances in a dielectric resonator. The sample is mounted between two metal plates. A variable-frequency signal from a network analyser is applied to the sample via a coaxial line and excites a field in it. The signal is returned to the analyser by a second coaxial line. The metal plates act as a waveguide operating below its cut-off frequency. This means that the high-frequency field emerging from the sample decreases rapidly with increasing radial distance. To reduce the effects of eddy-current losses in the metal plates, the sample is given a high length/diameter ratio.

The resonant modes in a cylindrical dielectric resonator are related to the modes in a cavity resonator. In a *cavity resonator* only transverse modes, designated TE_{mnp} and TM_{mnp} , can occur. The subscripts m , n and p are integers relating to the number of periods of the electric or magnetic field in the circum-

ferential, radial and axial directions respectively. Either the electric field (the TE_{mnp} mode, the subscripts relating to the magnetic field) or the magnetic field (the TM_{mnp} mode, the subscripts relating to the electric field) is a purely transverse field, that is to say with no axial component. The non-transverse magnetic or electric field does therefore have an axial component. In a cylindrical, non-metallized *dielectric resonator* only transverse modes with rotational symmetry ($m = 0$) can be excited. In dielectric resonators, how-

ever, there may also be non-transverse modes, which do not have rotational symmetry. In these modes, called HE_{mnp} or 'hybrid' modes, both the magnetic and electric fields have an axial component.

Of the many modes that can occur in a cylindrical resonator, only the transverse mode TE_{011} is of practical significance. In fig. 1b the electric and magnetic lines of force of the TE_{011} mode in the sample under test are shown schematically. The (transverse) electric lines of force are circles whose centres lie on the axis of the sample. The magnetic lines of force are approximately in the axial direction in the material and are rather like those of a magnetic dipole. Since each magnetic line of force lies in a plane through the axis and also extends beyond the material, the TE_{011} mode can easily be coupled to the network analyser by coaxial cables.

Fig. 1c shows a frequency spectrum measured with the arrangement of fig. 1a. Such a spectrum will only contain peaks due to the transverse modes TE_{0np} and TM_{0np} and the hybrid modes HE_{mnp} . It can be seen that the mode with the lowest frequency is the HE_{111} mode. Around the peak due to the TE_{011} mode in the figure an indication is given of how the Q of the material Q_m can be determined from the shape of the peak. Q_m is approximately equal to the ratio of the resonant frequency to the width Δf of the resonance peak, measured 3 dB below the peak. (A value of 3 dB corresponds to a power ratio of 0.5.) The permittivity can be calculated from the resonant frequency and the dimensions of the sample.

A dielectric resonator used in oscillator circuits has the shape of a flat cylinder. It is not clamped between metal plates, but is mounted with some axial and radial play in a metal case. The magnetic field can therefore extend outside the resonator in the axial direction; see fig. 2a. The mode indicated here is therefore called the $TE_{01\delta}$ mode. The coupling to the rest of the circuit is made with a microstrip line; see fig. 2b. To limit the losses in the ground plane of the microwave circuit the resonator is mounted at the bottom of the housing on a quartz-glass ring; see fig. 2c. The resonator has a height/diameter ratio of about 0.4, so that the differences in frequency between the $TE_{01\delta}$ mode and the neighbouring modes are relatively large. The side walls of the case are far enough away from the resonator to

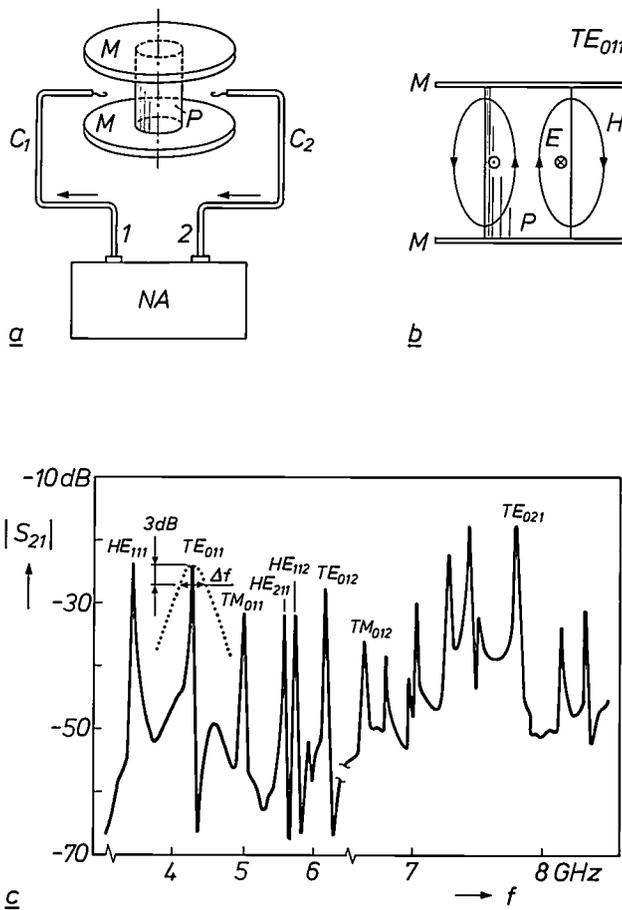


Fig. 1. Measurement of materials for dielectric resonators. a) Experimental arrangement. A cylindrical ceramic sample P is placed between two metal plates M . Microwave resonances are excited in the sample by a signal from a network analyser NA (Hewlett Packard 8410). The signal is supplied to the sample from port 1 of the analyser via a coaxial cable C_1 , and is returned via a second cable C_2 to port 2 of the analyser. Both coaxial cables are terminated in a loop in a plane perpendicular to the axis of the resonator. b) Diagram illustrating the TE_{011} mode in the sample. E electric line of force. H magnetic line of force. c) Example of a frequency spectrum measured with the network analyser. The modulus of the ratio S_{21} of the voltages at terminals 2 and 1 is plotted in dB as a function of frequency f . The spectrum is a combination of two measurements with different linear frequency scales. The modes corresponding to some of the peaks are indicated. The peak for the TE_{011} mode has been magnified in the horizontal direction: the dotted line. The Q -factor of the material can be determined from the width Δf of the peak at 3 dB below the peak.

- [2] M. W. Pospieszalski, On the theory and application of the dielectric post resonator, IEEE Trans. **MTT-25**, 228-231, 1977.
 [3] D. Hennings and P. Schnabel, Dielectric characterization of $Ba_2Ti_9O_{20}$ type ceramics at microwave frequencies, Philips J. Res. **38**, 295-311, 1983.
 [4] B. W. Hakki and P. D. Coleman, A dielectric resonator method of measuring inductive capacities in the millimeter range, IRE Trans. **MTT-8**, 402-410, 1960;
 W. E. Courtney, Analysis and evaluation of a method of measuring the complex permittivity and permeability of microwave insulators, IEEE Trans. **MTT-18**, 476-485, 1970.

minimize the excitation of eddy currents in them by the r.f. field. The case is nevertheless small, because the field decreases rapidly in the radial direction, since the upper and the lower plate also function here as a radial waveguide operating below its cut-off frequency.

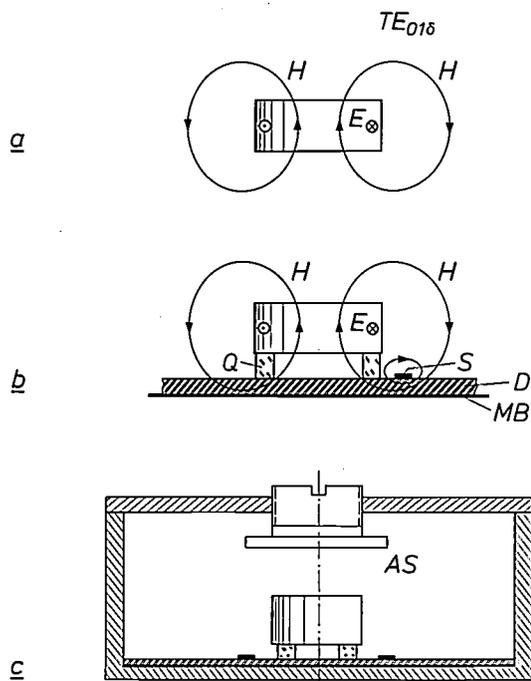


Fig. 2. a) Diagram of the $TE_{01\delta}$ mode in a dielectric resonator for oscillator circuits. (This mode should strictly be classified as $TE_{0\gamma\delta}$, since the magnetic field also extends from the resonator in the radial direction.) b) Coupling the resonator to a microstrip line in a microwave integrated circuit. Q quartz-glass ring. S strip of the microstrip line. D dielectric substrate. MB metal base. c) The dielectric resonator in its metal case, which contains the rest of the oscillator circuit (not shown). AS adjusting screw for fine tuning of the resonant frequency.

The resonant frequency is difficult to predict exactly from the dimensions of the resonator and the case. The frequency setting is therefore made with a tuning screw, which varies the distance over which the magnetic field of the $TE_{01\delta}$ mode extends in the axial direction. The resonant frequency can be set very accurately in this way.

Choice of material

It was discovered in the fifties that barium nonatitanate ($Ba_2Ti_9O_{20}$), because of its relatively high permittivity ϵ_r and very small temperature coefficient τ_{ϵ_r} , was a suitable material for stable ceramic capacitors [6]. Much later it turned out that the properties of this material also make it suitable for use in dielectric resonators for frequencies up to 10 GHz [6].

In the phase diagram of the quasi-binary system $BaO-TiO_2$, see fig. 3a, lines for the compound $BaTi_4O_9$ (80 mol. % TiO_2) and for pure TiO_2 lie on opposite sides of the line for the compound $Ba_2Ti_9O_{20}$ (81.818 mol. % TiO_2). These three materials all have a high permittivity (that of TiO_2 is very high) and a low loss factor; see Table I. The temperature coefficient τ_{f_R} of resonators made with TiO_2 and $BaTi_4O_9$ is relatively large, whereas it is almost zero for $Ba_2Ti_9O_{20}$. The compounds $BaTi_4O_9$, $Ba_2Ti_9O_{20}$ and TiO_2 are not miscible with one another, but they can be combined by sintering to form a dense type of ceramic with a porosity of less than 0.5%. It therefore seems sensible to obtain the small — but not zero — temperature coefficient required by making the two-phase ceramics $BaTi_4O_9 - Ba_2Ti_9O_{20}$ or $Ba_2Ti_9O_{20} - TiO_2$.

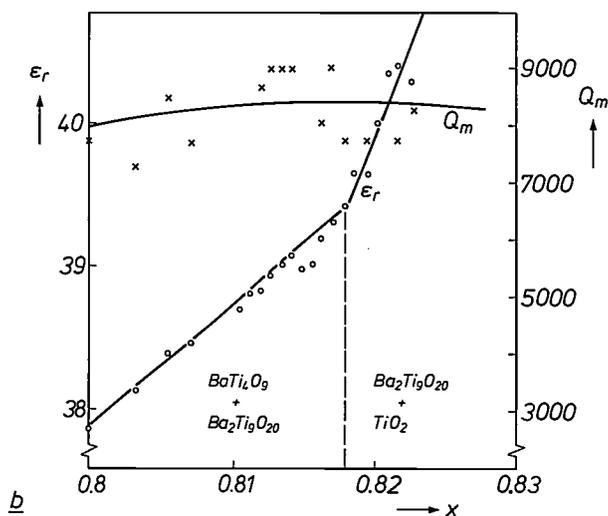
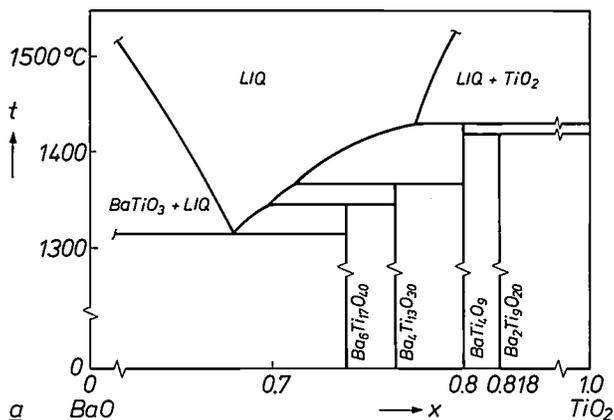


Fig. 3. Part of the phase diagram of the quasi-binary system $BaO-TiO_2$ [6]. t temperature. x molar quantity of TiO_2 . b) Relative permittivity ϵ_r and Q -factor Q_m as a function of x , as the result of a number of measurements with the arrangement of fig. 1a. The regions for the two-phase ceramics $BaTi_4O_9 - Ba_2Ti_9O_{20}$ and $Ba_2Ti_9O_{20} - TiO_2$ are separated by a dashed line.

Table I. Permittivity ϵ_r , temperature coefficient τ_{f_R} and loss factor $\tan \delta_m$, measured at 4 GHz for BaTi_4O_9 , $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ and TiO_2 samples.

	ϵ_r	τ_{f_R} [$^{\circ}\text{C}^{-1}$]	$\tan \delta_m$
BaTi_4O_9	37.5	20×10^{-6}	1.6×10^{-4}
$\text{Ba}_2\text{Ti}_9\text{O}_{20}$	39.3	≈ 0	1.1×10^{-4}
TiO_2	100	1550×10^{-6}	2×10^{-4}

Fig. 3b shows that samples made from such two-phase ceramics have a permittivity that is a monotonically increasing function of the molecular content x of TiO_2 . The Q-factor $Q_m = 1/\tan \delta_m$ is seen to be high and to change relatively little as a function of x . Fig. 4 shows the temperature coefficient τ_{f_R} as a function of temperature for samples of the two-phase ceramic $(\text{Ba}_2\text{Ti}_9\text{O}_{20})_{1-y}-(\text{BaTi}_4\text{O}_9)_y$, with differing mixture ratios y . For pure $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ the value of τ_{f_R} at 0°C happens to be exactly equal to zero. This material might therefore be used for making resonators that have a broad minimum at 0°C in the curve for the resonant frequency as a function of temperature, which approximates to a parabola. However, we want

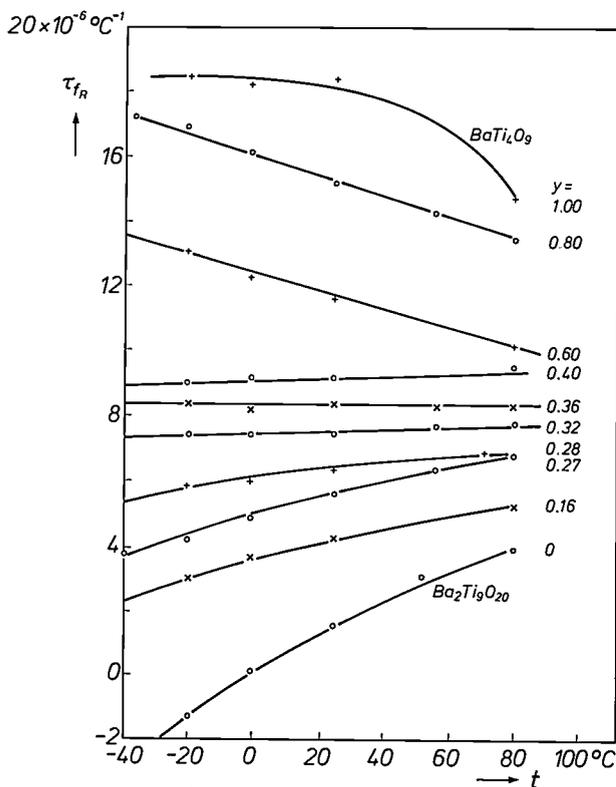


Fig. 4. Temperature coefficient τ_{f_R} for the resonant frequency f_R of ceramic samples as a function of temperature t . Each curve corresponds to a different value of the molar content y of BaTi_4O_9 in the two-phase ceramic $(\text{Ba}_2\text{Ti}_9\text{O}_{20})_{1-y}-(\text{BaTi}_4\text{O}_9)_y$.

the temperature coefficient of the resonant frequency to be independent of temperature and to have a small positive value. In this way we can compensate for negative temperature coefficients of other elements in the oscillator circuit.

Fig. 4 shows that this condition is satisfied by ceramic samples of $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ containing about 30% of BaTi_4O_9 . Resonators made of this material have been used in oscillator circuits that gave a frequency drift of only 150 kHz at 10 GHz in a temperature range from -20 to 100°C [7]. The product of Q_m and the resonant frequency f_R (in GHz) of ceramic resonator materials is generally found to be practically constant in the range from 2 to 15 GHz, and is therefore used as a figure of merit. With the new materials described here we have reached values of about 45 000 for this product.

The oscillator circuit

Coupling to the resonator

There are various types of microwave oscillator circuit. If a dielectric resonator is used to stabilize the resonant frequency, it is always connected to the active element in the circuit by a microstrip line, the 'coupling line'. The coupling line is terminated by a load impedance, usually a matched load — a resistance equal to the characteristic impedance Z_0 of the line, so that there is no reflection at the termination. The terminating impedance is connected between the strip of the coupling line and the ground plane of the microwave integrated circuit; see fig. 5a.

The resonator may be regarded as a tuned circuit consisting of a resistance, an inductance and a capacitance connected in parallel, as in fig. 5b. Unloaded, the circuit has a Q-factor Q_R and an angular resonant frequency ω_R , and is considered to be coupled to the rest of the circuit by an ideal transformer with a turns ratio of n . If the tuned circuit plus transformer is replaced by another equivalent tuned circuit connected directly into the circuit, the values R , L and C change to R_1 , L_1 and C_1 (see fig. 5c), and we can write:

$$R_1 = n^2 R.$$

The coupling factor β is defined as the ratio of this

[5] G. H. Jonker and W. Kwestroo, The ternary systems $\text{BaO-TiO}_2\text{-SnO}_2$ and $\text{BaO-TiO}_2\text{-ZrO}_2$, J. Am. Ceram. Soc. 41, 390-394, 1958.

[6] H. M. O'Bryan, J. Thomson and J. K. Plourde, A new BaO-TiO_2 compound with temperature-stable high permittivity and low microwave loss, J. Am. Ceram. Soc. 57, 450-453, 1974.

[7] C. Tsironis and D. Hennings, Highly stable FET DROs using new linear dielectric resonator material, Electron. Lett. 19, 741-743, 1983.

resistance to the characteristic impedance of the coupling line:

$$\beta = \frac{R_1}{Z_0} = \frac{n^2 R}{R_0}, \quad (4)$$

where R_0 is the load impedance. The magnitude of the coupling factor is largely determined by the distance between resonator and strip.

The impedance Z_1 is defined as the impedance between the strip and the ground plane at the position of the resonator and looking towards the load. This impedance is a function of the angular frequency ω of the microwave signals and can be expressed in terms of the quantities defined above:

$$Z_1(\omega) = Z_0 \left\{ 1 + \frac{\beta}{1 + jQ_R(\omega/\omega_R - \omega_R/\omega)} \right\}. \quad (5)$$

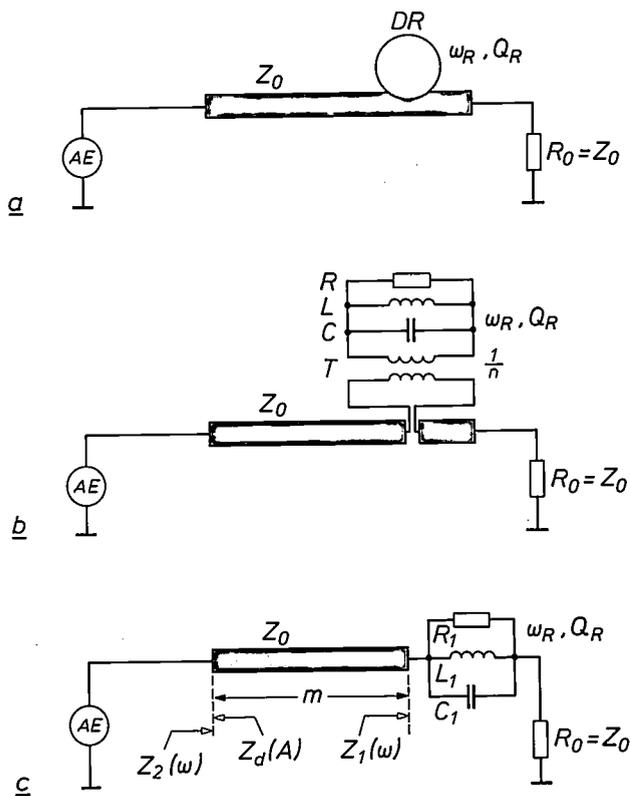


Fig. 5. Coupling of the dielectric resonator DR to the active element AE . *a*) Diagram showing the coupling line, a microstrip line of characteristic impedance Z_0 . The coupling line is terminated by a resistance R_0 , equal to Z_0 . ω_R, Q_R angular resonant frequency and Q -factor of the resonator. *b*) Circuit with the resonator replaced by a parallel tuned circuit of R, L and C . The tuned circuit is connected to the coupling line by a transformer T with turns ratio n . *c*) Circuit with the tuned circuit and transformer replaced by a tuned circuit of R_1, L_1 and C_1 . m length of the coupling line. $Z_1(\omega)$ the impedance between strip and ground plane near the resonator, looking towards R_0 ; ω angular frequency. $Z_2(\omega)$ as $Z_1(\omega)$, but at the input to the coupling line. $Z_d(A)$ input impedance of the active element; A amplitude of the oscillation.

Any non-variable impedance can be represented as a point on a Smith chart [8]; $Z_1(\omega)$ represents a circle on the chart (see fig. 6).

A Smith chart can be used to represent both the vector reflection coefficient, i.e. the voltage ratio of the incident and reflected microwave signals, and the real and imaginary components of the impedance normalized to the characteristic impedance. The voltage standing-wave ratio or VSWR can be read from these quantities in the diagram. This is the ratio of the voltages at the nodes and antinodes in the standing wave resulting from the incident and reflected signals.

The circle $Z_1(\omega)$ goes through M , the centre of the Smith chart, when $\omega = 0$ and $\omega = \infty$, because Z_1 is equal to the characteristic impedance Z_0 at these values of ω . The point for $\omega = \omega_R$, i.e. for resonance of the tuned circuit in fig. 5c, corresponds to the right point of intersection of the circle $Z_1(\omega)$ with the line h in the chart, since then $Z_1 = Z_0 + R_1$, from eq. (5). (Line h corresponds to impedances that are purely resistive.) The diameter of the circle is equal to $R_1 = \beta Z_0$; this diameter therefore increases as the resonator is more closely coupled to the microstrip line.

It is a property of the Smith chart that a point on the chart moves along a circle of centre M when the position where the impedance is being considered moves along a transmission line (in our case the microstrip line). If we want the impedance $Z_2(\omega)$ that the active element 'sees' at the input to the coupling line — the impedance at that point due to the coupling line, resonator and load — we find it by rotating the circle $Z_1(\omega)$ clockwise through an angle θ with respect to M . This angle is given by

$$\theta = \frac{2m}{\lambda} 360^\circ,$$

where m is the length of the coupling line and λ is the wavelength. The locus for the load of the active element corresponds almost exactly to the circle $Z_2(\omega)$ thus obtained.

The oscillation condition for the complete circuit can be expressed as

$$Z_d(A) + Z_2(\omega) = 0, \quad (6)$$

where the input impedance Z_d of the active element is a function of the amplitude A of the oscillation.

Oscillation can only occur if the complex input impedance of the active element contains a negative real component. (This is the case in Gunn and IMPATT diodes when the supply voltage satisfies certain conditions; in transistors the oscillation condition can only be satisfied by using negative feedback.) Within the limited frequency range determined by the $TE_{01\delta}$

mode of the resonator, only the real component of Z_d changes with varying amplitude; the reactance is virtually constant. On the Smith chart $-Z_d(A)$ therefore corresponds well to a line of constant reactance: an arc of a circle with its centre-point on the line v . The operating point of the circuit is the point of intersection of the lines for $Z_2(\omega)$ and $-Z_d(A)$.

The frequency and amplitude of the oscillation can be adjusted by moving the operating point. The location of the line $-Z_d(A)$ depends on the supply voltage of the active element and — in the case of transistors — on the coupling between the three terminals. The location of the locus $Z_2(\omega)$ depends on the length of the microstrip line between the active element and the resonator. The diameter of $Z_2(\omega)$ is determined by the coupling factor. The resonant frequency ω_R of the actual resonator depends on its dimensions and the setting of the tuning screw. An appropriate choice for the characteristic impedance — which determines the

scale of the Smith chart — will ensure that point B , which represents the point $-Z_d(A = 0)$, lies inside the circular locus. It is easy to see that the operating point will then always be in a stable situation.

Terminating the coupling line with a pure resistance equal to the characteristic impedance of the line has the advantage that components in the signal whose frequency differs from ω_R are attenuated. The resonator is then 'electrically inoperative' for components at such frequencies. It can be shown that terminating the coupling line with a reactance does not give one single stable operating point, but three different operating points. However, one of these is in an unstable situation. The other two, although stable, correspond to frequencies that are very close together, so that the ultimate frequency of the oscillation is not predictable. Terminating the microstrip line with a reactance does not therefore produce a circuit of any practical use.

The various types of oscillator circuit

Fig. 7 shows five oscillator circuits, each containing a three-terminal active element [9]. The active element can be a bipolar transistor or a field-effect transistor. The microstrip lines in the various circuits are terminated by their characteristic impedance Z_0 . The resistance that acts as a load in the circuits shown should be considered as equivalent to the actual transmitting or receiving circuits. The main difference between the microwave circuits is in the position of the load resistance. All these oscillator circuits were originally used with cavity resonators.

Fig. 7a shows the simplest type of oscillator circuit. The terminating resistance of the microstrip line acts not only as a damping resistance (R_D) but also as the load (R_L). The dielectric resonator is used as a 'reaction-type' resonator.

Fig. 7b shows a circuit with a 'transmission-type' resonator. The load is connected to the resonator by a second microstrip line. The resonator is not so tightly coupled to this line as to the first line. The signal goes from the active element via the resonator to the load, so that components at frequencies outside the pass-band of the resonator are strongly attenuated.

In the circuit of fig. 7c the damping resistance and the load are connected to different terminals of the transistor. There is a matching network for the load resistance. The resonator is coupled to the microstrip line in such a way that the maximum energy is reflected

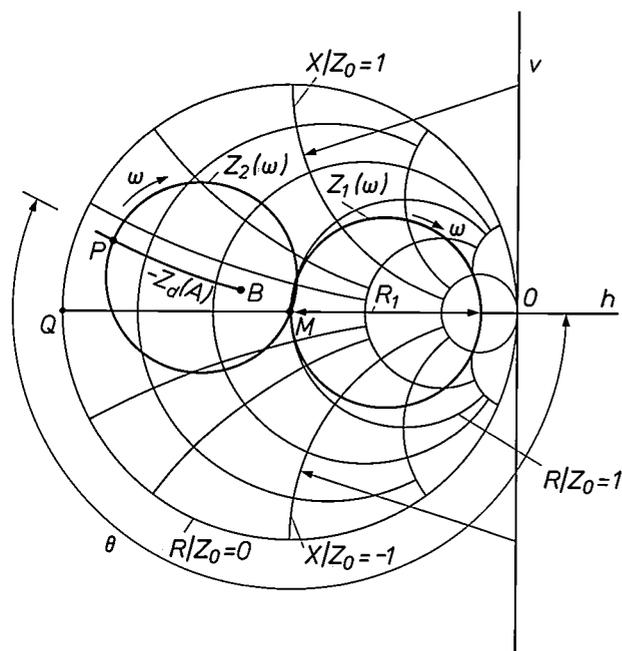


Fig. 6. Smith chart [8] for the coupling of the resonator to the active element. Every point inside the large circle represents an impedance $Z = R + jX$, normalized to the characteristic impedance Z_0 . Circles through O with their centres on line v correspond to a constant ratio X/Z_0 . X is positive above line h , and negative below it. Circles through O with their centres on line h correspond to a constant ratio R/Z_0 . For O : $R = \infty$, $X = \infty$. For M : $R = Z_0$, $X = 0$. For Q : $R = 0$, $X = 0$. The thicker circles represent the impedances $Z_1(\omega)$ and $Z_2(\omega)$; see fig. 5. The circle $Z_1(\omega)$ intersects h at M for $\omega = 0$ and $\omega = \infty$. The second intersection of the circle with h corresponds to $Z_1(\omega_R) = Z_0 + R_1$. θ angle through which the circle for $Z_1(\omega)$ has to be rotated around M to obtain $Z_2(\omega)$; for $\theta = 360^\circ$, m is equal to a half-wavelength. The line $-Z_d(A)$, see also fig. 5, is seen to correspond to a line of constant reactance. The point B corresponds to $-Z_d(A = 0)$ and lies within $Z_2(\omega)$. P , the intersection point of the lines $-Z_d(A)$ and $Z_2(\omega)$, is the operating point, which in this case has a stable position.

[8] P. H. Smith, Transmission line calculator, Electronics 12, No. 1 (January), 29-31, 1939;
L. V. Blake, Transmission lines and waveguides, Wiley, New York 1969.

[9] I. M. Clarke and B. F. van der Heijden, Dielectric resonator oscillators, Proc. Military Microwaves '84, London 1984, pp. 591-595.

at the resonant frequency. A standing wave is then set up in the microstrip line between the active element and the resonator. This circuit is therefore referred to as an oscillator with a 'reflection-type' resonator.

In the circuit of fig. 7d two of the terminals of the active element are connected to a microstrip line coupled to the resonator. These lines with the resona-

tor form a feedback path, and the device is therefore referred to as an oscillator with a 'feedback-type' resonator. The load forms part of the impedances Z_{31} or Z_{23} , which are connected to the third terminal of the transistor.

All these circuits have one resonator. However, two resonators can be used, which will give even better frequency stability; see fig. 7e. This circuit may be regarded as a combination of a reaction-type and a reflection-type resonator. The line with the reaction resonator is terminated by the load, the other by the damping resistance. A problem with this oscillator circuit is the mechanical tuning of the two resonators, since it only oscillates when both resonators are accurately tuned to the same frequency.

For completeness fig. 8 shows the ways in which the base, collector and emitter of a bipolar transistor, or the source, drain and gate of a field-effect transistor can be connected to terminals 1, 2 and 3 of the circuits in fig. 7. Fig. 8a relates to oscillators with the reaction-type or transmission-type resonator, fig. 8b to the oscillator with reflection-type resonator, and fig. 8c to the oscillator with feedback-type resonator. The first three circuits in fig. 7 also operate with Gunn or IMPATT diodes. Terminal 1 and the impedances Z_{12} and Z_{31} should then be omitted in fig. 7a and b. In fig. 7c terminal 1 should be omitted and the impedances Z_{12} and Z_{31} combined.

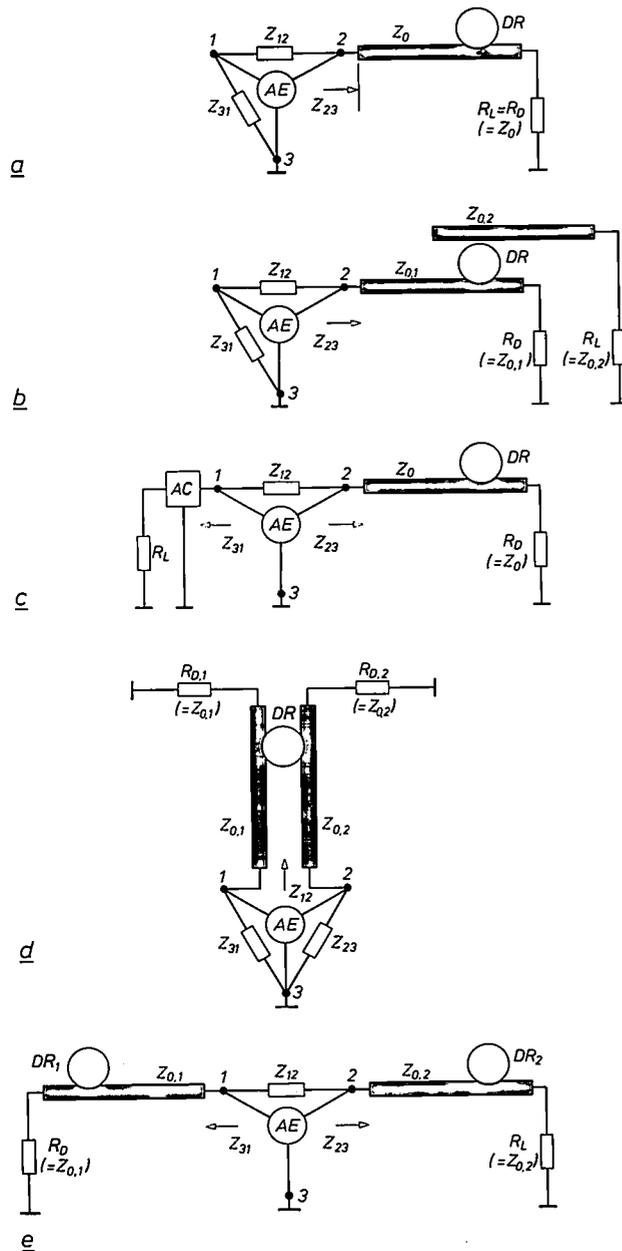


Fig. 7. The different oscillator circuits. a) Oscillator with reaction-type resonator. 1, 2, 3 terminals of the active element. Z_{12} , Z_{31} , Z_{23} impedances between the terminals. R_L load. R_D damping resistance. R_L and R_D are combined here to form a single resistance equal to the characteristic impedance Z_0 of the coupling line. b) Oscillator with transmission-type resonator. $Z_{0,1}$, $Z_{0,2}$ characteristic impedances of the coupling lines. c) Oscillator with reflection-type resonator. AC matching network for the load. d) Oscillator with feedback resonator. $R_{D,1}$, $R_{D,2}$ damping resistors in the coupling lines. e) Oscillator with two resonators: DR_1 and DR_2 .

Frequency stability

The resonant frequency of the oscillator circuit varies slightly in practice. The variations may be divided into long-term and short-term variations. The long-term variations are a consequence of temperature fluctuations; their magnitude depends on the temperature coefficient τ_{f_0} for the oscillator resonant frequency f_0 (defined in a similar way to τ_{f_R} , see eq. 1). The short-term variations are random and are referred to as phase noise.

The temperature coefficient τ_{f_0} depends on the thermal behaviour of the active element, the dielectric resonator, the microstrip line and the case in which the microwave circuit is mounted. With a suitable design, the effect of the microstrip line and the case on τ_{f_0} can be made small compared with the effect of other elements. τ_{f_0} then depends almost entirely on the thermal characteristics of the active element and the resonator.

The temperature coefficient of the oscillator circuit can be expressed as ^[10]:

$$\tau_{f_0} = \frac{F(\beta)}{Q_R} \frac{d\phi_d}{dt} + \tau_{f_R}, \quad (7)$$

where $F(\beta)$ is a function of the coupling factor and ϕ_d

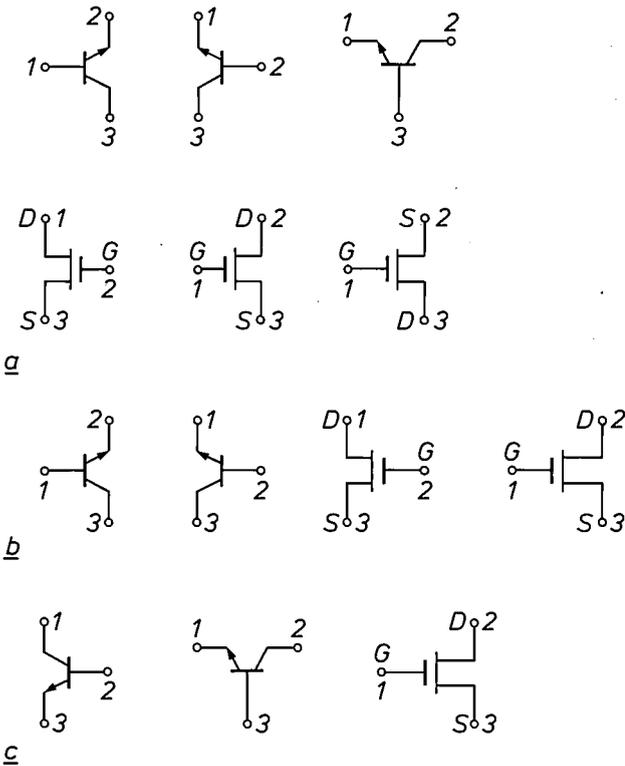


Fig. 8. The various ways in which a bipolar transistor or a field-effect transistor can be incorporated in the circuits of fig. 7. The numbering of the terminals corresponds to that in fig. 7. a) The active elements for the oscillators of fig. 7a and b. b) The active elements for the oscillator of fig. 7c. c) The active elements for the oscillator of fig. 7d.

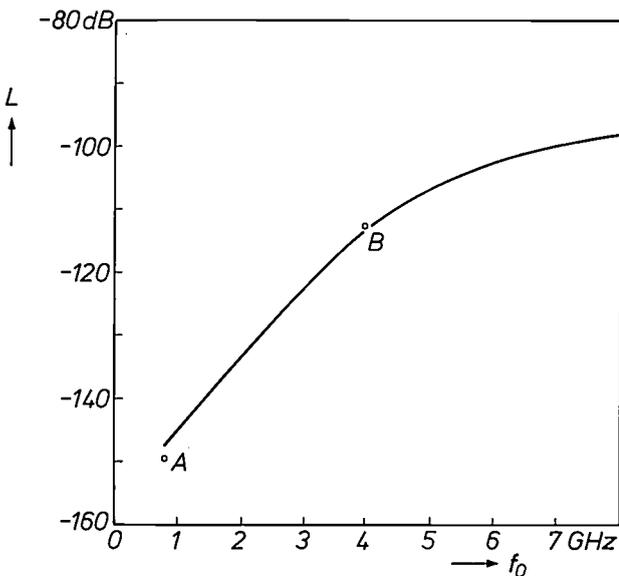


Fig. 9. Power density of the phase noise of microwave oscillators with dielectric resonators as a function of the oscillation frequency f_0 . The curve is the result of an extrapolation using eq. (9) and literature data [13] (A) together with results of our own measurements (B). Plotted along the vertical axis is the ratio L in dB of the power density in W/Hz of the (single-sideband) power density spectrum of the phase noise at a distance of 10 kHz from the resonant frequency, and the output power of the oscillator.

is the argument associated with the input impedance Z_d of the active element (see fig. 5c). The first term in (7) is negative for all active (semiconductor) elements and has a value in the range from -2×10^{-6} to $-10 \times 10^{-6} \text{ }^\circ\text{C}^{-1}$. Since $|\tau_{f_0}|$ should be as small as possible, τ_{f_R} should have a low positive value. The correct value of τ_{f_R} can be achieved, as we have seen, by using a mixed ceramic with a particular ratio of $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ to BaTi_4O_9 for the resonator material, or by building a sandwich resonator of two different materials [11].

The (single-sideband) power density of the phase noise (i.e. the noise power per unit bandwidth) is generally related to the output power P_O of the oscillator. This relative power density $L(f - f_0)$ of the phase noise at a 'distance' $f - f_0$ from the oscillator frequency is proportional to f_0^2 and the absolute temperature, and inversely proportional to the square of the Q-factor of the oscillator Q_O [12]. If we compare two spectra of the phase noise at different frequencies f_0 and f_0' , different output powers P_O and P_O' and different Q-factors Q_O and Q_O' at the same temperature, we obtain the following equation for the ratio of the relative power densities:

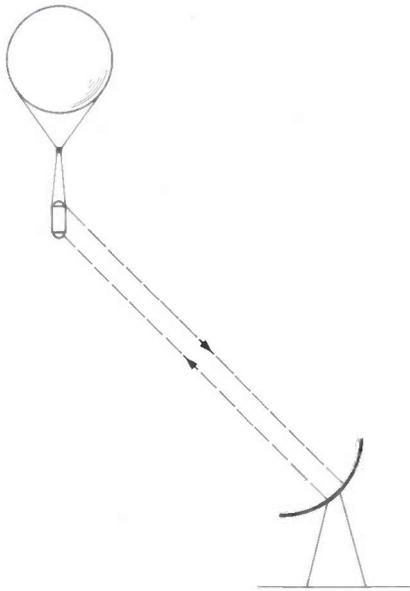
$$\frac{L(f - f_0)}{L'(f - f_0)} = \frac{f_0^2 Q_O'^2 P_O'}{f_0'^2 Q_O^2 P_O} \quad (8)$$

If in addition the coupling factors of the resonators are identical, the Q-factors of the oscillators in eq. (8) can be replaced by those of the resonators. It is also known that the maximum available output power of microwave oscillators is inversely proportional to the square of the oscillator frequency. Equation (8) can therefore be written as:

$$\frac{L(f - f_0)}{L'(f - f_0)} = \frac{f_0^4 Q_R'^2}{f_0'^4 Q_R^2} \quad (9)$$

This equation gives an idea of the phase noise to be expected at high frequencies. It has been used in fig. 9 together with data from the literature [13] (point A) and data from our own measurements of relative power densities of the phase noise at 4 GHz (point B). The Q-factors Q_R are comparable, because resonators of pure $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ were used in both cases. Fig. 9 and

[10] T. Makino and A. Hashima, A highly stabilized MIC Gunn oscillator using a dielectric resonator, IEEE Trans. MTT-27, 633-638, 1979;
C. Tsironis and V. Pauker, Temperature stabilization of GaAs MESFET oscillators using dielectric resonators, IEEE Trans. MTT-31, 312-314, 1983.
[11] C. Tsironis and P. Lesartre, Temperature stabilization of GaAs FET oscillators using dielectric resonators, Proc. 12th Eur. Microwave Conf., Helsinki 1982, pp. 181-186.
[12] K. Kurokawa, Injection locking of microwave solid-state oscillators, Proc. IEEE 61, 1386-1410, 1973.
[13] G. D. Alley and H.-C. Wang, An ultra-low noise microwave synthesizer, IEEE Trans. MTT-27, 969-974, 1979.



a



b

Fig. 10. *a)* Principle of locating the position of a radiosonde by means of a transponder, a combination of a transmitter and a receiver. The transponder, which is suspended from the balloon of the radiosonde, responds to a signal transmitted by a parabolic antenna on the ground. The position of the radiosonde can be determined from the angles at which the ground station receives a signal from the transponder, and from the time that elapses between the transmission and reception of the signal. *b)* The transponder, complete with antennas, which are located underneath the top and bottom shielding caps. The transmitted power is 200 mW and is sufficient for a range of 50 km. The transponder contains a microwave oscillator with dielectric resonator and is a product of the Philips GmbH supply centre Systeme und Sondertechnik in Bremen.

eq. (9) show that at high frequencies the phase noise increases steeply. It has also been found that bipolar transistors are very suitable for use in oscillator circuits with low phase noise.

Applications

Dielectric resonators will mainly be used to replace the cavity resonators now used for stabilizing microwave oscillators — for example in microwave receivers for satellite television^[14]. In addition, various applications have emerged that would be inconceivable without dielectric resonators. Two of them will be described below.

An oscillator for a transponder in a radiosonde

A transponder is a combination of a transmitter and a receiver. Nowadays a radiosonde can carry a microwave transponder for locating its position; see *fig. 10*. The transponder must be small and light, and also inexpensive — since the transponder is not usually recovered after the flight.

The transponder responds to the signal from a transmitter on the ground. If the frequency and an identifying modulation of the signal correspond to preset input data, the transponder replies by sending out a signal at the same frequency with the same modulation. The position of the radiosonde can be determined from the time between the transmission of the search signal and the reception of the reply signal and from the angles from which the reply signal is received.

The signal in the transponder, which has a frequency of 9.2 GHz, is obtained by frequency-doubling and then amplifying the 4.6-GHz output signal of an oscillator. The signal is modulated by switching the oscillator on and off in a fixed pattern. Three different circuits have been developed for the oscillator, with reaction, transmission or reflection-type resonators; see *fig. 7a-c*. The active element of the oscillator circuits is a bipolar transistor, which is connected as shown in the second diagram in *fig. 8a* or *b*. The circuits are designed such that terminals 1 and 2 are not connected ($Z_{12} = \infty$). In the circuits with a reaction-type or transmission-type resonator (*fig. 7a* and *b*) Z_{31} is a purely capacitive feedback impedance. In the circuit with the reflection-type resonator (*fig. 7c*) the matching network *AC* transforms the load resistance of 50 ohms into a high resistance between collector and emitter. This circuit also has a capacitive feedback impedance between terminals 1 and 3.

The value of the feedback impedance must be set very accurately, or the oscillator circuit will resonate in the $HE_{11\delta}$ or $TM_{01\delta}$ mode, whose resonances are at

frequencies close to that of the required $TE_{01\delta}$ mode; see fig. 1c. It is very difficult to set the frequency exactly, because the input impedance Z_d of the active element changes as a result of transients when the oscillator is switched on and off.

The three different oscillator circuits are illustrated in fig. 11. The type of circuit can be recognized from the shape of the lines on the ceramic substrate. The dielectric resonator of pure $Ba_2Ti_9O_{20}$ is mounted directly on the substrate and has a diameter of 14 mm

and a height of 5 mm. The height of the aluminium case is therefore no more than 14 mm; the length and width are both 50 mm. The complete oscillator weighs as little as 80 g. The oscillation frequency can be tuned through a band of ± 30 MHz by means of a screw in the cap.

The transistors used in the circuits are leadless chips designed for microwave frequencies. At a supply voltage of 9 V and an input current of 32 mA, the output power of the oscillators is 30 mW. Fig. 12a gives a plot of frequency as a function of temperature for the transmission-type resonator. The curve corresponds to a temperature coefficient of $-5.1 \times 10^{-6} \text{ }^\circ\text{C}^{-1}$. If the resonators had been made of $Ba_2Ti_9O_{20}$ with about 30% of $BaTi_4O_9$, the absolute value of the tem-

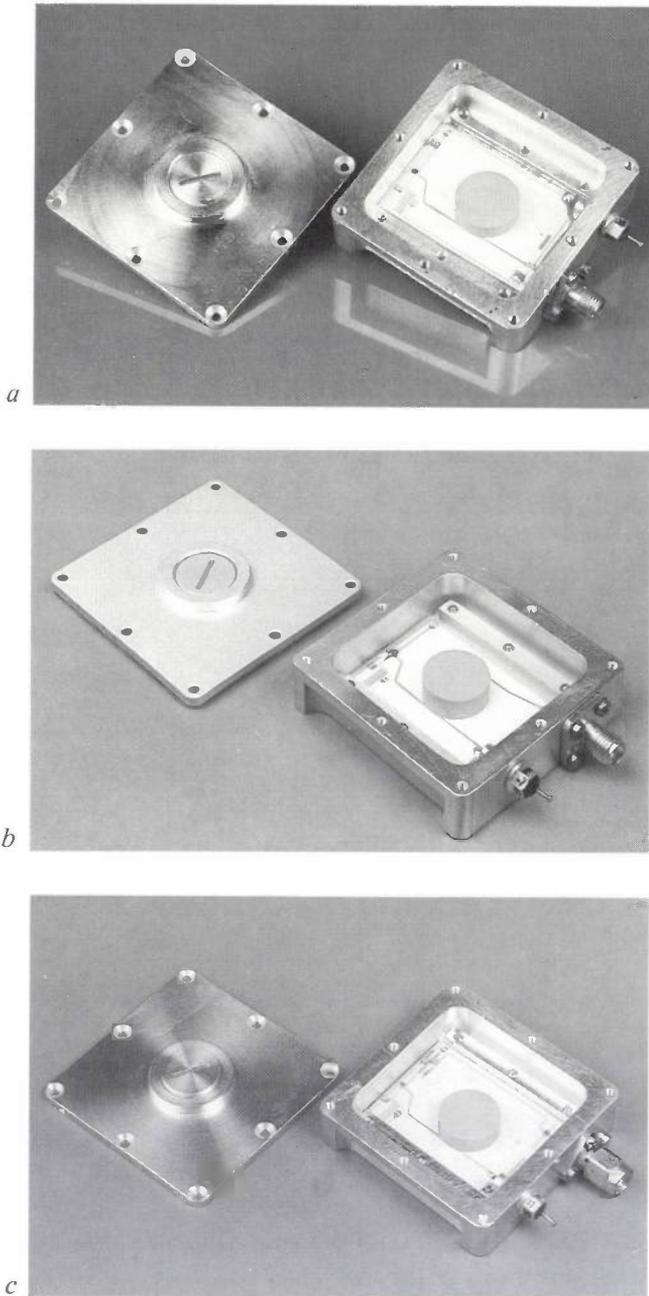


Fig. 11. The various oscillators for the transponder in fig. 10 that have been built and tested. a) Oscillator with reaction-type resonator. b) Oscillator with transmission-type resonator. c) Oscillator with reflection-type resonator. The tuning screw for fine tuning of the oscillator frequency is clearly visible in the cap of the housing of all three oscillators; see fig. 2c.

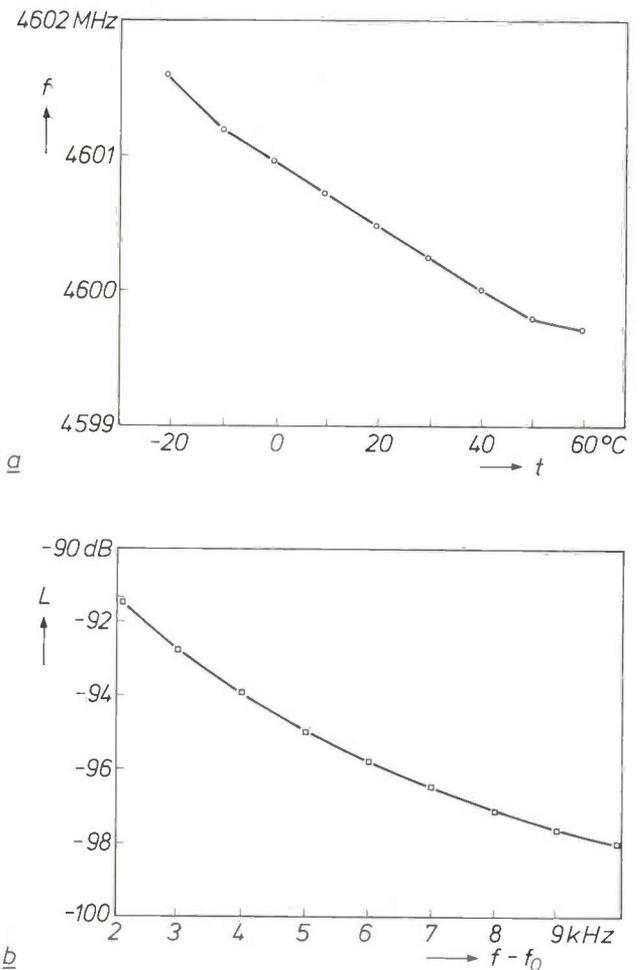


Fig. 12. a) Plot of the frequency f as a function of temperature t for the oscillator with transmission-type resonator in fig. 11b. (The corresponding temperature coefficient of $-5.1 \times 10^{-6} \text{ }^\circ\text{C}^{-1}$ can be improved by using a somewhat different composition for the resonator material.) b) The (single-sideband) power density spectrum of the phase noise of the oscillator with reaction-type resonator in fig. 11a. See also the caption to fig. 9.

[14] P. Harrop, P. Lesartre and T. H. A. M. Vlek, Low-noise 12GHz front-end designs for direct satellite television reception, Philips Tech. Rev. 39, 257-268, 1980.

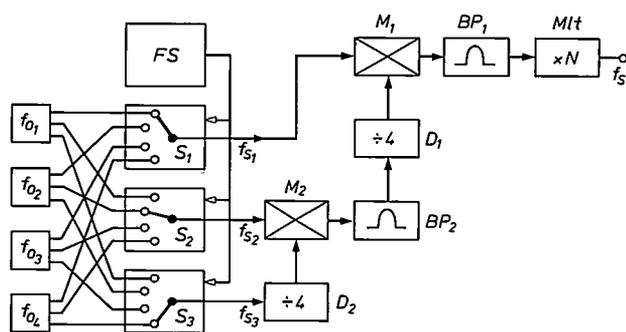


Fig. 13. Block diagram of the circuit of a direct microwave synthesizer. f_{01} to f_{04} microwave oscillators with different frequencies and stabilized by dielectric resonators. FS frequency selector. S_1 to S_3 four-position switches. f_{s1} to f_{s3} signals that can optionally have the frequency of one of the oscillators f_{01} to f_{04} . M_1 , M_2 mixers. D_1 , D_2 frequency dividers, for making the signal frequency four times smaller. BP_1 , BP_2 bandpass filters. Mlt multiplier that multiplies the frequency of the signal by N . f_s output signal.

perature coefficient would have been even smaller and less than $10^{-6} \text{ }^\circ\text{C}^{-1}$. Fig. 12b shows the result of noise-spectrum measurements on the circuit with the reaction-type resonator. The spectrum contains some amplitude noise as well as phase noise. The noise from the oscillator is so low that the result of the measurement almost corresponds to the noise spectrum of the highly sensitive spectrum analyser (Hewlett-Packard 8566 A).

The complete transponder, shown in fig. 10b, contains in addition to the oscillator circuit a number of batteries and the transmitting and receiving circuit with the accompanying antennas. The antennas are located underneath the protective caps at top and bottom. The transmission power is 200 mW, sufficient for a range of 50 km. The transponder is manufactured by the Philips GmbH supply centre Systeme und Sondertechnik, in Bremen.

Oscillators for direct frequency synthesis

Microwave synthesizers — circuits for synthesizing signals at microwave frequencies — are usually *indirect* synthesizers. This means that the frequency of the output signal is stabilized by a phase-locked loop (PLL) in the circuit. The operation of *direct* microwave synthesizers is based on the direct mixing of highly stable reference signals. Circuits in which the

reference signals are generated by oscillators stabilized by cavity resonators tend to be rather large. Compact circuits that give a stable output signal can be built using oscillators with dielectric resonators^[15].

Fig. 13 shows a possible block diagram for such a circuit. There are four oscillators, and the output signal can have 64 different frequencies. The oscillators are connected to the output by switches, mixers, bandpass filters, dividers and a multiplier. The frequency f_s of the output signal is given by

$$f_s = N \left(f_{s1} + \frac{f_{s2} + f_{s3}/4}{4} \right),$$

where f_{s1} , f_{s2} and f_{s3} are each set equal to one of the frequencies f_{01} , f_{02} , f_{03} or f_{04} of the oscillators. The multiplication factor N of the multiplier can be small when dielectric resonators are used. Large values of N are necessary when crystal-stabilized oscillators are used, since these oscillators operate at a frequency of about 100 MHz. A high value of N increases the noise component in the output signal. Circuits using microwave oscillators stabilized by dielectric resonators can therefore deliver an output signal that combines high stability with very low noise.

[15] G. Lütkeke, Dielectric resonator stabilised reference oscillators, IEE Colloq. on Microwave programmable sources and frequency synthesis, London 1981, pp. 3/1-3/4.

Summary. Dielectric resonators can be used to stabilize the frequency of the output signal of microwave integrated oscillators. Oscillators of this kind will be widely used in the near future in microwave receivers for satellite-television broadcasts. The required high Q-factor, low temperature-dependence and small dimensions of dielectric resonators only became feasible with the availability of mixed ceramics of $\text{Ba}_2\text{Ti}_9\text{O}_{20}$ and BaTi_4O_9 or TiO_2 . These materials can be evaluated for their usefulness by generating microwave resonances in cylindrical samples and measuring the frequency spectra. In practice, dielectric resonators are discs. They are mounted in a metal case and can be mechanically tuned to the required frequency. The resonators have a small positive temperature coefficient to compensate for the negative temperature coefficients of other elements in the oscillator circuit. A locus for the impedance of coupling line, resonator and load of the active element in the oscillator circuit as a function of oscillation frequency can be drawn on the Smith chart. The intersection of this locus with the line for the input impedance of the active element gives the operating point. The position of the operating point may be stable or unstable. The oscillator circuits are designed for operation with reaction-type, transmission-type, reflection-type or feedback resonators. Stable oscillator circuits have been designed for transponders used in radiosondes. The use of dielectric resonators in oscillators for direct frequency synthesizers has also been investigated.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	<i>A</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	<i>B</i>
Philips Natuurkundig Laboratorium, Postbus 80000, 5600 JA Eindhoven, The Netherlands	<i>E</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	<i>R</i>

H. Boulard, Y. Kamp, H. Ney & C. J. Wellekens <i>B, H</i>	Speaker-dependent connected speech recognition via dynamic programming and statistical methods	Biblioth. Phonet. 12	115-148	1985
J. B. H. Peek <i>E</i>	Digitale Signalverarbeitung, ein Fachgebiet mit wachsender Bedeutung	E & M 103	163-169	1985
L. M. G. Feijs & J. H. Obbink <i>E</i>	Process models: methods as programs	ESPRIT '85, Elsevier Science, Amsterdam	577-591	1986
G. F. M. Beenker, T. A. C. M. Claasen & P. J. van Gerwen <i>E</i>	Design of smearing filters for data transmission systems	IEEE Trans. COM-33	955-963	1985
M. J. H. van de Steeg, H. L. Peek, J. G. C. Bakker, J. A. Pals, B. G. M. H. Dillen & J. M. A. M. Oppers <i>E</i>	A frame-transfer CCD color imager with vertical antiblooming	IEEE Trans. ED-32	1430-1438	1985
H. D. Hinz & H. Löbl <i>H</i>	Electrophoretic recording of electronically stored radiographs	IEEE Trans. MI-4	39-43	1985
F. H. P. M. Habraken*, R. H. G. Tijhaar*, W. F. van der Weg* (* Univ. Utrecht), A. E. T. Kuiper & M. F. C. Willemsen <i>E</i>	Hydrogen in low-pressure chemical-vapor-deposited silicon (oxy)nitride films	J. App. Phys. 59	447-453	1986
E. F. Stikvoort <i>E</i>	Digital dynamic range compressor for audio	J. Audio Eng. Soc. 34	3-9	1986
P. C. M. Gubbens*, A. M. van der Kraan* (* Interuniv. Reactor Inst., Delft) & K. H. J. Buschow <i>E</i>	Crystal field in Tm_2Ni_{17}	J. Magn. & Magn. Mater. 54-57	483-484	1986
P. Haaker, E. Klotz, R. Koppe, R. Linde & H. Möller <i>H</i>	A new digital tomosynthesis method with less artifacts for angiography	Med. Phys. 12	431-436	1985
P. Delsarte, Y. Genin & Y. Kamp <i>L</i>	Pseudo-Carathéodory functions and Hermitian Toeplitz matrices	Philips J. Res. 41	1-54	1986
K. H. J. Buschow & D. B. de Mooij <i>E</i>	Structural and magnetic characteristics of several ternary compounds of the type GdX_2Si_2 and UX_2Si_2 ($X = 3d, 4d$ or $5d$ metal)	Philips J. Res. 41	55-76	1986
J. Haisma, U. K. P. Biermann, J. Peersman & P. C. Zalm <i>E</i>	Resistance modification of indium-oxide layers by laser annealing and ion bombardment respectively	Philips J. Res. 41	77-81	1986
U. Neitzel (CHF Müller Med.-Tech. Syst., Hamburg), J. Kosanetzky & G. Harding <i>H</i>	Coherent scatter in radiographic imaging: a Monte Carlo simulation study	Phys. Med. & Biol. 30	1289-1296	1985
W. Hoving*, P. M. L. O. Scholte*, P. Dorenbos*, G. A. Fokkema*, E. A. G. Weits*, F. van der Woude* (* Univ. Groningen), I. Vincze (Hungarian Acad. Sci., Budapest) & K. H. J. Buschow <i>E</i>	Packing and chemical effects in amorphous Fe-Zr and Fe-B alloys	Phys. Rev. B 32	8368-8371	1985

G. Duggan, H. I. Ralph & K. J. Moore <i>R</i>	Reappraisal of the band-edge discontinuities at the $Al_xGa_{1-x}As$ -GaAs heterojunction	Phys. Rev. B 32	8395-8397	1985
J. Landman*, W. D. Luedtke*, R. N. Barnett*, C. L. Cleveland*, M. W. Ribarsky* (* Georgia Inst. Technol., Atlanta, GA), E. Arnold, S. Ramesh, H. Baumgart, A. Martinez & B. Khan <i>N</i>	Faceting at the silicon (100) crystal-melt interface theory and experiment	Phys. Rev. Lett. 56	155-158	1986
A. A. Jafari (City Univ., New York), T. N. Saadawi & M. Schwartz <i>N</i>	Blocking probability in two-way distributed circuit-switched CATV	Proc. GLOBECOM '85, New Orleans 1985	868-875	1985
A. H. Reader, F. W. Schapink* & S. Radelaar* (* Univ. Technol. Delft) <i>E</i>	Grain growth in heavily implanted LPCVD polysilicon layers during annealing	Proc. Microsc. Semicond. Mater. Conf., Oxford 1985	151-156	1985
G. Duggan, H. I. Ralph & R. J. Elliot (Dept. Theor. Phys., Oxford) <i>R</i>	Quantisation effects and interface recombination in GaAs-(AlGa)As quantum wells	Proc. MRS Conf., Strasbourg 1985	37-40	1985
G. Duggan, H. I. Ralph, K. S. Chan* & R. J. Elliott* (* Dept. Theor. Phys., Oxford) <i>R</i>	Exciton binding energy in quantum wells including the heavy-light hole interaction	Proc. MRS Conf., Strasbourg 1985	47-51	1985
B. A. Joyce, J. H. Neave, P. J. Dobson, K. Woodbridge & J. Zhang (Imp. College, London) <i>R</i>	RHEED oscillations and surface structure of MBE-grown GaAs-(Al,Ga)As	Proc. MRS Conf., Strasbourg 1985	135-139	1985
H. Dimigen, H. Hübsch, P. Willich & K. Reichelt (Kernforschungsanlage Jülich) <i>H</i>	Stoichiometry and friction properties of sputtered MoS_x layers	Thin Solid Films 129	79-91	1985
W. Benecke*, H. Dimigen, P. Eichinger*, H. G. Feller*, U. Jahns*, R. Klinger* (* Tech. Univ. Berlin), K. Kobs, R. Leutenecker*, H. Ryssel*, H.-P. Spöhrle* & K. H. Zeller* (* Fraunhofer-Inst., München) <i>H</i>	Tribologisches Verhalten Stickstoff-implantierter Stähle	Tribologie — Reibung/Verschleiss/Schmierung, Springer, Berlin	335-365	1985

Contents of Philips Telecommunication Review 44, No. 1, 1986

- R. J. Mulder: Office automation and the ISDN (pp. 2-13)
 R. Heemskerk: SPELL project engineering tools for SOPHO S (pp. 14-25)
 J. P. M. van Kleef: The computer voice (VPU) assisting PABX operators with impaired eyesight (pp. 26-35)
 J. van Duuren: The U_3 digital two-wire PABX interface (pp. 36-43)

Contents of Electronic Components & Applications 7, No. 4

- P. Anders & T. van Kampen: Parallel processing and pipelining usher DSP into the future (pp. 202-207)
 P. Melville: Thermal resistance of SO and PLCC packages (pp. 208-216)
 C. Franx: ZrO_2 oxygen sensor (pp. 217-222)
 R. Zavrel: State-of-the art ICs simplify SSB-receiver design (pp. 223-228)
 M. ten Have: Speech synthesis — the complete approach with the PCF8200 (pp. 229-238)
 A. H. H. J. Nillesen: Line-locked digital colour decoding (pp. 239-245)
 J. J. de Jong: In-line system for surface-mount assembly accepts large circuit boards (pp. 246-249)
 Insulated encapsulations for easier mounting of power semiconductors (pp. 250-255)
 P. F. de Jongh & S. B. Worm: YJ1600 magnetron for microwave heating up to 6 kW (pp. 256-260)

Silicon cold cathodes

G. G. P. van Gorkom and A. M. E. Hoeberechts

The release of electrons from the cathode of a cathode-ray tube is usually a thermal process (thermionic emission). For some applications there are also cathodes that are not heated, called cold cathodes. The energy that the electrons require to enable them to escape into the vacuum is supplied in some other way, for example by photons (the photoelectric effect), by a strong external electric field (field emission) or by a strong internal electric field in a semiconductor with a reverse-biased p-n junction. This last case applies in the silicon cold cathodes that are the subject of this article. At first such cathodes seemed to be of little practical use. Further investigation, however, making use of modern semiconductor technology, has resulted in improvements that now offer good prospects of practical applications.*

Introduction

It has been known for nearly thirty years that electron emission into a vacuum can occur from a reverse-biased p-n junction at or near the surface of a semiconductor [1]. The electric field in the depletion layer formed on both sides of the p-n junction by the removal of electrons and holes may be so strong locally that 'avalanche breakdown' occurs, in which electrons and holes are created. The resultant acceleration ('heating') of the electrons may be strong enough for some of the electrons to be emitted from the semiconductor surface.

This effect has attracted considerable interest for some time, because it seemed to offer the possibility of making a cold cathode. Compared with the more conventional thermionic cathodes, a cold cathode would have certain advantages, e.g. it could carry a higher current density and it could be switched more easily. Most of the early experiments were on p-n junctions in silicon [2], but other semiconductors such as silicon carbide (SiC) were investigated later [3]. As far as we know, however, these investigations did not lead to practical applications, probably because the efficiency

of the emission process was insufficient and the current density of the silicon cathodes was too low on account of the large dimensions of the emissive area. Problems encountered with SiC cathodes included their poor reproducibility and their unusual geometry.

In the typical geometry the plane of the p-n junction intersected the semiconductor surface, so that the electrons were emitted from a very small region around the line of intersection. It was also found possible to produce electron emission from a p-n junction parallel to the surface [4]. Since the early seventies another

[1] J. A. Burton, Electron emission from avalanche breakdown in silicon, *Phys. Rev.* **108**, 1342-1343, 1957.

This effect was discovered at about the same time by P. Zalm and F. van der Maesen in an investigation initiated by H. B. G. Casimir.

[2] R. J. Hodgkinson, Hot-electron emission from semiconductors, *Solid-State Electron.* **5**, 269-272, 1962; D. J. Bartelink, J. L. Moll and N. I. Meyer, *Phys. Rev.* **130**, 972-985, 1963;

M. Waldner, Hot electron emission from silicon surfaces, *J. Appl. Phys.* **36**, 188-192, 1965.

[3] See R. V. Bellau and A. E. Widdowson, Some properties of reverse-biased silicon carbide p-n junction cold cathodes, *J. Phys. D* **5**, 656-666, 1972.

[4] See for example the article by Bartelink *et al.*, note [2]. A modern version of this experiment has been described by I. Shahriary, J. R. Schwank and F. G. Allen in *J. Appl. Phys.* **50**, 1428-1438, 1979.

effect that has been extensively studied, in addition to electron emission in vacuum, is electron injection from silicon in SiO_2 layers, partly because of the resultant degradation of MOS transistor characteristics and partly because of a possible practical application as a non-volatile memory [6]. These studies showed that high current densities and high injection efficiencies can be achieved [6].

These results prompted us to look once again at the potential uses of cold silicon cathodes, but this time making use of modern semiconductor technology. In the investigations described in this article various Si cathodes were designed and made, and properties such as the energy distribution of the emitted electrons, the emission efficiency, the current density and the total current were examined. On the basis of the results we attempted to improve the materials, the doping and the geometry so as to satisfy the requirements for practical application.

Our experiments showed that cathodes with the p-n junction parallel to the silicon surface gave the best results [7]. High current densities can be obtained when the p-n junction is only a little way (say 10 nm) below the surface and the emitting region is very small, e.g. 1 μm in diameter for a region with rotational symmetry. To lower the work function of the electrons the silicon surface is coated with a very thin layer of caesium. In this way cathodes were made that gave a current density of 1500 A/cm^2 at an efficiency of 1.5%. This current density is much higher than can be achieved with more conventional cathodes.

In this article we shall first give a general description of the electron emission from a p-n junction. We shall then discuss some aspects of the manufacture of silicon cold cathodes, in particular the geometry adopted and the doping required for the silicon. Next we shall deal with some of the emission characteristics: efficiency, energy distribution, current density, brightness and maximum current. We conclude with a brief consideration of the potential applications of these cathodes.

Electron emission at a p-n junction

To obtain efficient electron emission from a p-n junction in silicon a heavily doped p-type layer and an even more heavily doped very thin n-type layer are applied. Reverse biasing the junction then results in an electron-energy diagram as shown in *fig. 1*. As can be seen in this figure, the presence of a very strong electric field in the depletion layer has the effect of accelerating or 'heating' the electrons (and the holes too, but we do not need to consider these further here). The hot electrons can generate new electron-

hole pairs by collisions with valence electrons. In their turn, these can excite other electrons from the valence band, so that an avalanche of electrons is produced. The collision processes and the interaction with optical phonons have the effect of reducing the energy of the electrons ('cooling' them). The balance between heating and cooling gives an energy distribution with a characteristic electron temperature that is much higher than the temperature of the silicon lattice. Electrons that have a kinetic energy near the surface higher than or equal to the work function have a chance of escaping into the vacuum.

It follows from all this that a cold silicon cathode emits 'hot' electrons. In this respect such a cathode corresponds more to the thermionic cathodes than to the other cold cathodes, e.g. those with a negative electron affinity (NEA cathodes) [8] and field emitters [9].

In an NEA cathode, e.g. p-type GaAs coated with caesium, the bottom of the conduction band has a higher energy than the vacuum level [8]. Electrons in the conduction band do not then require any kinetic energy to escape into the vacuum. Unlike the silicon cathode, an NEA cathode therefore emits 'cold' electrons. By the same token it could be said that a field-emission cathode emits 'lukewarm' electrons. In a field emitter, either a metal or a very heavily doped semiconductor, the electrons are subjected to a strong external electric field and 'tunnel' through the potential barrier at the surface [9]. Most of the emitted electrons then have an energy not very different from the Fermi energy in the material.

To emit electrons at high efficiency a silicon cathode must meet some special requirements. For example, the doping and geometry have to be such that a very

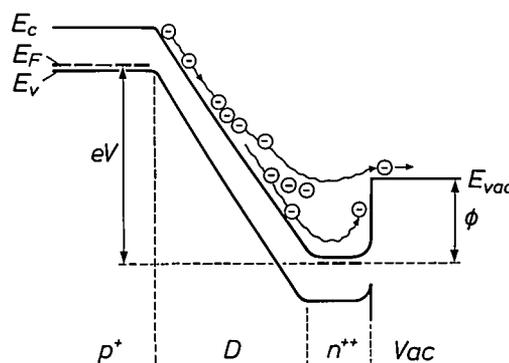


Fig. 1. Energy-level diagram for the electrons in a p-n junction close to the surface, and with a reverse bias voltage V . E_v top of the valence band. E_F quasi-Fermi level. E_c bottom of the conduction band. E_{vac} potential energy of an electron in vacuum. Electrons coming from the p^+ region are accelerated by the strong electric field in the depletion layer D . On the way to the very thin n^+ channel they lose some of their energy by interaction, mostly with optical phonons but possibly by collision with valence electrons, so that more and more of the new electrons arrive in the conduction band (avalanche breakdown). Electrons whose kinetic energy at the surface is equal to or higher than the work function ϕ have a chance to escape into the vacuum Vac .

strong electric field can be generated in a thin depletion layer. This field, however, must never become so strong as to cause a tunnel breakdown (the Zener effect) in the semiconductor material, since the energy of tunnelling electrons is low compared with the work function^[10]. In the n-type layer at the surface the electrons will be cooled down by collisions, and this layer must therefore be very thin, but without the series resistance becoming too high. The best results are achieved when the work function is low, of course. The silicon surface is therefore coated with a monoatomic layer of caesium. The cathode must also have reliable electric contacts and it must have high stability. All these aspects play an important part in the production of cold silicon cathodes.

Manufacture

Manufacturing a cold cathode based on the emission mechanism described above consists in making p-type and n-type regions in silicon that have been doped and structured in such a way as to produce a thin depletion layer and a very small active electron-emitting area. Electrical connections are also required for connecting the cathode into an electrical circuit. It is essential to obtain the appropriate geometry and doping, so that avalanche breakdown only occurs at the places where it is required.

Experiments with different geometries showed that cathodes with a p-n junction that intersects the surface are not suitable for practical applications. This is because the electron emission from such cathodes has a very broad energy distribution^[11], which is disadvantageous for application in devices such as cathode-ray tubes. After being coated with a layer of caesium, these cathodes also present stability problems because positively charged Cs atoms migrate over the surface as a result of the strong electric field in the depletion layer and just outside it. We shall therefore confine ourselves here to cathodes with a p-n junction parallel to the surface.

An example of a geometry that has given good results^[7] is shown schematically in *fig. 2*. The electrical connections required are also indicated. The circuit on the silicon slice actually contains two parallel diodes: a small electron-emitting diode and a contact diode, which also provides protection from electrical overload.

The cathode contains p-type and n-type regions with different structures and dopings. It was made by photolithographic methods currently used in semiconductor technology. To obtain a low series resistance the process starts with an epitaxial p-type layer on a heavily doped (p⁺) substrate. By means of a phos-

phorus diffusion a heavily doped (n⁺) region about 3 μm thick is made for the electrical contact with the voltage source. The active area consists of a p⁺ layer produced by implantation of boron ions. Above this layer is a very heavily doped (n⁺⁺) channel, which is produced by implantation of arsenic or antimony ions. The dopant concentrations of the p⁺ region and the n⁺⁺ channel might typically be 5 × 10¹⁷ and 4 × 10¹⁹ cm⁻³. The implantation is carried out at very low energy, so that the n⁺⁺ channel can be made very thin (about 10 nm). The electron-emitting area has rotational symmetry here, with a diameter of less than 10 μm, but it can also take the form of an ellipse or other shape.

At the edge of the diode the depletion layer comes to the surface. This edge is protected, e.g. by a 1-μm SiO₂ layer or, better still, by two SiO₂ layers with a

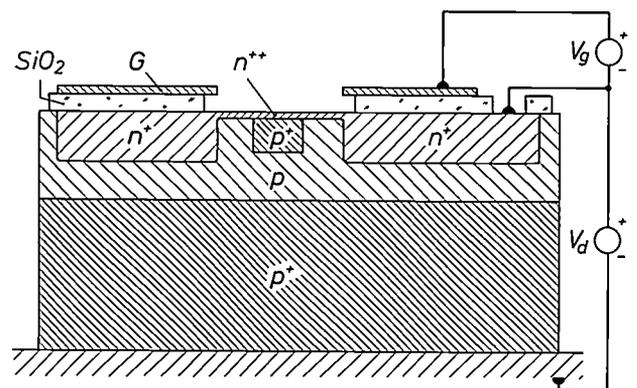


Fig. 2. Schematic cross-section of a silicon cold cathode and the electrical connections. A p⁺ substrate about 400 μm thick carries p-type and n-type regions with different structures and dopings. The electron-emitting p⁺ region here has rotational symmetry with a diameter < 10 μm. On top of it is a very thin n⁺⁺ channel about 10 nm thick. Electrical contact with the voltage source V_d is made by means of a thicker n⁺ region (3 μm). The p⁺ substrate is connected to this voltage source by a gold contact. A structured SiO₂ layer about 1 μm thick provides protection for the emitting diode and carries a conducting layer G (gate), which is connected to the gate-voltage source V_g.

[6] A good general account of electron emission from silicon in SiO₂ layers is given by T. H. Ning, *Solid-State Electron.* **21**, 273-282, 1978.

[7] This has been described by J. F. Verwey and B. J. de Maagt in *Solid-State Electron.* **17**, 963-971, 1974 for a p-n junction that intersects the Si/SiO₂ interface, and by T. H. Ning, C. M. Osburn and H. N. Yu in *J. Appl. Phys.* **48**, 286-293, 1977 for a p-n junction parallel to the Si/SiO₂ interface.

[8] G. G. P. van Gorkom and A. M. E. Hoeberechts, An efficient silicon cold cathode for high current densities, I: Experimental data and main results, *Philips J. Res.* **39**, 51-60, 1984.

[9] More information on NEA cathodes can be found in J. van Laar and J. J. Scheer, *Philips Tech. Rev.* **29**, 54-66, 1968, or in six articles in *Acta Electron.* **16**, 215-271, 1973.

[10] See for example A. van Oostrom, Field emission of electrons and ions, *Philips Tech. Rev.* **33**, 277-292, 1973.

[11] G. G. P. van Gorkom and A. M. E. Hoeberechts, Performance of silicon cold cathodes, *J. Vac. Sci. & Technol. B* **4**, 108-111, 1986.

[12] G. G. P. van Gorkom and A. M. E. Hoeberechts, Electron emission from depletion layers of silicon p-n junctions, *J. Appl. Phys.* **51**, 3780-3785, 1980.

layer of Si_3N_4 sandwiched between them. On top of the insulation there is a conducting layer, about $0.5 \mu\text{m}$ thick, of a material such as polysilicon or aluminium. This layer is called the 'gate', for historical reasons, since the present Si cathode originated from the MOS transistor. It prevents charge from building up on dielectric layers, screens the strong electric fields at the p-n junction near the edge, and may also act as a first electrode in an electron-optical system.

After appropriate structures have been applied to the silicon slice, it is sawn into rectangular chips of say $2 \text{ mm} \times 2 \text{ mm}$. Fig. 3 shows a part of a chip in a

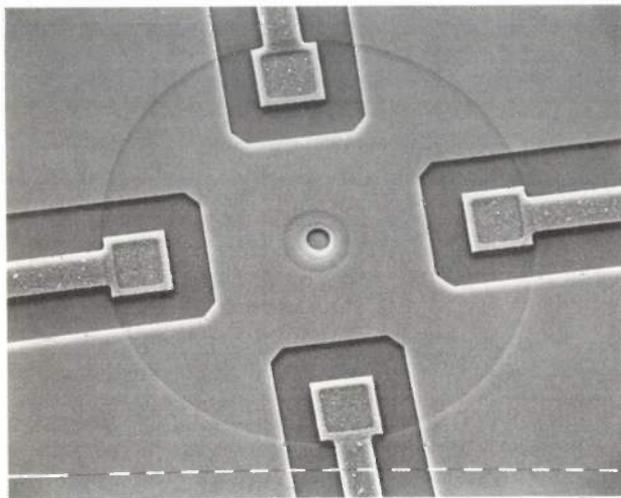


Fig. 3. Scanning electron micrograph of part of a chip carrying a silicon cold cathode (each dash represents $10 \mu\text{m}$). The electron emission comes from a non-visible region below the surface (see fig. 2) at the centre. The circle around it corresponds to the edge of a hole in the gate. The other concentric circles correspond to steps in the oxide layer, made among other reasons for applying the contact diode. The electrical connection to the contact diode is quadrupled here for the electron optics.

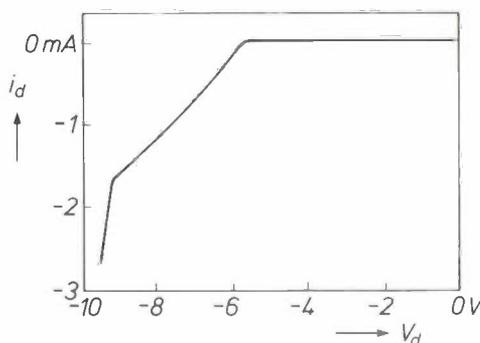


Fig. 4. Current-voltage characteristic of a silicon cold cathode. The measured diode current i_d is plotted against the diode voltage V_d . Negative values of V_d and i_d correspond to a voltage and current for a reverse-biased diode. For a low reverse bias a negligible current is measured. From about -6 V an avalanche breakdown occurs in the electron-emitting region, causing the reverse current to rise steeply, approximately linearly with the reverse bias. At about -9 V the contact region breaks down and the reverse current increases even more steeply with the reverse bias.

photomicrograph made with a scanning electron microscope. A chip is soldered to a specially cleaned substrate, which also provides the electrical contact with the p^+ substrate. Aluminium leads, between 25 and $40 \mu\text{m}$ thick, are ultrasonically bonded to the outer part of the n^+ region and to the gate.

The current-voltage characteristic of a silicon cold cathode depends on the geometry and the dimensions of the active area. A typical characteristic is given in fig. 4. When the reverse bias is increased a value is reached at which there is avalanche breakdown in the active area. A diode current then flows, accompanied by the emission of electrons into the vacuum. A further increase in the reverse bias results in a virtually linear increase of the diode current (and the electron emission) until there is also breakdown in the contact region. The operating voltage of the cathode must be between the two breakdown voltages. In fig. 4 the operating voltage is between about -6 and -9 V .

We also made another version of a silicon cold cathode in which the breakdown region is less accurately defined [12]. This version is easier and quicker to make and is good enough for studying stability problems and for use in experimental television tubes. An example of a television tube with such a cathode will be discussed later on.

After the cathode has been mounted in the vacuum system in which it is to operate, a monoatomic layer of caesium is deposited on its surface. As this layer is required to be of extreme purity, the caesium source is first thoroughly outgassed to remove oxidizing gases such as O_2 , H_2O and CO_2 . The deposition of a monolayer does not present many problems. Since the adsorption energy for caesium on caesium is fairly low (about 0.7 eV), and that for caesium on silicon is much higher (about 2.0 eV) the caesium has a strong tendency to form a monolayer on the surface; the excess caesium merely evaporates.

Emission characteristics

An important feature of silicon cold cathodes is the efficiency of the emission process, i.e. the current due to the emitted electrons divided by the total current. It is also important to know the energy distribution of the emitted electrons, and the values obtainable for the current density, brightness and total current. We shall now briefly discuss some results obtained with our cathodes.

Efficiency

Only a small proportion of the conduction electrons in the depletion region have sufficient energy at the surface for them to escape into the vacuum. This im-

plies that the emission current i_{vac} into the vacuum is much smaller than the diode current i_{d} of the non-emitted electrons. The efficiency η is thus given by

$$\eta = \frac{i_{\text{vac}}}{i_{\text{d}} + i_{\text{vac}}} \approx \frac{i_{\text{vac}}}{i_{\text{d}}}. \quad (1)$$

The highest values of η are obtained with cathodes in which the active layer is doped in such a way that the electric field-strength in the depletion layer has a maximum value just below the threshold at which the undesirable effect of tunnel breakdown sets in (about 10^6 V/cm). The effect of the reduction of the work function on the efficiency is greatest when the surface is covered with a monoatomic layer of caesium. The coating lowers the work function by 2.5 to 3 eV, which increases the efficiency by a factor of between 400 and 1000, depending on the purity of the surface. The optimum thickness of the n^{++} channel (fig. 2) is found to be about 10 nm. Although a thinner channel would allow more electrons to escape, the series resistance to the flow of non-emitted electrons to the contact diode would become too large.

Efficiencies of between 1 and 2% have been measured for cathodes with the geometry, doping and caesium layer mentioned above, in which modern surface analysis revealed the presence of some residual contaminations on the surface, mainly carbon and oxygen. After *in situ* removal of this oxygen by electron bombardment an efficiency of about 5% can be achieved [13]. An even higher efficiency (7.25%) has been measured after sputter-cleaning of the surface by bombardment with positive argon ions [14], but this made the n^{++} channel too thin.

The results of the efficiency measurements can be qualitatively interpreted in terms of a 'lucky-electron' model, as described for electron injection from silicon in SiO_2 [15][16]. To reach the vacuum the electrons must possess an energy at least equal to the work function (fig. 1). However, the only ones of these electrons that will be able to escape are those that are 'lucky' enough to lose no energy by interaction with the crystal lattice on their way to the surface; interactions with optical phonons are especially likely and collisions with valence electrons can also occur. For an electron that has sufficient potential energy at a distance x below the surface the probability P of being emitted is given by:

$$P = A \exp(-x/\lambda), \quad (2)$$

where A is a normalization constant and λ is an effective mean free path [6]. Characteristic values derived from experiments on Si/SiO₂ are $A = 2.9$ and $\lambda = 9$ nm [16]. To arrive at the efficiency we have to multiply P by the fraction ε of the number of electrons whose potential energy is equal to or greater than the work function:

$$\eta = \varepsilon A \exp(-x/\lambda). \quad (3)$$

The reduction of the work function by 2.5 to 3 eV because of the monolayer coating of caesium results in a marked increase in ε and a decrease in x , so that the efficiency rises considerably.

Another theory that can be used for describing the emission characteristics of cold silicon cathodes is based on Boltzmann's transport equation for charge carriers in an electric field [17], but this will not be dealt with here.

Energy distribution

We have said that the electrons in a silicon cathode are heated to a high temperature. The value of the effective 'thermal' electron energy kT_e can be determined by measuring the energy distribution of the emitted electrons [7][11]. The distribution found for an Si cathode with no caesium coating is usually in good agreement with a distribution calculated for thermal emission [18]. The distribution function $N(E_k)$ of the kinetic energy E_k of the emitted electrons is then given by

$$N(E_k) \propto E_k \exp(-E_k/kT_e). \quad (4)$$

The quantity kT_e can be derived from the half-value width ΔE_k of the distribution curve:

$$kT_e = \Delta E_k/2.45. \quad (5)$$

A characteristic example of a measured energy distribution with the associated theoretical distribution is given in fig. 5a. The half-value width in this case is 0.92 eV. Eq. (5) shows that this corresponds to $kT_e = 0.38$ eV and $T_e = 4400$ K.

The reduction in the work function by coating with a monolayer of caesium is associated with a broadening of the energy distribution, with a clearly observable deviation from the theoretical energy distribution; see fig. 5b. We shall not consider the origin of these effects further here. However, the deviation from the theoretical distribution is not too large in this case, so that eq. (5) is a reasonable approximation for determining kT_e . A value of 1.20 eV is derived for the half-value width. This means that $kT_e = 0.49$ eV and $T_e = 5700$ K.

The measured energy distribution is much narrower than that of a p-n junction that intersects the surface [11], but it is nevertheless significantly broader

[12] A. M. E. Hoeberechts and G. G. P. van Gorkom, Design, technology and behaviour of a silicon avalanche cathode, *J. Vac. Sci. & Technol. B* 4, 105-107, 1986.

[13] This was found in experiments by J. Zwier and J. H. A. Vasterink of these Laboratories.

[14] A. van Oostrom, L. J. M. Augustus, G. G. P. van Gorkom and A. M. E. Hoeberechts, to be published.

[15] J. F. Verwey, R. P. Kramer and B. J. de Maagt, Mean free path of hot electrons at the surface of boron-doped silicon, *J. Appl. Phys.* 46, 2612-2619, 1975.

[16] See the article by Ning *et al.*, note [6].

[17] G. A. Baraff, Distribution functions and ionization rates for hot electrons in semiconductors, *Phys. Rev.* 128, 2507-2517, 1962.

[18] R. D. Young, Theoretical total-energy distribution of field-emitted electrons, *Phys. Rev.* 113, 110-114, 1959.

than that of conventional cathodes. In some cases this may be a disadvantage, as in television camera tubes, which should really have a narrow energy distribution to give a rapid response^[19]. At relatively low accele-

rating voltages there may be chromatic aberration. For many applications, however, including television picture tubes, the energy spread causes no problems. The negative effect on the brightness is more than compensated by the much higher current density obtainable with a silicon cathode.

Current density and brightness

The current density j_{vac} for electron emission into vacuum is given by

$$j_{\text{vac}} = \eta j_{\perp}, \quad (6)$$

where j_{\perp} is the current density perpendicular to the surface of the diode; see fig. 6. To achieve high current densities it is therefore necessary not only to have a high η but also the largest possible value of j_{\perp} . However, the limiting factor for j_{vac} is not j_{\perp} but j_{\parallel} , the diode current density parallel to the surface. This is because the diode current of the non-emitted electrons has to be squeezed into the very thin n^{++} channel. At an optimum channel thickness of about 10 nm the emitting area must be very small to avoid 'current crowding' of the electrons. It has been found empirically that, for stable emission, j_{\parallel} must be no larger than about 2.5×10^6 A/cm² [7] [10].

The relation between j_{\perp} and j_{\parallel} is determined by the areas of the regions through which the diode current flows, perpendicular and parallel to the surface. These are easily derived with the aid of fig. 6: for j_{\perp} it is the area of a circle of radius r , for j_{\parallel} it is the area of a cylindrical shell of radius r and the channel thickness d as the height. The ratio j_{\perp}/j_{\parallel} is therefore given by:

$$j_{\perp}/j_{\parallel} = \frac{2\pi r d}{\pi r^2} = 2d/r. \quad (7)$$

For the smallest emitters so far produced, r is only 0.5 μm . With $d = 10$ nm and $j_{\parallel} = 2.5 \times 10^6$ A/cm² (the maximum permissible value) a value of 10^5 A/cm² has been reached for j_{\perp} , so that from equation (6) j_{vac} is equal to $\eta \times 10^5$ A/cm². This means that with $\eta = 1.5\%$ a current density of 1500 A/cm² into the vacuum is obtained, which can be very stable under adequate vacuum conditions (to be discussed later).

The brightness of the electron emission is defined as the current density divided by the solid angle within which the electrons are emitted. The brightness B is proportional not only to the current density j_{vac} but also to the ratio of the potential electron energy at an applied field V and the effective ('thermal') electron energy kT_e :

$$B = \frac{e j_{\text{vac}} V}{\pi k T_e}. \quad (8)$$

Eq. 18 shows that a brightness of 9×10^6 A/cm²sr

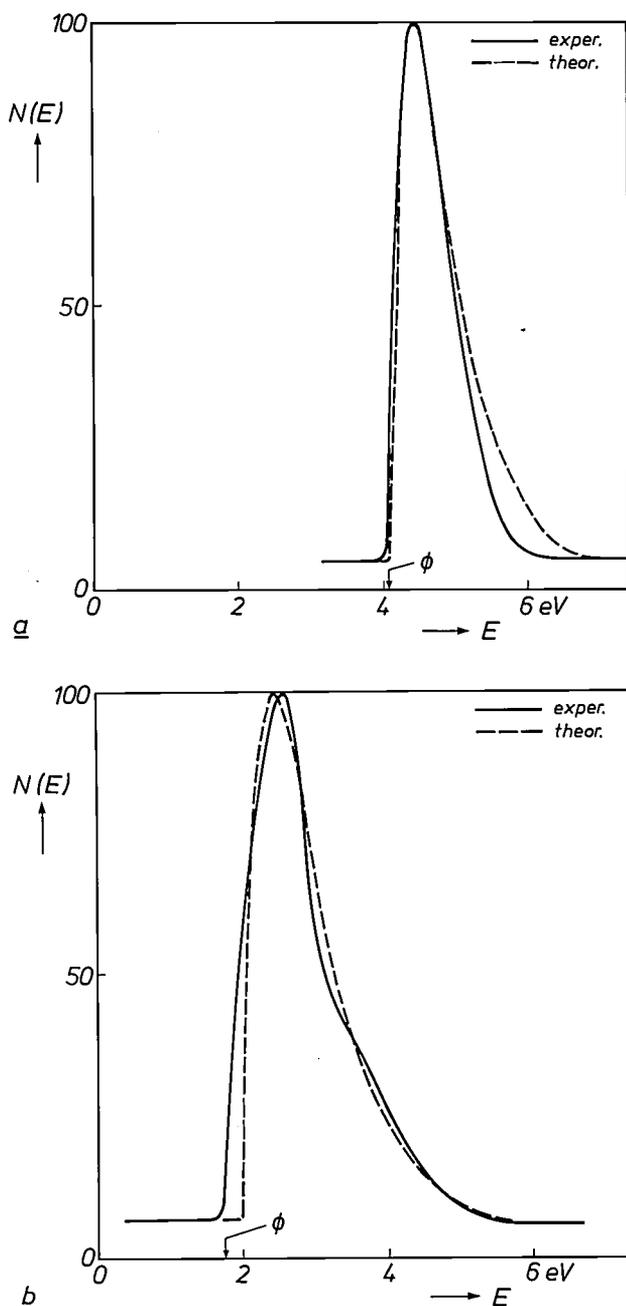


Fig. 5. Characteristic energy distribution of the emitted electrons for two silicon cold cathodes, one (a) coated (unintentionally) with about 0.02 of a monolayer of caesium, and the other (b) coated with a true monolayer of caesium. The distribution function $N(E)$ is given in relative units; E is the sum of the work function ϕ and the kinetic energy E_k . The dashed lines give the theoretical distributions that gives the best fit with the measured curves. These theoretical distributions were obtained by calculating $N(E_k)$ from equation (4), taking kT_e as 0.38 eV (a) and 0.49 eV (b). In b the theoretical curve is shifted slightly towards the higher energy to obtain the best fit. For the cathode with the high work function (with little caesium) a reasonably good agreement is found between the measured and calculated curves. The monolayer of caesium not only gives a low work function but also gives a broader energy distribution, which clearly departs from the best-fitting theoretical distribution.

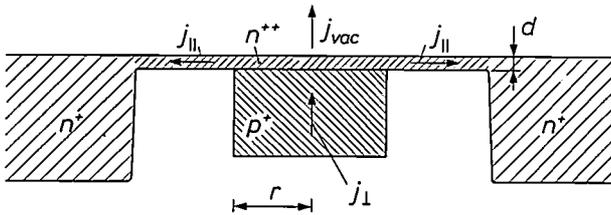


Fig. 6. Schematic geometry for the electron emission of a silicon cold cathode. The arrows indicate the directions of the electron currents. The current density of the electron emission into the vacuum (j_{vac}) is equal to the product of the emission efficiency and the diode current density perpendicular to the surface (j_{\perp}). The diode current density parallel to the surface (j_{\parallel}) is much larger than j_{\perp} , because the thickness d of the n^{++} channel is much smaller than the radius r of the emitting p^{+} area.

Table I. Comparison of the characteristic values for the current density j_{vac} and intrinsic brightness B (at 10 kV) for some conventional thermionic cathodes and a silicon cold cathode. The values for the silicon cathode were obtained for an emitting area of diameter 1 μm .

Cathode	j_{vac} (A/cm ²)	B (A/cm ² sr)
oxide	0.5	2×10^4
W	5	7×10^4
LaB ₆	20-50	$4 \times 10^5 - 1 \times 10^6$
Si	1500	9×10^6

can be obtained from cathodes with $kT_e = 0.5$ eV and $j_{vac} = 1500$ A/cm² at a voltage of 10 kV.

The current density and brightness values obtained compare favourably with those of conventional thermionic cathodes. This can be seen in Table I, which gives the intrinsic value for the brightness of the cathodes. When the cathodes are used in electron guns with a triode structure, as oxide, tungsten and LaB₆ cathodes usually are, the effective (useful) brightness may be much lower owing to the additional energy spread caused by the strong interaction between the electrons [20]. However, the silicon cathodes will usually be used in a diode structure, in which there is much less electron interaction and the effective brightness is roughly equal to the intrinsic value. The Table does not include figures for field emitters. These can give current densities of the order of 10^4 A/cm² and correspondingly high intrinsic brightness values. However, the diameter of the virtual source is extremely small in this case (about 5 nm), so that in practice the effective brightness is often much lower, e.g. because of vibration of the point source. The values are therefore not readily comparable with the values in Table I.

[19] See J. H. T. van Roosmalen, A new concept for television camera tubes, Philips Tech. Rev. 39, 201-210, 1980.

[20] See H. Boersch, Experimentelle Bestimmung der Energieverteilung in thermisch ausgelösten Elektronenstrahlen, Z. Phys. 139, 115-146, 1954 and K. H. Loeffler, Energy-spread generation in electron-optical instruments, Z. Angew. Phys. 27, 145-149, 1969.

Total emission current

In conventional thermionic cathodes the total current can be increased by enlarging the emitting area, and keeping the current density constant. The situation is rather more complex for silicon cold cathodes. As we saw earlier, to achieve high current densities we need a very small emitting area. This means that a compromise has to be found between the obtainable current density and the total current. The final choice will depend on the nature of the application.

From eq. (1) the total emission current i_{vac} is given by:

$$i_{vac} \approx \eta i_d. \quad (9)$$

The value of i_{vac} can be increased by increasing the diode current at a given η . Here again, however, the limitation is the maximum value of j_{\parallel} , since i_d , and hence i_{vac} , are proportional to j_{\parallel} :

$$i_{vac} \approx \eta i_d = 2\pi\eta r d j_{\parallel}. \quad (10)$$

For a cathode with $\eta = 1.5\%$, $r = 0.5$ μm , $d = 10$ nm and $j_{\parallel} = 2.5 \times 10^6$ A/cm² (the maximum permissible value) a total current of about 10 μA is reached.

For some applications this kind of value will be sufficient. If a higher value is required, it can be obtained by increasing the emitting area, although with some loss of current density and brightness. Increasing r from 0.5 μm to say 3 μm , without changing η , d and j_{\parallel} , has the effect of increasing the total emission current from about 10 μA to about 60 μA . This increase reduces the current density and the brightness by a factor of six, since, as equations (6), (7) and (8) show, these quantities are inversely proportional to r . If an even higher emission current is required, it is better to use a different geometry for the emitting area, such as a meander structure. In principle such a structure will give a higher total emission for a given (fixed) current density of about 100 A/cm². The limitations then are the heat generation in the chip and the arrangements necessary for removal of that heat.

Potential applications

The silicon cold cathodes described in this article have features that make them interesting for applications in devices such as cathode-ray tubes. The applications we have particularly in mind are those where advantage can be taken of the high current density and brightness. Another useful aspect is that the electron emission can be modulated easily and rapidly. A voltage of only a few volts will vary the diode current and hence the emission current through a large range (fig. 4). Since the capacitance of the active area is small, very high switching rates can be used. We

expect to reach switching rates of at least the order of 1 GHz.

Fig. 7 shows an example of a possible application in a small television picture tube. An even smaller version of this tube could be used as a monitor in a camera with a solid-state sensor. Since the cold cathode requires little power (say 10 mW) it can be run from a battery, and with no warm-up time required it will give 'instant switch-on'. The high current density will give very small spots with relatively simple electron optics. Another advantage is that the gate of the cathode (fig. 2) can be used as the first electrode of the electron-optical system. This integrated electrode can be divided into say eight segments, which can have individual or paired external connections. 'Dynamic beam shaping' can be obtained by applying a small voltage to pairs of segments, where this small voltage is dependent on the instantaneous position of the spot on the screen. This allows spot corrections right into the corners of a television picture.

Fig. 8 shows a test pattern on the screen of the tube in fig. 7, with the video signal applied direct to the silicon cathode. The resolution of the picture is good: all the lines are clearly resolved. There are however some deflection errors because the electron lens is not properly matched to the deflection coils in this experimental tube. At a frequency of 5 MHz the modulation depth is better than 50%.

In principle one chip can carry several cathodes, so that multibeam tubes can be made.

A disadvantage is that the caesium monolayer on a silicon cathode is very sensitive to oxidizing gases and electric fields. In the long run this could affect the stability. We should therefore conclude with a few words about the stability of the electron emission.

Stability

A direct threat to the stability of the electron emission is the possibility of migration and desorption of caesium atoms. To remove non-emitted electrons it is necessary to have an electric field in the thin n^{++} channel. Because of the series resistance, this field is particularly strong at high current densities (it may reach 4000 V/cm). This field, parallel to the surface, also acts on the positively charged caesium atoms, which may therefore start to drift across the surface. This causes a local build-up of caesium atoms, and because of the low adsorption energy of caesium on caesium, some desorption may take place even at room temperature. At high current densities there can also be direct desorption of caesium atoms as a result of interaction with the escaping electrons. Our investigations into stability have shown, however, that these undesirable effects can be almost completely eliminated if



Fig. 7. Above: Small experimental television picture tube in which the electron source is a silicon cold cathode. The screen has a diagonal of about 10 cm. Below: Detail of part of the tube with the cathode on a type SOT 14 mount.

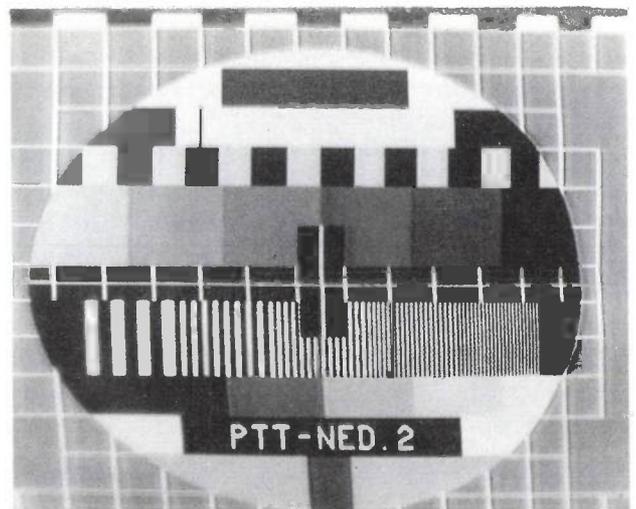


Fig. 8. Test pattern on the screen of the tube in fig. 7 with the video signal (swing about 5 V, 10 mA) applied direct to the silicon cathode. The picture has a high resolution: at 5 MHz the modulation depth is better than 50%. The deflection errors in the picture can be avoided by improving the matching between the electron lens and the deflection coils.

the current density parallel to the surface ($j_{||}$) is not allowed to exceed the value of about 2.5×10^6 A/cm² mentioned earlier.

Another threat comes from the oxidizing residual gases (O₂, H₂O, CO₂) that are present in any vacuum system. Adsorption of these gases on the caesium surface increases the work function, so that fewer electrons are emitted. The better the vacuum the less this effect will affect the stability of the electron emission. Fortunately, it so happens that oxygen can also desorb from the caesium layer by interaction with the escaping electrons. In certain conditions the caesium layer may also be attacked by bombardment from positive ions generated by the electron beam. To avoid this an ion trap can be used, i.e. a device that stops any positive ions while letting the electrons pass.

In experiments under good conditions very long lives have been observed. We can therefore expect that silicon cold cathodes will be applied in practice, especially if their specific characteristics are fully utilized. Further research will clarify the position here.

In the manufacture of the cathodes valuable contributions were made by P. Drummen and members of

the IC department of these laboratories. Technical assistance was provided by P. J. M. Peters, P. A. Dessens and B. P. Cadier. The cathodes were dismantled and incorporated in various tubes in the tube technology department under the supervision of J. van Esdonk.

Summary. A silicon cold cathode emits electrons from an area in which an avalanche breakdown is created by the strong electric field produced at a reverse-biased p-n junction. The emission characteristics depend mainly on the geometry, the doping and the purity of the surface. The best results are obtained with cathodes that have a very shallow p-n junction (about 10 nm below the surface) and a very small emitting area (e.g. a diameter less than 3 μm). Coating the surface with a monolayer of caesium gives electron-emission efficiencies up to 5%. The cold cathodes emit hot electrons; the effective electron temperature may be as high as 5700 K. Current densities and brightness values can be very high: 1500 A/cm² and 9×10^6 A/cm²sr at a voltage of 10 kV. The total emission current can be increased by increasing the emitting area, but at the expense of current density. Besides their very high current density and brightness, silicon cathodes have other desirable features, e.g. for application in cathode-ray tubes: they are easily and rapidly modulated, they use little energy and are easy to manufacture. The prospects for long useful life look promising.

Applications of light guides in process control

The automation of an industrial process will often improve the reproducibility of the end-product. If at the same time the characteristics of the product can be brought closer to the desired characteristics, then automation has also improved the quality. The control unit must of course be provided with the appropriate process data.

Sometimes such data can be collected by using *light*, especially if they relate to quantities such as pressure, temperature, coating thickness or the velocity and direction of a movement, which modify some parameter of the light such as the intensity. It is usually easy to express pressure and temperature in terms of a mechanical displacement. A temperature change, for example, can be measured by determining the change in the length of a glass or metal rod. (It may be that light is emitted in the process. In pyrometry, for instance — where the temperature of an incandescent object is determined from its colour — an optical fibre can be used to conduct the light to a more accessible location for the pyrometer, and the measurement can be made there.)

The measured data are then transferred from the process to the process controller by optical fibres or rods, or 'light guides'. The use of light and light guides has a number of advantages. Measurements can be made at high temperatures, and since there are no electrically conducting materials, measurements can be made in environments such as a chemical reactor heated by an r.f. generator, or in a location where strong electromagnetic fields, e.g. from welding equipment, could upset conventional measurements. Generally speaking, no problems will be encountered in atmospheres containing highly flammable, explosive or corrosive substances.

The term optical fibres usually refers to the thin glass fibres of the type used in optical communications [1]. Philips Research Laboratories, however, can now make light guides with a much larger diameter, ranging from 50 μm to 10 mm. Some of the techniques used for making them go back 20 or 30 years, but of course in that time the tolerances have been reduced from a few millimetres to a few microns.

In this article we shall look briefly at two applications of light guides in industrial process measurements. The first is a displacement measurement in which optical fibres are used. In the second, thicknesses of coatings are measured with the aid of much thicker light guides.

Fig. 1a is a very schematic representation of a Michelson interferometer in which the light paths are not straight lines but are formed by highly flexible single-mode optical fibres [2]. In the present version of the interferometer the beam is split by a beam-splitter that consists of a solid housing containing four ball lenses and a partially transmitting mirror. The ball lenses collimate the light emerging from an optical fibre, or focus a collimated beam on to the end of a fibre. A beam-splitter of this type requires extremely accurate alignment of the components.

The interferometer works as follows. The light beam from the laser is divided into two separate beams in the beam-splitter. One of the beams goes to a fixed reference plane, and the other to the plane whose position is to be measured. The beams reflected from these two planes arrive, at least in part, at the detector, where they interfere. For interference to occur the light in the entire system must have a defined phase. It is therefore necessary to use a single-mode optical fibre, since all the modes in a multimode fibre have a different phase velocity, and 'the' phase of the

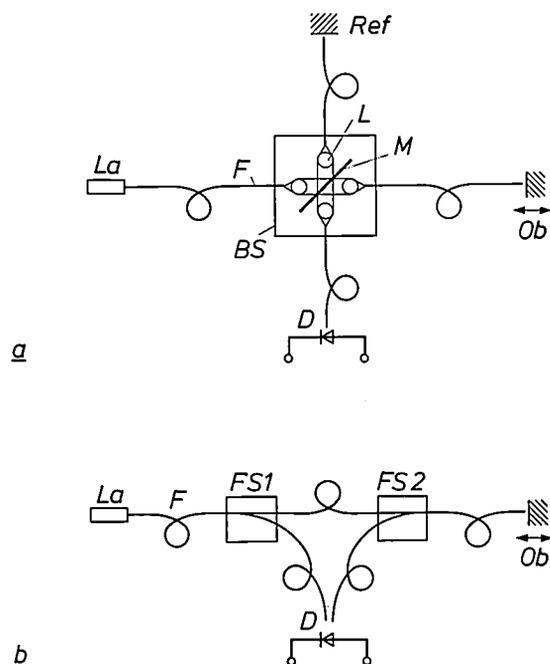


Fig. 1. Schematic diagram of a Michelson interferometer in which the light path is determined by optical fibre. *a*) Version with a conventional beam-splitter. *La* laser. *D* detector. *Ref* reference plane. *Ob* object whose movement is measured (in the direction indicated). *F* fibres. *BS* beam-splitter. *b*) Version with two fibre-splitters *FS1* and *FS2*. Other symbols as under (*a*).

light is not defined over the whole cross-section of the fibre.

When the plane whose position is to be measured is displaced, the detector 'sees' successive maxima and minima in the intensity of the incident light. The displacement Δl required for going from one maximum to another is $\Delta l = m\lambda/2$; where λ is the wavelength of the light and m is an integer. The distance between a maximum and an adjacent minimum is therefore $\lambda/4$. By counting m it is thus possible with an interferometer to measure very small displacements. The great advantage of this interferometer is that the optical fibres that guide the light beams are very flexible, and their length can be chosen as required. (In our case the distance between the measurement position and the processing device was about 3 m.) The only limitation applies to the difference in length of the two fibres between the beam-splitter, the reference plane and the object. This difference must be smaller than the coherence length of the laser light — or there will be no interference. For this application we used a simple HeNe laser with a coherence length of a few centimetres. The use of optical fibres enables us to choose a light path such that measurements can be made on moving parts at places that would be inaccessible to a conventional interferometer. The interferometer described here has been used, for example, for examining ultrasonic welds of gold wires for IC manufacture. The displacements are such that a large number of maxima and minima ($m \approx 20$) can be observed in an ultrasonic period, which may moreover be fairly short (15-20 μ s). Processing the measurement data is therefore no easy matter.

As we noted earlier, the version of the beam-splitter described above requires extremely accurate alignment of the components. A much simpler version can be obtained by making use of 'fibre splitters'. In these components the beam is split by a special technique in which the actual splitting takes place in the glass fibre. When these components are more widely available, the version in fig. 1b will be preferable. In the first fibre splitter (FS1) the light is split into a reference beam and a beam that goes to the object, where it is reflected. In its turn the reflected light is split in FS2. Some of this light is fed to the detector, where it interferes with the reference beam. This version eliminates the tedious process of alignment. The various components can easily be interconnected by optical fibres.

Another example of an unusual application of light guides is the coating-thickness monitor, which is used for measuring the thickness of an infrared-reflecting coating while it is being applied to the inside of the outer glass envelope of low-pressure sodium lamps. This coating, typically of $\text{In}_2\text{O}_3:\text{Sn}$, transmits the

visible light but reflects the infrared. In this way the temperature required for optimum sodium pressure in the discharge tube is maintained [3]. This coating is applied at a temperature of about 500 °C.

The monitor (fig. 2a) consists of two light guides with an outer diameter of 2.5 mm and a core diameter of 1.9 mm. They are about 1 m long, with a temperature difference of about 500 °C between the measurement end and the other end, which is at room temperature. At one end of each light guide there is a connector that couples the light guide to a 'bundle' of optical fibres. These bundled fibres provide the connections to the light source and the processing equipment,

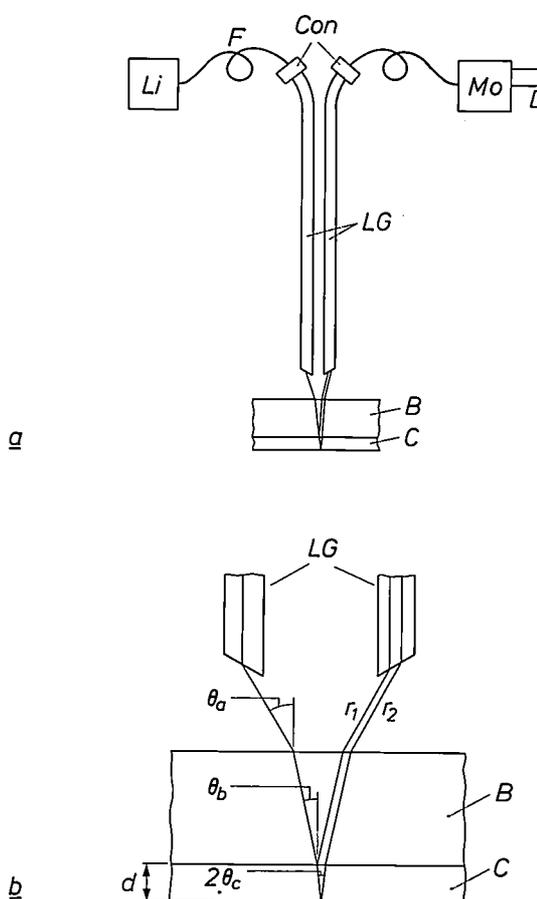


Fig. 2. a) Diagram illustrating the measurement of the thickness of an infrared-reflecting coating in a low-pressure sodium lamp. *Li* light source. *F* optical fibres. *Con* demountable connectors. *Mo* monochromator. *D* detector. *LG* light guides (diameter 2.5 mm) with beveled end faces. *B* glass envelope of the lamp. *C* infrared-reflecting coating. b) Detail of the path of rays (not to scale). θ_a , θ_b , θ_c angle of light beam to the normal in air, glass and coating. r_1 , r_2 beams reflected by the coating. d) thickness of the coating. Other symbols as in (a).

[1] K. Mouthaan, Optical communication systems with glass-fibre cables, Philips Tech. Rev. 36, 178-181, 1976.

[2] Some differences between single-mode and multimode optical fibres are mentioned in the article by A. J. A. Nicia, Philips Tech. Rev. 42, 245-261, 1986.

[3] See pages 226-231 of H. Köstlin and G. Frank, Thin-film reflection filters, Philips Tech. Rev. 41, 225-238, 1983/84.

which is located a few metres away. The light from a light source is conducted to the measurement location by one of the light guides. The emergent beam is reflected from the coating being measured, picked up by the other light guide and conducted to the detector. To ensure that the optical axes of the two parallel light guides are directed towards the same point on the coating, the end faces of the guides are bevelled, as shown in fig. 2.

The principle of the measurement again depends on interference. Since the light beam is reflected from both sides of the selective reflecting coating, two reflected beams are obtained, which interfere because of the difference in optical pathlength (fig. 2*b*). If the thickness of the coating increases, the detector will again detect variations in the intensity of the incident light. Snell's law states that

$$\sin \theta_c = \frac{n_a}{n_c} \sin \theta_a,$$

where θ is the angle between the light beam and the normal, and n is the refractive index. The subscripts a and c indicate that the quantity relates to air or the infrared-reflecting coating. Since $n_a = 1$, then for small values of θ :

$$\cos \theta_c \approx 1 - \frac{\sin^2 \theta_a}{2 n_c^2}.$$

It can be seen from fig. 2*b* that the optical path differ-

ence l between the two reflected beams is $2n_c d / \cos \theta_c$. Bearing in mind that maxima in the intensity will occur at values of d at which $l = m\lambda$, and using the expression above for $\cos \theta_c$, we find that the increase Δd in the thickness of the coating for the detector signal to go from one maximum to the next (i.e. when m increases by one), is given by

$$\Delta d = \frac{\lambda}{2 n_c} \left(1 - \frac{\sin^2 \theta_a}{2 n_c^2} \right).$$

The distance between a maximum and an adjacent minimum is half this value, of course. In practice it is equal to about $\lambda/6$.

The method described here has already been used successfully for a few weeks in the production of low-pressure sodium lamps. The development of a new version suitable for measuring temperatures up to 800 °C is now in an advanced stage. This version will make it possible to measure the thickness of coatings applied by certain CVD (chemical vapour deposition) processes.

P. J. W. Severin
A. P. Severijns

Dr P. J. W. Severin and A. P. Severijns are with Philips Research Laboratories, Eindhoven.

1937

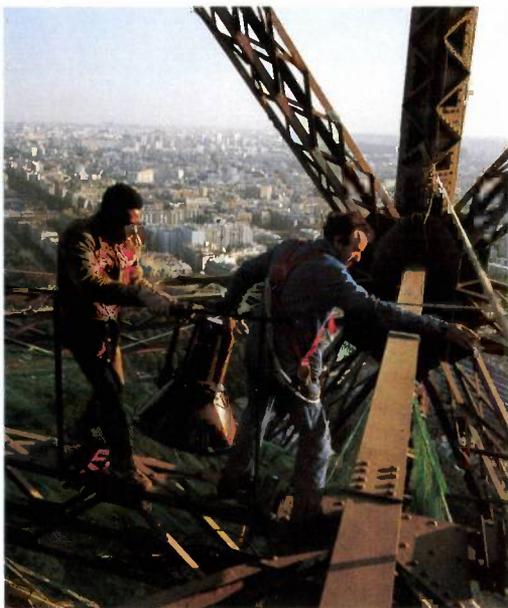
THEN AND NOW

1987

Lighting the Eiffel Tower

Philips have long played a part in the lighting of all kinds of structures. One particularly fine example is the Eiffel Tower in Paris, which was given a highly up-to-date lighting system in 1937 at the time of the Paris International Exhibition of Arts and Technology. A number of searchlights were included as a special feature. These had Philips high-pressure mercury lamps with water cooling. The illustration [*] at the lower right gives an impression of this installation, which was modified many times through the years.

Some time ago Philips installed a completely new lighting system, with nearly 300 fittings, most of them containing a 1000-watt high-pressure sodium lamp (photo lower left). This gives a much better impression of the silhouette and the detail of the construction of the tower (large photo), and reduces the glare for



visitors climbing the tower in the evenings. The new system also cuts the electricity consumption by a half, and maintenance is now very much easier because a computer can identify breakdowns in any section of the system. The illumination is switched on automatically when the daylight level falls below 10 to 15 lux.

[*] From: Philips Technical Review, December 1937.

Laser diagnostics for low-pressure mercury discharges

P. van de Weijer and R. M. M. Cremers

In optical measurement systems the laser is perhaps the most convenient type of light source. The special features of the laser that set it apart from other light sources are the high intensity of laser light and the distance over which the light is coherent. There are many examples of the utilization of these special features. The high value of the coherence length enables distances to be determined extremely accurately, as in the COLATH lathe and the Compact Disc player^[]. The high intensity in a small wavelength range is very useful in spectroscopy. The authors of the article below have used a laser for determining the lifetime and density of excited mercury atoms in a low-pressure mercury-vapour discharge — the origin of the light in the well-known fluorescent lamp.*

Introduction

The principal application of low-pressure mercury-vapour discharge is in fluorescent lamps. The most familiar example is probably the tubular fluorescent lamp. Smaller types, such as the Philips SL* and PL* lamps, have also been on the market for some time. Although these lamps are relatively expensive, their advantages compared with incandescent lamps are their longer life and their higher light output per unit of energy.

A fluorescent lamp consists of a glass tube filled with a mixture of mercury vapour and an inert gas^[1]. The mercury vapour pressure is determined by the temperature of the tube wall and is about 1 Pa. The inert-gas pressure is a few hundred Pa. When an electric current is passed through this gas mixture the mercury atoms can be ionized by collisions with electrons and thus maintain the discharge, or they can be raised to an excited state and then decay to the ground state, accompanied by the emission of radiation. In a low-pressure mercury-vapour discharge lamp this process produces ultraviolet radiation, which is converted into visible light by a fluorescent powder on the tube wall (hence the name fluorescent lamp). The inert gas is not excited by collisions with electrons, or only slightly so, since the primary excitation energy of the inert gas is much higher than that of mercury. The inert gas therefore makes no direct contribution to the light output in fluorescent lamps; its function is to

reduce the mean free path of the electrons. If there were no inert gas in the discharge, the mean free path of the electrons would be very large, and most of the electrons would lose their energy to the wall of the tube before they could transfer it to the mercury atoms.

The radiation output of the discharge depends mainly on the number of mercury atoms that can be excited by collisions with electrons. The conditions in the discharge must be such that these excited atoms can lose their energy as radiation and not via mechanisms such as collision with other particles. Important parameters in these processes are the radius of the discharge tube, the densities of the mercury vapour and the inert gas, and the magnitude of the current through the discharge. Models based on fundamental processes in these gas discharges have been developed for calculating the effect of the parameters on the radiation output of the discharge^[2]. The excitation and decay mechanisms of the mercury atoms play an important part in these models. Methods of determining the density and lifetime of excited mercury atoms have been studied at Philips Research Laboratories. The object of these measurements is to test the models and to improve them.

A mercury atom in the ground state (6^1S_0) has two electrons with opposite spins in the outer shell, the 6s shell. As a result of collisions with electrons in the dis-

[*] T. G. Gijsbers, COLATH, a numerically controlled lathe for very high precision, Philips Tech. Rev. 39, 229-244, 1980; special issue 'Compact Disc Digital Audio', Philips Tech. Rev. 40, 149-180, 1982.

Dr P. van de Weijer is with Philips Research Laboratories, Eindhoven; Ing. R. M. M. Cremers, formerly with these Laboratories, is with the Glass Main Supply Group, Philips NPB, Eindhoven.

charge one of the 6s electrons may be excited to a higher energy level. After the ground state, the lowest excited states, the 6P energy levels, have the highest population in the discharge. The radiative transitions from the 6^1P_1 and the 6^3P_1 levels to the ground state, at wavelengths of 185 nm and 254 nm respectively, form the main contributions to the radiation output of the mercury discharge. Fig. 1 shows the energy-level diagram with the energy levels and radiative transitions that are significant in our investigation. The transitions from the 6^3P_0 and 6^3P_2 levels to the ground state are optically forbidden. These levels are called metastable levels. Although these metastable levels make no direct contribution to the radiation output, they do play an important part in the experiments described in this article.

The photons generated in the permitted transitions to the ground state may be absorbed again by a mercury atom in the ground state and will thus not contribute immediately to the radiation output. This process is known as radiation trapping or imprisonment of radiation. The lifetime of an excited state that is of importance for the radiation output is therefore not the natural lifetime. The important parameter is the time that elapses between the moment when the excited state is reached and — after repeated absorption and emission of the UV photon — the moment at which the UV photon reaches the wall. At this moment the fluorescent powder on the wall emits a visual photon. This time is called the effective radiation lifetime (τ). As this effective lifetime increases, the probability that the atom in the excited state will lose its energy through collisions with other particles also increases. This is why the effective lifetime of the excited state — as well as the density of the excited mercury atoms — has an important bearing on the radiation output of the lamp. To measure these two quantities, lifetime and density, we have made use of a laser.

The energy level whose lifetime we want to measure is populated by means of a laser pulse via a higher level. The time-dependent fluorescence signal, which is a result of the transition from the level of interest to the ground state, then gives the effective lifetime we wish to know. Compared with the usual method, which consists in directly populating the level of interest from the ground state, this method has the advantage that the measurement of the fluorescence signal is not perturbed by scattered laser light, since the wavelength of the radiation used for populating the level is different from that of the fluorescence radiation.

The density of the atoms in an excited state can be determined by using the 'hook' method. Interference patterns in the neighbourhood of an absorption line are measured, giving the refractive index as a function

of wavelength. This can be used to determine the density of the atoms at the lower energy level of the transition.

This method has certain advantages compared with the more conventional absorption method. The shape and width of the line do not have to be known, and the dynamic range is greater than that of the absorption method. On the other hand, the absorption method is somewhat more sensitive, which implies that absorption measurements give more accurate results than the hook method at low densities.

If the density at a particular energy level and the effective lifetime of that level are known, the radiation output of that level — the ratio of the two quantities — can be calculated and compared with the directly measured output. This comparison provides a check on the reliability of the different experiments involved in this procedure.

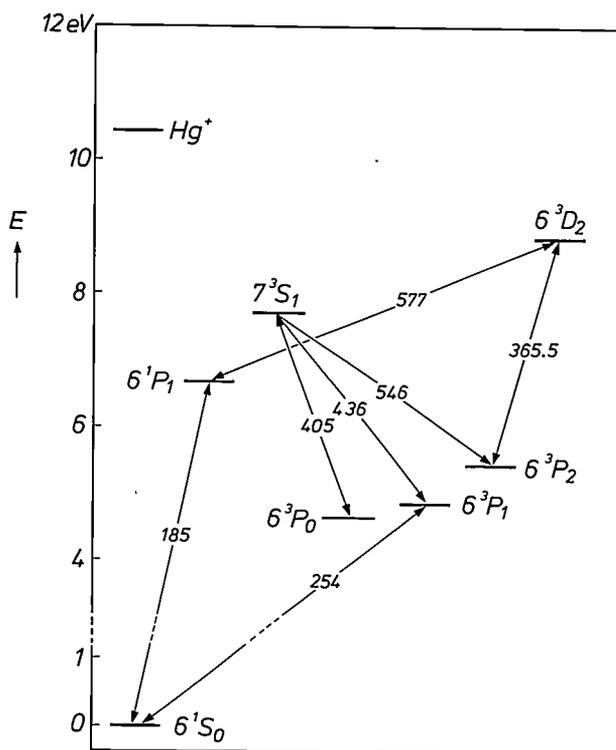


Fig. 1. Energy-level diagram of a mercury atom; only the levels of interest in the present study are shown. The transitions that make the major contributions to the radiation output of the low-pressure mercury-vapour discharge are those from the 6^3P_1 and the 6^1P_1 levels to the ground state 6^1S_0 . This is ultraviolet radiation with wavelengths of 254 nm and 185 nm respectively. The other transitions play a part in the determination of the lifetime and density of mercury atoms at the 6^3P and 6^1P_1 levels. The wavelengths are indicated in nm.

[1] In the TL lamp type this is a straight cylinder, in an SL* lamp it is a doubly folded cylinder, while a PL* lamp consists of two parallel cylinders connected at one end.

[2] The model we used is based on the work of M. A. Cayless (Br. J. Appl. Phys. 14, 863-869, 1963). The present version of the model was developed by F. A. S. Ligthart and D. B. de Mooij (Bull. Am. Phys. Soc. 24, 126, 1979).

We shall first discuss the determination of the effective lifetime of the two principal excited mercury levels (6^3P_1 and 6^1P_1), and then we shall consider the hook method for determining the excited-state density and our proposed improvements. Finally, the radiation output of the discharge calculated from our measured quantities will be compared with the experimental data.

Lifetime measurements

The lifetime of mercury atoms at the 6^3P_1 level

The radiative transition from the 6^3P_1 level to the ground state ($\lambda = 254$ nm) makes the most important contribution to the radiation output of the low-pressure mercury/inert-gas discharge. For the effective lifetime of this level to be determined it must first be populated. Normally this is done by irradiating the mercury atoms in a vapour cell with periodically interrupted 254-nm radiation from a low-pressure mercury discharge. The intensity of the fluorescent radiation, which also has a wavelength of 254 nm, is recorded as a function of time, and this gives the required lifetime^{[9]-[7]}. This procedure has some disadvantages, however. The absorption of the 254-nm radiation in the mercury-vapour cell is so high that this radiation can only excite atoms in a small part of the cell. The excitation profile therefore differs from the 6^3P_1 density profile in a discharge lamp. This need not be such a difficulty, however, because the theory indicates that the variation of the fluorescence signal can be described after some time by a time constant that corresponds to the lifetime of the level^[8]. One of the reasons for making a further study of the lifetime of this level was that we wanted to see whether the shape of the excitation profile really does not affect the value of the lifetime. An advantage of our method of measurement is that excitation and fluorescence radiation have different wavelengths, so that it is relatively easy to prevent scattered excitation radiation from affecting the fluorescence signal.

Another reason for our study was that earlier measurements had been carried out on mercury vapour in a cell, whereas we are also interested in mixtures of mercury vapour and inert gases, because of their application in fluorescent lamps.

The principle of our method is as follows. A low-pressure mercury/inert-gas mixture is irradiated for 10 ns with a laser pulse of 405 nm wavelength. This has the effect of exciting mercury atoms from the metastable 6^3P_0 level to the 7^3S_1 level (fig. 1). The atoms decay from this level to the 6^3P levels with a known time constant of 8 ns. The result is that the 6^3P_1 level becomes temporarily overpopulated. The

time-variation of the decay of atoms from this level to the ground state can be determined from the variation of the fluorescence signal. This signal is recorded with a photomultiplier and an integrator. Since excitation and fluorescence radiation have different wavelengths, the scattered laser light can be separated by a filter from the fluorescence radiation to be detected. To obtain the most uniform irradiation of the discharge, we made the diameter of the laser beam larger than the diameter of the discharge tube. There is very little decrease in the intensity of the laser beam after it has passed through the discharge tube. This is due not only to the greater intensity of the laser compared with the excitation radiation of a discharge lamp, but is also a consequence of the fact that the atoms are now not excited from the ground state but from the less densely populated 6^3P_0 level. There is much less absorption at 405 nm than at 254 nm. In this way the discharge is irradiated as uniformly as possible and as a result the excitation profile in the 6^3P_1 level is an exact copy of the radially symmetrical density profile at the 6^3P_0 level.

This method does have a disadvantage, however, compared with the method mentioned earlier. For atoms to be available at the 6^3P_0 level the measurements have to be made in a discharge. But the discharge current has to be very low, to prevent atoms at the 6^3P_1 level from losing their energy except through non-radiative transitions (e.g. by collisions with electrons). An incidental advantage of this low discharge current is that the laser-induced fluorescence signal from the 6^3P_1 level is many times larger than the emission signal caused by the discharge itself.

This effect follows from the fact that the density in the 6^3P_1 level and in the metastable levels^[9] depends in a different way on the discharge current. At high currents the densities of the atoms at the 6^3P levels are comparable because they are coupled by exciting and de-exciting collisions with electrons. At low currents, i.e. at low electron densities, this coupling is absent. At this low electron density the metastable levels are more densely populated than the 6^3P_1 level, because atoms at these levels can only lose their energy by diffusion to the wall, which is a much slower process than a radiative transition — the way in which atoms at the 6^3P_1 level lose their energy. The result is that the number of mercury atoms transferred from the metastable 6^3P_0 level via the 7^3S_1 level to the 6^3P_1 level is large compared with the number of mercury atoms already present at the 6^3P_1 level in the steady-state discharge.

To test our method of measurement we first of all determined the natural lifetime of the 6^3P_1 level in a pure mercury discharge. When the vapour density of the mercury is made low, by keeping the wall temperature low, no radiation trapping occurs and the time

constant of the fluorescent radiation is the natural lifetime of that level. The result of our measurements, 120 ± 2 ns, is in agreement with the results of others (values reported in the literature show a spread from 114 to 122 ns ^{[5]-[7]}). If we increase the wall temperature of the discharge tube, thus increasing the mercury vapour pressure, radiation trapping occurs and we measure the effective lifetime of the level. This effective lifetime is determined not only by the mercury pressure but also by the dimensions of the tube. If we take the discharge tube to be an infinitely long cylinder, then the effective lifetime of the atoms in an excited state is a function of the tube radius. In *fig. 2* the effective lifetime that we measured is plotted as a function of k_0R , where R is the radius of the discharge tube and k_0 the absorption coefficient for the 254 nm radiation, a measure of the mercury vapour pressure. This figure also shows the effective lifetimes measured ^{[3]-[7]} and calculated ^{[10][11]} by others. The agreement between the results of the different experiments is good. In the older experiments an asymmetric excitation profile was used, whereas we took particular care to ensure that the excitation profile corresponded as closely as possible to the symmetrical 6^3P_1 density profile in a fluorescent lamp. The good agreement achieved in spite of this difference seems to confirm the correctness of the theory that predicts that the measured effective lifetime will be independent of the excitation profile.

Seeking further support for this conclusion we carried out a number of experiments with local excitation profiles (e.g. using narrower laser beams), but none of these changes produced any essential change in the results. Differences between the results were no more than 10%, which is within the experimental error. It

was the first time that this theoretical prediction had been supported by experimental data.

The next step was to measure the effective lifetime of the 6^3P_1 level in mercury/argon and mercury/krypton discharges. The presence of the inert gas shortens the effective lifetime; the wings of the 254-nm absorption line are broadened by collisions between mercury and inert-gas atoms, enabling the radiation to escape

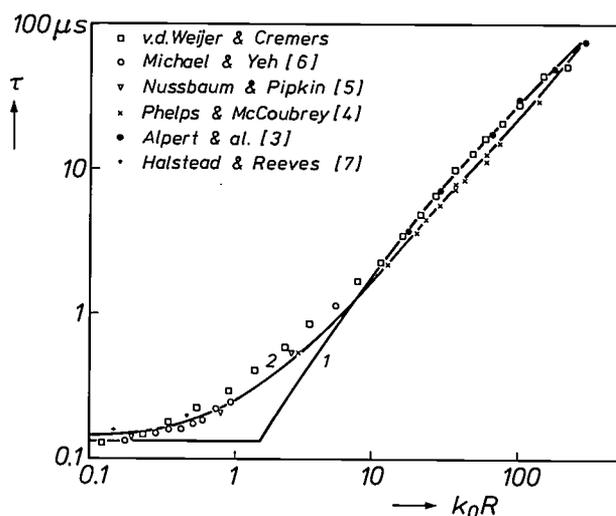


Fig. 2. The effective radiation lifetime (τ) of the 6^3P_1 level of mercury atoms in a mercury discharge (tube diameter 13 mm, length 500 mm, wall temperature varying from -15°C to $+84^\circ\text{C}$) as a function of k_0R (R is the radius of the discharge tube and k_0 the absorption coefficient ^[10], which at low pressure is proportional to the mercury vapour pressure). The symbols indicate the measured results, the continuous curves represent the calculated values. Curve 1 was calculated from a theory that is only valid for large values of k_0R ^[10], curve 2 from a theory that is valid for all optical depths ^[11].

[3] D. Alpert, A. O. McCoubrey and T. Holstein, Imprisonment of resonance radiation in mercury vapor, *Phys. Rev.* 76, 1257-1259, 1949.

[4] A. V. Phelps and A. O. McCoubrey, Experimental verification of the 'incoherent scattering' theory for the transport of resonance radiation, *Phys. Rev.* 118, 1561-1565, 1960.

[5] G. H. Nussbaum and F. M. Pipkin, Correlation of photons in cascade and the coherence time of the 6^3P_1 state of mercury, *Phys. Rev. Lett.* 19, 1089-1092, 1967.

[6] J. V. Michael and C. Yeh, Absolute cross section for $\text{Hg}(^3P_1) + \text{H}_2$ and the imprisonment lifetimes for $\text{Hg}(^3P_1)$ at low opacity, *J. Chem. Phys.* 53, 59-65, 1970.

[7] J. A. Halstead and R. R. Reeves, Determination of the lifetime of the mercury 6^3P_1 state, *J. Quant. Spectrosc. & Radiat. Transfer* 28, 289-296, 1982.

[8] C. van Trigt, Analytically solvable problems in radiative transfer. IV, *Phys. Rev. A* 13, 726-733, 1976.

[9] M. Koedam and A. A. Kruithof, Transmission of the visible mercury triplet by the low pressure mercury-argon discharge; concentration of the 6^3P states, *Physica* 28, 80-100, 1962.

[10] P. J. Walsh, Effect of simultaneous Doppler and collision broadening and of hyperfine structure on the imprisonment of resonance radiation, *Phys. Rev.* 116, 511-515, 1959.

[11] A. V. Phelps, Effect of the imprisonment of resonance radiation on excitation experiments, *Phys. Rev.* 110, 1362-1368, 1958.

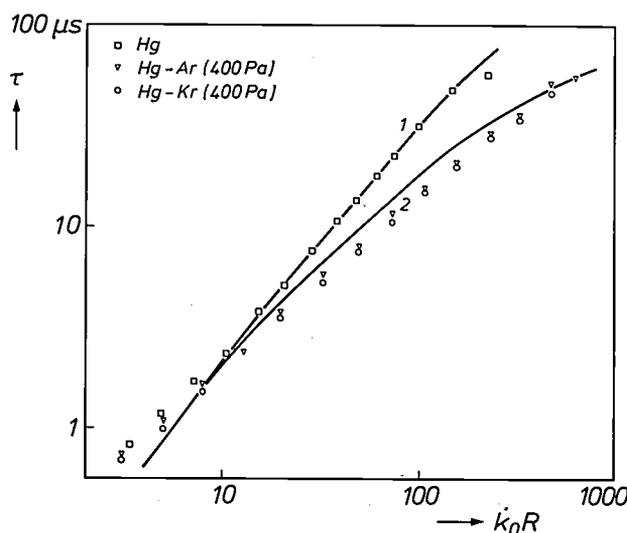


Fig. 3. The effective radiation lifetime (τ) of the 6^3P_1 level of mercury atoms in a mercury discharge, a mercury/argon discharge and a mercury/krypton discharge as a function of k_0R . The pressure of the inert gas is 400 Pa in both cases and the diameter of the tube is 36 mm. The continuous lines relate to the lifetime calculated from the theory due to P. J. Walsh ^[10], for a mercury discharge (1) and for a mercury/argon discharge (2).

more easily from the discharge. At low mercury pressures, when little radiation trapping occurs, the presence of the inert-gas atoms has little if any effect, because the radiation output is determined mainly by the centre of the line. At high mercury vapour pressures the radiation does not escape so easily from the centre of the line, and the wings of the line also make an increasing contribution to the radiation output. For this reason the effect of the inert gas will become greater with rising mercury vapour pressure. Fig. 3 shows the effective lifetime for both discharges as a function of k_0R . The data for the mercury/argon and the mercury/krypton discharges differ from each other by no more than the experimental error of 10%, which means that the line broadening caused by collisions with krypton atoms is the same as that caused by collisions with argon atoms. This was not previously known. Comparison of the results with the theoretical curves does not come out so favourably as it does in the case of a pure mercury discharge, but the difference is nowhere greater than 30%. In the range of direct relevance to fluorescent lamps — effective lifetime 1 to 3 μs — the difference between theory and experiment is much less than 30%. Here, therefore, the theory gives a reasonable description of the effective lifetime of atoms at the 6^3P_1 level.

The lifetime of atoms at the 6^1P_1 level

Until recently no experimental or theoretical data about the effective lifetime of the 6^1P_1 level were available. Model calculations used the theory on which the calculation of the lifetime of the 6^3P_1 level is based. It is assumed in that theory, however, that in a transition of atoms from the ground state to an excited state and the subsequent decay to the ground state after the natural lifetime, the photons have wavelengths covering the whole width of the absorption profile (complete spectral redistribution [12]). This assumption is satisfied for the 6^3P_1 level (natural lifetime 120 ns) but not for the 6^1P_1 level (natural lifetime 1.3 ns). As a result of this assumption, which is incorrect for the 6^1P_1 level, the calculated lifetime is too short. This can be explained in the following way. The probability that a photon will be absorbed is highest at the centre of an absorption line. If the emitted photon has a wavelength close to the centre of the line — within the Doppler width — there is a high probability that this photon will also be absorbed. If complete spectral redistribution occurs, a photon that is absorbed at the centre of the line can be emitted in the wings, where the probability of escape is greater.

Until recently only one experiment for measuring the effective lifetime of atoms at the 6^1P_1 level was known [13]. A dye laser tuned to a wavelength of

577 nm was used to measure absorption at the 6^1P_1 - 6^3D_2 transition in a low-pressure mercury discharge (fig. 1). The density at the 6^1P_1 level can be determined from the absorption. The lifetime of the 6^1P_1 level can be found from the variation of the density as a function of time in the afterglow of the discharge [14]. This is only possible when the discharge current is so low that only radiative decay takes place from the 6^1P_1 level to the ground state. An additional problem arises here, however. The effective lifetime can only be determined from these absorption measurements when the 6^1P_1 level in the steady-state discharge is populated largely by collisions between mercury atoms

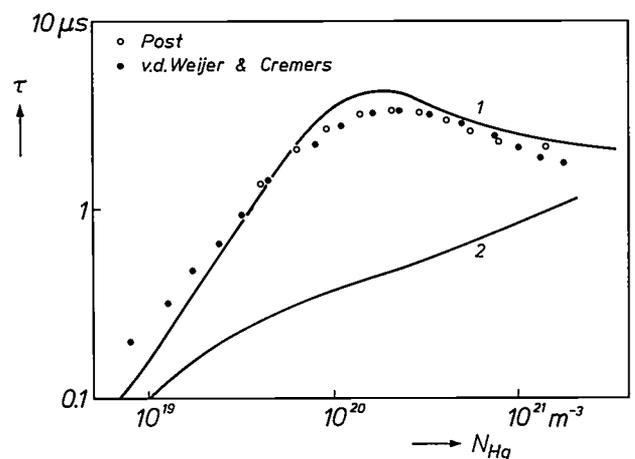


Fig. 4. The effective radiation lifetime (τ) of the 6^1P_1 level of mercury atoms in a mercury discharge as a function of the mercury vapour density N_{Hg} . The diameter of the discharge tube is 25 mm. The mercury vapour densities correspond to wall temperatures of 0 °C to 80 °C. In this figure the results of the afterglow experiments [13] are compared with those of our measurements. Curves 1 and 2 are calculated lifetimes. Complete spectral redistribution is assumed for curve 2, but not for curve 1.

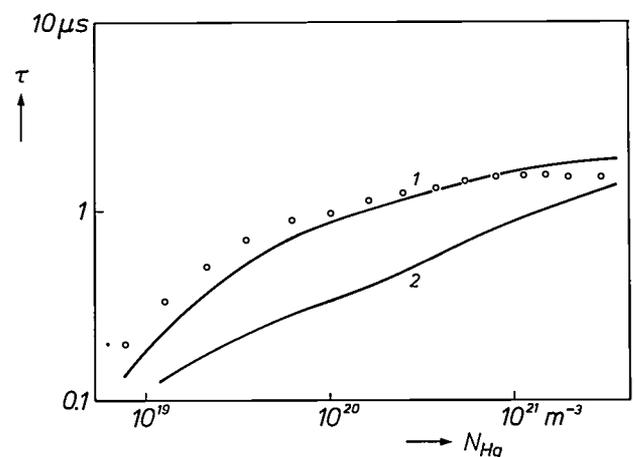


Fig. 5. The effective radiation lifetime (τ) of the 6^1P_1 level of mercury atoms in a mercury/argon discharge with a diameter of 36.5 mm, as a function of the mercury vapour density N_{Hg} . The circles indicate our experimental results, the continuous lines 1 and 2 are calculated lifetimes [13]. For the calculation of curve 2 complete spectral redistribution was assumed, but not for curve 1.

in the ground state and electrons. In general this is true in the case of low-pressure mercury discharges at low current densities, but not in the case of low-pressure mercury/inert-gas discharges. In such discharges this level is populated by a two-step excitation from the ground state via the 6^3P levels. The variation of the density at the 6^1P_1 level in the afterglow of the discharge is then completely determined by the relaxation times of the 6^3P levels, which are much longer than the effective lifetime of the 6^1P_1 level.

We measured the effective lifetime of the 6^1P_1 level in a low-pressure mercury/inert-gas discharge in the same way as for the lifetime measurements for the 6^3P_1 level. We irradiated a low-pressure mercury/inert-gas discharge with a laser pulse of 10 ns duration and 365.5 nm wavelength, thus causing a transition from the metastable 6^3P_2 level to the 6^3D_2 level. From this level there is radiative decay to the 6^1P_1 level, producing temporary overpopulation there. The 6^1P_1 density in the steady-state discharge is orders of magnitude smaller than the densities in the 6^3P levels. The laser pulse therefore causes a considerable overpopulation of the 6^1P_1 level. The decay of the 185-nm fluorescent radiation again gives the effective lifetime of that level.

As a check we determined the effective lifetime in a low-pressure mercury-vapour discharge. Fig. 4 shows our results and the experimental results of H. A. Post^[13] as a function of the mercury vapour density. Unlike the lifetime of atoms at the 6^3P_1 level, the lifetime of atoms at the 6^1P_1 level has a maximum as a function of the mercury vapour pressure. This points to a difference in the mechanisms of radiation trapping in the two transitions. Comparison of the experimental results with values calculated from the theory (curves 1 and 2) in fig. 4 shows that in fact no complete spectral redistribution takes place. For the calculation of curve 2 a complete spectral redistribution was assumed, but not for the calculation of curve 1. Good agreement exists between the calculated and the measured lifetimes of atoms in the 6^1P_1 level in a mercury-vapour discharge. We next applied our method of measurement to a mercury/argon discharge (see fig. 5). Here again the measured relation between mercury pressure and the lifetime of atoms in the 6^1P_1 level is described well by Post's theory in which no complete spectral redistribution is assumed.

Density measurements

The hook method

We determined the density in the excited 6P levels in low-pressure mercury-vapour discharges by means of the 'hook' method. This method was introduced in

1912 by D. Roschdestwensky^[15]. The principle of the method is as follows. A monochromatic beam of light passes through an interferometer consisting of two beam splitters and two mirrors (fig. 6a). The beam is split by the first beam splitter into two beams of equal intensity, which are made to coincide again at the second beam splitter by the two mirrors. If there is no difference in pathlength between the two beams, no interference pattern appears. If the second beam splitter is at a small angle ϕ to the other mirrors, an optical path difference is introduced that causes an interference pattern to appear on the screen (fig. 6b). There

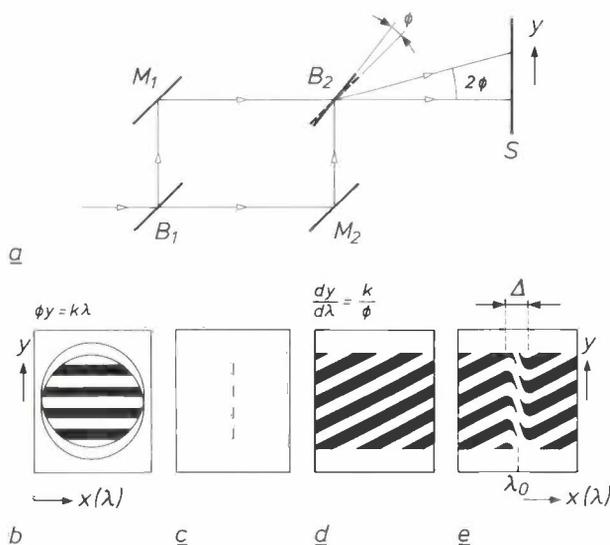


Fig. 6. Diagram showing the principle of the hook-method experiment. a) A monochromatic light source of wavelength λ is used for setting up the interferometer. The light beam is split by the beam splitter B_1 into two beams, which are made to coincide at the beam splitter B_2 by the mirrors M_1 and M_2 . The beam splitter B_2 is at an angle of ϕ to the other mirrors. An interference pattern is produced on the screen S . b) The lines of the pattern on the screen have a position given by $y = (k/\phi)\lambda$, where $k = 1, 2, 3, \dots$ c) When an image of the pattern is produced by a lens that focuses only in the horizontal direction, a line with a (cosine)² intensity distribution appears on the screen. d) The interference pattern shown here occurs when the monochromatic light source is replaced by a broadband light source and the light is dispersed after the interferometer. e) Diagram of a hook spectrum. The distance between the two positions where the slope of the lines becomes zero is Δ .

[12] If there are no collisions between the excited atoms and other atoms, the thermal motion of the atoms may make the wavelength of the emitted photon differ from that of the absorbed photon by no more than an amount equal to the Doppler width. If collisions do occur, which is the case when the natural lifetime is long enough, complete spectral redistribution takes place. This means that the wavelength of the emitted photon may lie anywhere within the absorption profile of the line.

[13] H. A. Post, Radiative transport at the 184.9-nm Hg resonance line I, Phys. Rev. A 33, 2003-2016, 1986.

[14] The afterglow of the discharge is a state in which the voltage across the discharge is switched off. Excitation of higher states via collisions of electrons with mercury atoms in the ground state has almost completely vanished within about a microsecond. After this has taken place the decay mechanisms of the excited levels can be studied.

[15] D. Roschdestwensky, Anomale Dispersion im Natriumdampf, Ann. Phys. 39, 307-345, 1912.

is a maximum in the light intensity when $\phi y = k\lambda$, where y is the position on the screen, k is an integer and λ is the wavelength of the laser light. This interference pattern is focused in one direction such that a line appears on the screen with a $(\cosine)^2$ intensity distribution (fig. 6c). The monochromatic light source is now replaced by a broadband source and a grating is placed between the interferometer and the screen for spectral analysis of the light. An interference pattern now appears on the screen (fig. 6d) with lines (fringes) whose slope ($dy/d\lambda$) is given by

$$\frac{dy}{d\lambda} = \frac{k}{\phi} \quad (1)$$

The lines have different slopes, but if k is large ($k > 1000$) the difference between the slopes of neighbouring lines is so small that they seem to be parallel. Next, we put the discharge under investigation, with refractive index n and length l , into one arm of the interferometer. In the other arm we put a vacuum tube of the same length l . This gives an additional optical path difference of $(n - 1)l$. The slope ($dy/d\lambda$) of the lines in the interference pattern is now given by

$$\frac{dy}{d\lambda} = \frac{k}{\phi} + \frac{l}{\phi} \frac{dn}{d\lambda} \quad (2)$$

The determination of the slope will give us information about the wavelength-dependence of the refractive index. A change in the refractive index will occur in the neighbourhood of an absorption line. The interaction between the radiation and the medium in the discharge tube is dependent on the wavelength of the radiation. If this is a long way from the wavelength at which absorption occurs (λ_0), particles in the medium will absorb little energy and there will be vibrations of almost imperceptible amplitude. As the wavelength approaches the absorption wavelength, the amplitude of these vibrations will increase. It reaches a maximum when $\lambda = \lambda_0$, where all the radiation energy is absorbed by the particles. The vibrating particles affect the propagation velocity of the wave, and this effect can be observed as a change in the refractive index of the medium. The refractive index changes not only when $\lambda = \lambda_0$, but also in the neighbourhood of the absorption line. A mathematical description of this mechanism is given by Sellmeier's equation. Under certain conditions the relation between the wavelength and the refractive index in the neighbourhood of an absorption line is [16]

$$n - 1 = \frac{r_0 N f \lambda_0^3}{4\pi(\lambda - \lambda_0)} \quad (3)$$

where r_0 is the classical radius of the electron, N the density at the lower level of the transition and f the

oscillator strength of the transition. Combining (2) and (3) gives

$$\frac{dy}{d\lambda} = \frac{k}{\phi} - \frac{r_0 N f \lambda_0^3 l}{4\pi\phi(\lambda - \lambda_0)^2} \quad (4)$$

for the slope of the lines.

This slope becomes zero at two values of the wavelength. The position at which this happens is called a 'hook'. The distance between two of these points on opposite sides of an absorption line is called Δ . The density N at the lower level of the transition is then

$$N = \frac{\pi k \Delta^2}{f r_0 \lambda_0^3 l} \quad (5)$$

and can be calculated if k and Δ can be measured in the interference pattern and f is known. The value of k can be determined from the undisturbed part of the interference pattern a long way from the absorption line [17]. An example of an interference pattern obtained with the hook method is given in fig. 7.

A diagram of the arrangement used in our experiments is shown in fig. 8. The light source was a dye laser, pumped by a pulsed nitrogen laser. The spectral bandwidth of the dye laser was 0.01 nm, narrow enough for adjusting the interferometer (see fig. 6b and c). To record the interference spectrum (see fig. 6d and e) we need a broadband light source. We therefore remove the grating of the dye laser, which is responsible for wavelength selection when the laser is used at narrow bandwidths, and substitute a mirror for it. The bandwidth is then determined by the bandwidth of the dye (10-20 nm).

The light source formerly used in the hook-method experiment was a lamp in combination with an interference filter. The bandwidth of this light source is of the order of 10 nm. Since the coherence length of such light sources is small, both arms of the interferometer must be made equal to an accuracy of a few tens of microns when the interferometer is set up (fig. 6b). This obviously made setting up a tedious process.

With the narrow-band dye-laser the two arms of the interferometer only have to be equal to an accuracy of a few centimetres to give the interference pattern of fig. 6b. The introduction of the laser has obviously greatly simplified the hook-method experiment. In recent years there has consequently been a marked increase in the use of this method.

After the broad-band laser beam has passed through the interferometer, a spectral analysis of the light is necessary. In this respect our arrangement differs from that used by other experimenters. Formerly a commercial spectrometer was used for this. We only used a grating, however. The advantage of using a grating is that the optical components can be adjusted in such a way for each range of wavelengths that the angle be-

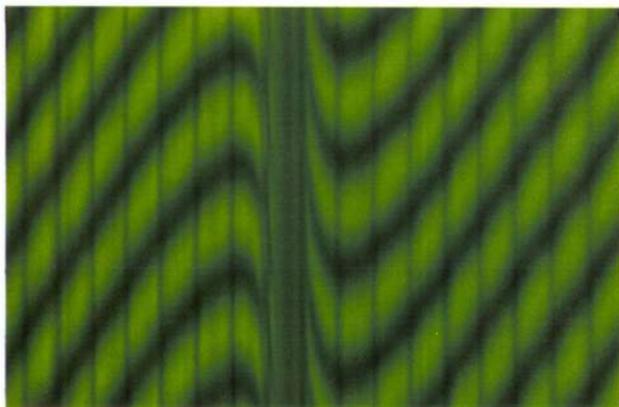


Fig. 7. Example of an interference pattern as observed near the absorption line at 546 nm (the transition from 7^3S_1 to 6^3P_2). The vertical pattern of fringes is due to interferences that occur in the dye laser between two partially reflecting surfaces, and is used for calibrating the wavelength scale. The distance between two vertical fringes corresponds to 0.03 nm.

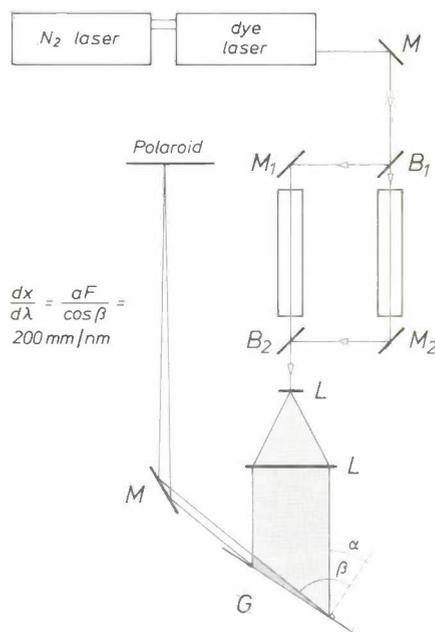


Fig. 8. The arrangement for the hook-method experiment. A dye laser pumped by a pulsed nitrogen laser acts as the light source. The discharge tube is mounted in one arm of the interferometer, and the other arm contains a tube of the same length in which a vacuum is maintained. After the interferometer the two beams are combined to form an image on a two-dimensional detector (a Polaroid camera) by two lenses L (focal length $F = 4$ m) via a grating G ($a = 2400$ lines/mm), which disperses the light. α is the angle between the normal to the grating and the incident beam, β is the angle between the normal to the grating and the diffracted beam^[18]. If this angle is made as close as possible to 90° , the dispersion is at a maximum. The dispersion obtained in this way, $dx/d\lambda = aF/\cos\beta$, is 200 mm/nm. The interference pattern is sufficiently large to be able to establish the distances in the pattern to the required accuracy.

tween the normal to the grating and the diffracted beam (β in fig. 8) is about 90° ^[18]. This gives a high dispersion. We can also use lenses with a large focal length F for focusing the beam, which makes the dispersion even higher. These two factors finally result in a dis-

persion that is generally of the order of 200 mm/nm; see fig. 8. This is almost two orders of magnitude better than the dispersion that can be achieved with a good spectrograph. It is thus no longer necessary for the interference pattern to be enlarged before the required quantities (k , Δ) can be determined with sufficient accuracy. Our arrangement is not only cheaper, but components can be adjusted much more rapidly, since the effects on the interference pattern are immediately visible. A disadvantage is that the dispersion has to be redetermined from the experimental geometry for every wavelength of interest. With β close to 90° , a small inaccuracy in this angle can give a large inaccuracy in the result.

Fortunately, we can turn to another method of determining the dispersion. The vertical interference fringes, which can be observed in fig. 7 and which we first found to be a nuisance, now serve a useful purpose. This vertical interference pattern is the result of interference between two partially reflecting surfaces in the dye laser used as the light source in our experimental arrangement. The distance $\Delta\lambda$ between two vertical fringes is given by

$$\Delta\lambda = \frac{\lambda^2}{2L}, \tag{6}$$

where L is the optical pathlength between the two partially reflecting surfaces. By counting the number of vertical interference fringes between two emission lines of known wavelength, we can determine the distance L . The distance between two vertical lines in the interference pattern is used to calibrate the wavelength scale.

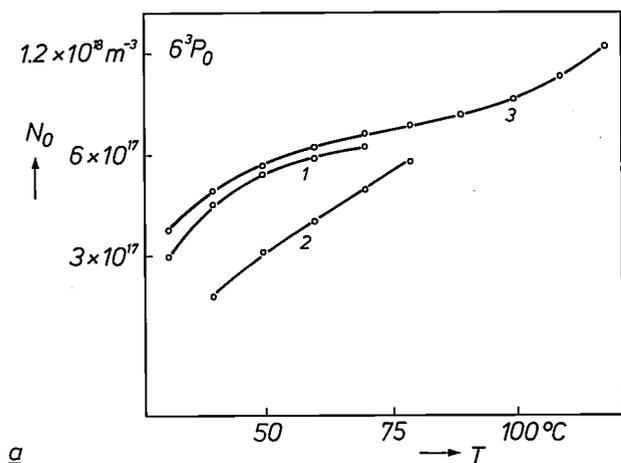
The density of mercury atoms at the 6^3P_1 and 6^1P_1 levels

We used the hook method to determine the densities of atoms in the excited states in a mercury/argon discharge. In fig. 9 the measured values are plotted as a function of the wall temperature, which determines the mercury vapour pressure. In this figure our results for the 6^3P levels are compared with the results of

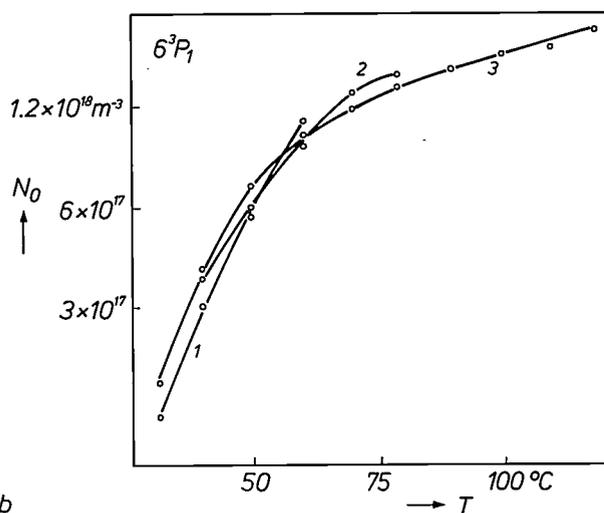
[16] W. C. Marlow, Hakenmethode, Appl. Opt. 6, 1715-1724, 1967. The conditions for the validity of Sellmeier's equation are: $\lambda - \lambda_0$ must be large compared with the absorption width of the line; λ_0 must be remote from all other absorption lines; the density at the lower level of the transition must be large compared with the density at the upper level. These three conditions are fully satisfied in the low-pressure mercury-vapour discharges studied here, so that in our case this method can readily be used for determining the density.

[17] For parallel fringes at λ and $\lambda + \Delta\lambda$, then $k\Delta\lambda = p(\lambda + \Delta\lambda)$, where p is the number of fringes in $\Delta\lambda$. If $\Delta\lambda$ is small compared with λ , then $k = p\lambda/\Delta\lambda$.

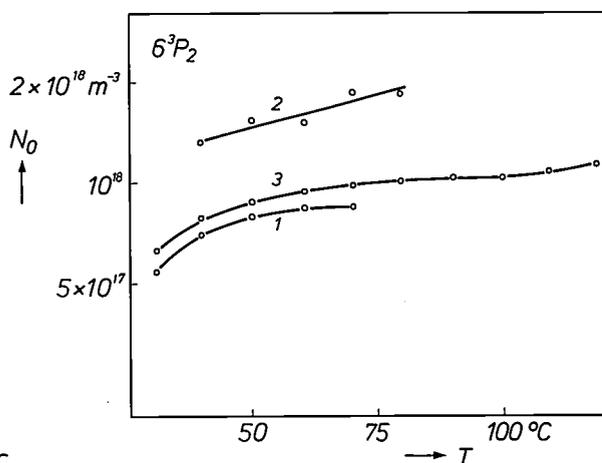
[18] The angle between the normal to the grating and the incident beam is α . In the 'reflected' light beam, maxima are produced because of diffraction. These maxima occur at an angle β to the normal to the grating. β satisfies $(\sin\alpha + \sin\beta) = m\lambda/a$, where $m = 1, 2, 3, \dots$



a



b



c

Fig. 9. The density of mercury atoms at the three 6P levels in a mercury/argon discharge as a function of the wall temperature (diameter of the discharge tube 36 mm, discharge current 400 mA, argon pressure 400 Pa). 1 Absorption measurements by M. Koedam and A. A. Kruithof^[21]. 2 Measurements by F. A. Uvarov and V. A. Fabrikant, made by the hook method^[19]. 3 Our measurements, also made by the hook method.

others^{[9][19]}. The results for the density in the 6^3P_1 level are in good agreement with each other. Strangely, our results for the density in the two metastable levels show more agreement with the absorption measurements than with the measurements made with the hook method. This agreement is in reality not so good as suggested in fig. 9. In both methods, absorption measurements and hook method, the product of the oscillator strength of the transition and the density at the lower level of the transition are determined. The oscillator strengths of the transitions required for calculating the density at the 6^3P levels were not so well known in 1962 as they are now. With the current known values for the oscillator strengths the absorption measurements for the 6^3P_0 and the 6^3P_2 levels come out 15% lower.

The reason for the difference between the results from the two hook-method experiments is not clear. Different values for the oscillator strengths and a greater experimental error of 20% are probably not sufficient to explain the difference of a factor of two in density.

The results for the density at the 6^1P_1 level, given in fig. 10, show an experimental error of 20%. No other measurements of the density at this level are known to us for the discharge conditions for our measurements, so that a comparison of experimental results is not possible. The experimental results can however be compared with densities at that level calculated with the aid of models. This comparison showed a difference that we at first misinterpreted. It turned out later that in calculating the density at the 6^1P_1 level incorrect assumptions had been made about the way in which the effective lifetime of this level had been determined (complete spectral redistribution, see page 66). When the correct effective lifetime of the 6^1P_1 level is used, the calculated densities are in good agreement with the measured densities, as can be seen in fig. 10.

The radiation output of low-pressure mercury-vapour discharges

Now that we know the lifetime and the density of atoms in the two principal excited states, we can calculate the radiation output of the low-pressure mercury discharge. The radiation output of a gas discharge at a particular wavelength can be expressed in terms of the number of photons emitted per unit time and per unit volume by the gas discharge, and is thus equal to the ratio of the mean density — across the diameter of the discharge tube — at the upper level of the transition to its effective lifetime (\bar{N}/τ). We calculated the output at a wavelength of 254 nm as a

function of the wall temperature for two discharges for which directly measured values were also available.

One was a pure mercury-vapour discharge [20]. These measurements and the values that we calculated from the lifetime and density measurements are given in fig. 11a. The agreement between experiment and calculation can be said to be good. Within the experi-

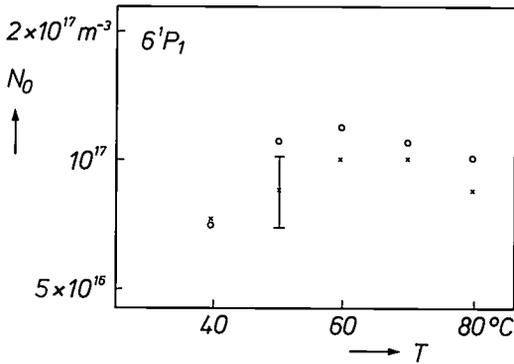


Fig. 10. The density N_0 at the 6^1P_1 level as a function of the mercury vapour pressure in a mercury/argon discharge. The diameter of the discharge tube was 15 mm, the discharge current was 500 mA and the argon pressure was 400 Pa. \circ Calculated densities, \times experimental data. The vertical line indicates the experimental error.

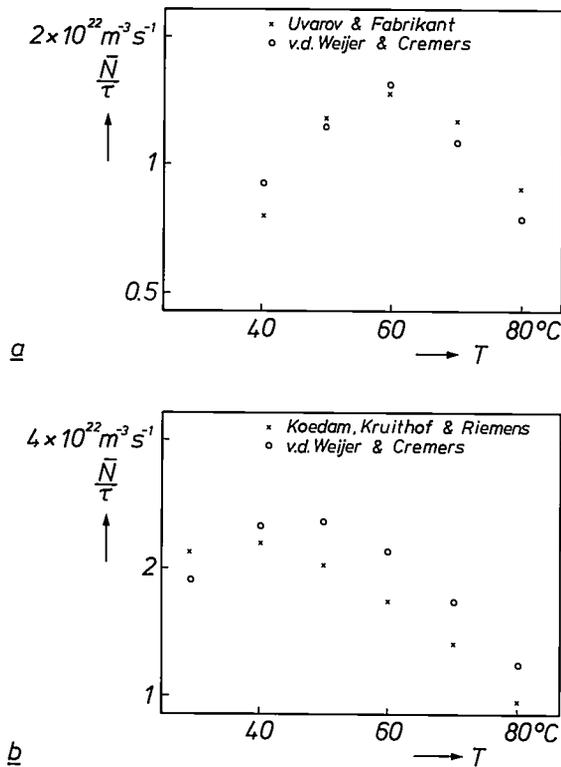


Fig. 11. a) The measured radiation outputs [20] at a wavelength of 254 nm as a function of the wall temperature T of the tube of a mercury discharge (diameter 38 mm, discharge current 450 mA), and the values that we calculated from the density N and the lifetime τ . b) The measured radiation output [21] at a wavelength of 254 nm of a mercury/argon discharge (diameter 36 mm, argon pressure 400 Pa, discharge current 400 mA) and our calculated values.

mental error of 15% — there were three different experiments — the results are the same.

The other discharge for which a comparison can be made between measured and calculated radiation output values is a mercury/argon discharge [21]. The agreement between measured and calculated results is not so good as for the pure mercury discharge. For wall temperatures higher than 50 °C the calculated values are 25% above the measured values. We believe that the agreement between measured and calculated values is sufficient to permit the conclusion that the quantities we measured are relevant parameters for the radiation output.

The radiation output at a wavelength of 185 nm was also measured for the mercury/argon discharges mentioned above. However, the density at the 6^1P_1 level in these discharges is so small that the hook method is not sensitive enough for a density measurement with an acceptable error. In discharge conditions where the 6^1P_1 density is high enough to permit measurements by the hook method (smaller diameter and hence greater current density, see fig. 10) the 185-nm radiation output was not measured. Because of this it is not possible at present to compare the 185-nm radiation output calculated from the density and lifetime of the 6^1P_1 level with a measured output as was done for the 254-nm line.

The main result of our measurements of the lifetime and density of atoms at the 6^1P_1 level is the improved understanding of the radiation-trapping mechanisms that play a part in this transition.

[19] F. A. Uvarov and V. A. Fabrikant, The absolute concentrations of excited atoms in the positive column of a mercury discharge, *Opt. & Spectrosc.* **18**, 433-437, 1965.

[20] F. A. Uvarov and V. A. Fabrikant, Experimental determination of the effective probability of photon emission by plasma atoms, *Opt. & Spectrosc.* **18**, 323-327, 1965.

[21] M. Koedam, A. A. Kruithof and J. Riemens, Energy balance of the low-pressure mercury-argon positive column, *Physica* **29**, 565-584, 1963.

Summary. A laser is used for measuring the lifetime and density of mercury atoms at the 6^3P_1 and 6^1P_1 levels in a low-pressure mercury/inert-gas discharge. The lifetime measurements are made by irradiating mercury atoms with laser light whose wavelength is such that the mercury atoms decay to the 6^3P_1 and the 6^1P_1 energy levels from a higher excited state. The lifetimes can be determined from the fluorescence signal of the subsequent transition to the ground state. The measured lifetimes agree reasonably well with values reported in the literature and have provided a better understanding of the excitation and decay mechanisms in a low-pressure mercury/inert-gas discharge. The density of the atoms at these levels was determined by the 'hook' method, in which the wavelength-dependent refractive index of the discharge in the neighbourhood of an absorption line depends on the number of atoms at the lower level of the transition. The densities at these levels and their lifetimes were used for determining the radiation output of low-pressure mercury/inert-gas discharges. Comparison of these calculated radiation outputs with outputs measured directly show fairly good agreement.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	A
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	B
Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	E
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	R

J.-M. Nicolas & J.-L. Bernatets	L	Atténuation correction using Mellin transform	Acoustical imaging Vol. 14, A. J. Berk- hout <i>et al.</i> (eds), Plenum, New York	459-469	1985
J.-M. Nicolas	L	A new method for the analysis of the frequency dependence of the backscattering coefficient in tissue-like media	Acoustical imaging Vol. 14, A. J. Berk- hout <i>et al.</i> (eds), Plenum, New York	761-764	1985
B. Post*, P. P. Gong* (* <i>Polytech. Inst. New York</i>), L. Kern & J. Ladell	N	The role of crystal rotation axis in experimental three- and four-beam phase determination	Acta Crystallogr. A42	178-184	1986
J.-P. Arragon, O. Chantelou, F. Fonsalás, M. le Quéau & G. Marie	L	Transmission d'images de télévision à haute définition compatible avec le standard D2-MAC/paquet	Acta Electron. 27	19-32	1985
J. Geninet (<i>Portenseigne, Fontenay-sous-Bois</i>)	L	Distribution par câble de programmes de télévision diffusés par satellite	Acta Electron. 27	33-36	1985
C. Jankowiak, M. Lamnabhi & J.-P. Arragon	L	Transmission sur réseaux câblés des signaux numériques de la famille X-MAC/paquet	Acta Electron. 27	37-54	1985
C. Kermarrec, J. Faguet, A. Collet & C. Mayousse	L	Récepteur d'émissions de télévision diffusées par satellite à 12 GHz intégré sur une seule puce de GaAs	Acta Electron. 27	65-76	1985
E. Rammos	L	Une antenne plane à lignes sur substrat suspendu pour les applications de réception satellite à 12 GHz	Acta Electron. 27	77-83	1985
J.-P. Michel	L	Convertisseurs vidéo 8 bits numérique/analogique et analogique/numérique	Acta Electron. 27	85-100	1985
T. Nillesen (<i>Philips Elcoma Div., Eindhoven</i>)	L	Line-locked digital colour decoding	Acta Electron. 27	101-107	1985
M. Lamnabhi	L	Filtrage numérique du son en radiodiffusion et télédiffusion directes par satellite	Acta Electron. 27	109-117	1985
C. B. Dekker, A. J. E. M. Janssen & P. J. van Otterloo	E	The contour plot method for noise reduction in digital video	Acta Electron. 27	119-131	1985
L. Bourassin (<i>La Radiotechnique, Suresnes</i>)	L	Introduction à la vidéographie	Acta Electron. 27	133-137	1985
A. Guénot (<i>RTC-Compélec, Paris</i>)	L	CIDAC: circuit intégré d'acquisition et de contrôle	Acta Electron. 27	139-144	1985
J. Haisma	E	Optical thin films produced by nonvacuum techniques	Appl. Opt. 24	2666-2673	1985
R. J. M. Zwiers & G. C. M. Dortant	E	Aspherical lenses produced by a fast high-precision replication process using UV-curable coatings	Appl. Opt. 24	4483-4488	1985

- | | | | | | |
|--|------|--|--|-----------|------|
| G. M. Loiacono, M. F. Shone, G. Mizell, R. C. Powell*, G. J. Quarles* & B. Elouadi* (* Oklahoma State Univ., Stillwater, OK) | N | Tunable single pass gain in titanium-activated lithium germanium oxide | Appl. Phys. Lett. 48 | 622-623 | 1986 |
| J. Maluenda, G. Martin, H. Schink* & G. Packeiser* (* Siemens, München) | L | Homogeneity qualification of GaAs substrates for large scale integration applications | Appl. Phys. Lett. 48 | 715-717 | 1986 |
| J. Mink (Philips Elcoma Div., Eindhoven) & B. H. Verbeek | E | Asymmetric noise and output power in semiconductor lasers with optical feedback near threshold | Appl. Phys. Lett. 48 | 745-747 | 1986 |
| J. Aarts, W. M. Gerits & P. K. Larsen | E | Observations on intensity oscillations in reflection high-energy electron diffraction during epitaxial growth of Si(001) and Ge(001) | Appl. Phys. Lett. 48 | 931-933 | 1986 |
| M. Taillepié & S. Gourrier | L | Current drift mechanism in In _{0.53} Ga _{0.47} As depletion mode metal-insulator field-effect transistors | Appl. Phys. Lett. 48 | 978-980 | 1986 |
| P. Blood, E. D. Fletcher, P. J. Hulyer & P. M. Smowton | R | Emission wavelength of AlGaAs-GaAs multiple quantum well lasers | Appl. Phys. Lett. 48 | 1111-1113 | 1986 |
| B. A. Joyce, J. H. Neave, P. J. Dobson & P. K. Larsen | R, E | Quantum wells and superlattices | Chem. Br. 21 | 1080-1084 | 1985 |
| M. Levent-Villegas & M. Soulard (TRT, Brive-la-Gaillarde) | L | Evaluation and influence of active and passive element variations on GaAs paraphase amplifier design | Coll. on ICs above 1 GHz, fabrication and circuit design, London 1986 | 4 pp. | 1986 |
| M. Binet | L | Test jig for digital GaAs IC mounted in 24 pins leadless chip carriers | Coll. on Microwave packaging, London 1986 | 3 pp. | 1986 |
| H. Sari & L. Desperben | L | Récupération de rythme adaptative pour égaliseurs numériques | Coll. sur le Traitement du signal et ses applications, Nice 1985 | 557-561 | 1985 |
| M. Lamnabhi & M. Louafi | L | Boucles à verrouillage de phase et synchronisation pour les transmission par salves | Coll. sur le Traitement du signal et ses applications, Nice 1985 | 589-594 | 1985 |
| S. Martin & M. Duseaux | | Kinetics of LEC GaAs properties under annealing | Defect recognition and image processing in III-V compounds, J. P. Fillard (ed.), Elsevier Science, Amsterdam | 287-294 | 1985 |
| H. Sari, S. Moridi, L. Desperben & P. Vandamme (CNET LAB, Lannion) | L | Interaction between an adaptive equalizer and a decision-feedback carrier recovery loop | Digital communications, E. Biglieri & G. Prati (eds), Elsevier Science, Amsterdam | 115-130 | 1986 |
| W. S. Newman | N | Adaptive control of a total artificial heart | Dynamic Systems: modelling and control, Vol. 1, M. Donath (ed.), Am. Soc. Mech. Eng., New York | 65-71 | 1986 |
| T. J. Mousley & K. P. McCarthy | R | Model of HF interference for application to data transmission | Electron. Lett. 22 | 70-71 | 1986 |
| F. H. P. M. Habraken*, W. F. van der Weg* (* Univ. Utrecht), J. Remmerie*, H. E. Maes* (* IMEC, Heverlee) & A. E. T. Kuiper | E | Physico-chemical and electrical characterization of LPCVD silicon oxynitride layers | ESPRIT'85, Elsevier Science, Amsterdam | 191-201 | 1986 |
| P. C. M. Gubbens*, A. M. van der Kraan* (* Interuniv. Reactor Inst., Delft) & K. H. J. Buschow | E | Quadrupole pair order in cubic TmGa ₃ studied by ¹⁶⁹ Tm Mössbauer spectroscopy | Hyperfine Interactions 29 | 1343-1346 | 1986 |
| T. A. C. M. Claasen & W. F. G. Mecklenbräuker (Vienna Inst. Technol.) | E | Adaptive techniques for signal processing in communications | IEEE Commun. Mag. 23 (No. 11) | 8-19 | 1985 |
| J. M. Shannon & B. J. Goldsmith | R | Bulk unipolar camel diodes formed using indium implantation into silicon | IEEE Electron Dev. Lett. EDL-6 | 583-585 | 1985 |

- A. J. E. M. Janssen & T. A. C. M. Claasen *E* On positivity of time-frequency distributions IEEE Trans. **ASSP-33** 1029-1032 1985
- R. Bradley (*N.A.P. Consumer Electron., Knoxville, TN*), J. F. Goldenberg & T. S. McKechnie *N* Ultra-wide viewing angle rear projection television screen IEEE Trans. **CE-31** 185-193 1985
- P. J. Patt *N* Design and testing of a coaxial linear magnetic spring with integral linear motor IEEE Trans. **MAG-21** 1759-1761 1985
- J. E. Knowles Packing factor and coercivity in tapes; a Monte Carlo treatment IEEE Trans. **MAG-21** 2576-2582 1985
Correction IEEE Trans. **MAG-22** 68 1986
- H. K. Kuiken *E* Edge effects in crystal growth under intermediate diffusive-kinetic control IMA J. Appl. Math. **35** 117-129 1985
- M. Taillepiéd, S. Gourier & J. P. Chané *L* Electrical characterization of the insulator/ $\text{Ga}_{0.47}\text{In}_{0.53}\text{As}$ interface Insulating films on semiconductors, J. J. Simonne & J. Buxo (eds), Elsevier Science, Amsterdam 85-88 1986
- M. J. Verkerk, A. van Eenbergen & E. Bruninx *E* ESCA study of the adhesion failure of galvanic nickel layers Int. J. Adhesion & Adhesives **6** 42-44 1986
- A. J. Lowe*, M. J. Powell & S. R. Elliot* (** Univ. Cambridge*) *R* The electronic properties of plasma-deposited films of hydrogenated amorphous SiN_x ($0 < x < 1.2$) J. Appl. Phys. **59** 1251-1258 1986
- C. J. van der Poel, D. J. Gravesteijn, W. G. V. M. Rippens, H. T. L. P. Stockx & C. M. J. van Uijen *E* Phase-change optical recording in TeSeSb alloys J. Appl. Phys. **59** 1819-1821 1986
- A. G. Dirks, T. Tien & J. M. Towner *E, S* Al-Ti and Al-Ti-Si thin alloy films J. Appl. Phys. **59** 2010-2014 1986
- H. J. Cornelissen & H. J. H. Merks-Eppingbroek *E* Electron density in low-pressure Na-Ne discharges deduced from laser absorption experiments J. Appl. Phys. **59** 2324-2331 1986
- G. B. McMillan (*Intel Corp., Santa Clara, CA*), J. M. Shannon, J. B. Clegg & H. Ahmed (*Univ. Cambridge*) *R* Characterization of shallow ($R_p < 20$ nm) As- and B-implanted and electron-beam annealed silicon J. Appl. Phys. **59** 2694-2703 1986
- W. M. Wagenaars*, A. J. M. Houtsmma* & R. A. J. M. van Lieshout* (** Inst. Perception Res., Eindhoven*) Subjective evaluation of dynamic compression in music J. Audio Eng. Soc. **34** 10-18 1986
- P. J. Schoenmakers & F. C. C. J. G. Verhoeven *E* Effect of pressure on retention in supercritical-fluid chromatography with packed columns J. Chromatogr. **352** 315-328 1986
- H. M. van Noort, B. C. M. Meenderink & A. Molenaar *E* *In situ* ^{119}Sn conversion electron Mössbauer study of the surface of tin layers as deposited by an electroless process J. Electrochem. Soc. **133** 263-265 1986
- K. H. J. Buschow, D. B. de Mooij, J. L. C. Daams & H. M. van Noort *N* Phase relationships, magnetic and crystallographic properties of Nd-Fe-B alloys J. Less-Common Met. **115** 357-366 1986
- T. T. M. Palstra*, G. J. Nieuwenhuys*, J. A. Mydosh* (** Univ. Leiden*), R. B. Helmholtz (*ECN, Petten*) & K. H. J. Buschow *E* Neutron diffraction and magnetostriction of cubic $\text{La}(\text{Fe}_x\text{Al}_{1-x})_{13}$ intermetallic compounds J. Magn. & Magn. Mater. **54-57** 995-996 1986
- W. Steiner*, M. Reissner* (** Tech. Univ. Wien*) & K. H. J. Buschow *E* Magnetic and ^{57}Fe Mössbauer investigations on $\text{Y}(\text{Fe},\text{Ir})_2$ J. Magn. & Magn. Mater. **54-57** 1079-1080 1986
- R. Grössinger*, X. K. Sun*, R. Eibler*, K. H. J. Buschow & H. R. Kirchmayr* (** Tech. Univ. Wien*) *E* Temperature dependence of anisotropy fields and initial susceptibilities in $\text{R}_2\text{Fe}_{14}\text{B}$ compounds J. Magn. & Magn. Mater. **58** 55-60 1986
- G. de With & H. Verweij *E* Properties and shaping of lightweight ceramics based on phosphate-bonded hollow silica microspheres J. Physique **47** (Colloque C1) C1/359-C1/363 1986
- G. de With *E* Note on the relation between the compressive strength of debased alumina and its use as hot-pressing die material J. Physique **47** (Colloque C1) C1/667-C1/671 1986
- R. Ward, A. R. Franklin, I. H. Lewin, P. A. Gould & M. J. Plummer *R* A 1:1 electron stepper J. Vac. Sci. & Technol. **B 4** 89-93 1986

- | | | | | | |
|---|---|---|--|-----------|------|
| A. M. E. Hoeberechts & G. G. P. van Gorkom | E | Design, technology, and behavior of a silicon avalanche cathode | J. Vac. Sci. & Technol. B 4 | 105-107 | 1986 |
| G. G. P. van Gorkom & A. M. E. Hoeberechts | E | Performance of silicon cold cathodes | J. Vac. Sci. & Technol. B 4 | 108-111 | 1986 |
| C. T. Foxon, J. J. Harris, R. G. Wheeler (<i>Yale Univ., New Haven, CT</i>) & D. E. Lacklison | R | Observation of strong localization effects in (AlGa)As-GaAs twodimensional electron gas structures at low magnetic fields | J. Vac. Sci. & Technol. B 4 | 511-514 | 1986 |
| J. M. Woodcock, J. J. Harris & J. M. Shannon | R | Controlled growth of GaAs layers for monolithic hot electron transistors | J. Vac. Sci. & Technol. B 4 | 609-611 | 1986 |
| L. Kuijpers, M. Janssen, J. de Witt | E | Perspektiven für Kapazitätsregelung bei Kühl- und Gefriergeräten | Ki Klima-Kälte-Heizung 14 | 87-91 | 1986 |
| D. J. Oostra*, A. Haring*, A. E. de Vries* (<i>* FOM, Amsterdam</i>), F. H. M. Sanders & G. N. A. van Veen | E | Etching of silicon by SF ₆ induced by ion bombardment | Nucl. Instrum. & Methods Phys. Res. B 13 | 556-560 | 1986 |
| M. F. C. Willemsen, A. M. L. Theunissen & A. E. T. Kuiper | E | Hydrogen profiling by proton-proton scattering | Nucl. Instrum. & Methods Phys. Res. B 15 | 492-494 | 1986 |
| A. A. van Gorkum & T. G. Spanjer | E | A generalized comparison of spherical aberration of magnetic and electrostatic electron lenses | Optik 72 | 134-136 | 1986 |
| M. E. A. Corthout & H. B. M. Jonkers | E | The transformational development of a new point containment algorithm | Philips J. Res. 41 | 83-174 | 1986 |
| P. Piret | B | Concatenated codes for mobile automatic telephony | Philips J. Res. 41 | 175-190 | 1986 |
| W. P. Lord | N | Interactive graphics based software tools for rapid prototyping | Philips J. Res. 41 | 191-204 | 1986 |
| R. N. Bates & A. G. Stove | R | Millimetre-wave radar | Philips J. Res. 41 | 206-218 | 1986 |
| R. A. Brown, R. J. Dewey & C. J. Collier | R | Some features of signal demodulation resulting from the practical implementation of a direct conversion radio receiver | Philips J. Res. 41 | 219-231 | 1986 |
| H. A. Post | E | Radiative transport at the 184.9-nm Hg resonance line. I. Experiment and theory | Phys. Rev. A 33 | 2003-2016 | 1986 |
| H. A. Post, P. van de Weijer & R. M. M. Cremers | E | Radiative transport at the 184.9-nm Hg resonance line. II. Extensive experiments | Phys. Rev. A 33 | 2017-2024 | 1986 |
| S. Barth*, E. Albert*, G. Heiduk*, A. Möslang*, A. Weidinger*, E. Recknagel* (<i>* Univ. Konstanz</i>) & K. H. J. Buschow | E | Local magnetic fields in ferromagnetic intermetallic compounds of cubic Laves-phase type | Phys. Rev. B 33 | 430-436 | 1986 |
| D. M. Krol, K. B. Lyons*, S. A. Braver* & C. R. Kurkjian* (<i>* AT&T Bell Labs, Murray Hill, NJ</i>) | E | High-temperature light scattering and the glass transition in vitreous silica | Phys. Rev. B 33 | 4196-4202 | 1986 |
| R. Moreh (<i>Brookhaven Nat. Lab., Upton, NY</i>), O. Shahal (<i>Nucl. Res. Center, Beer-Sheva</i>) & G. Kimmel | N | Orientation of nitrate molecules in graphite-HNO ₃ residue compounds | Phys. Rev. B 33 | 5717-5720 | 1986 |
| J. G. Kloosterboer & G. J. M. Lippits | E | Network formation by photopolymerization: high-precision replication of video discs | Polym. Preprints 26 | 351-352 | 1985 |
| W. Schreiner | N | Towards improved alignment of powder diffractometers | Powder Diffraction 1 | 25-33 | 1986 |
| A. J. M. Kaizer | E | The modelling of the nonlinear response of an electrodynamic loudspeaker by a Volterra series expansion | Proc. AES Conv., Montreux 1986 (preprint) | 22 pp. | 1986 |
| W. Newman | N | Nonlinear normalization of magnetic structure with applications to computer aided design of stepper motors | Proc. Ann. Symp. on Incremental motion control systems & devices, Urbana-Champaign, IL, 1985 | 8 pp. | 1985 |
| W. Newman | N | Design optimization of a high scanner motor for an infrared imaging device | Proc. Ann. Symp. on Incremental motion control systems & devices, Urbana-Champaign, IL, 1986 | 15 pp. | 1986 |
| C. Zardini*, J. D. Pistre*, F. Rodes*, J. L. Aucouturier* (<i>* Univ. Bordeaux I, Talence</i>), M. Monneray & H. Baudry | L | Ultrasonic bonding of heavy aluminium wires on copper thick films | Proc. Eur. Hybrid Microelectron. Conf., Stresa 1985 | 231-237 | 1985 |

- | | | | | | |
|---|-------------|---|---|-----------|------|
| H. Baudry & G. Kersuzan (<i>TRT, Le Plessis Robinson</i>) | <i>L</i> | Thick film hybrids on copper-clad invar: an evaluation | Proc. Eur. Hybrid Microelectron. Conf., Stresa 1985 | 454-464 | 1985 |
| A. Fihel & H. Sari | <i>L</i> | Reduced-bandwidth 16 QAM: a power and bandwidth efficient modulation technique for digital microwave radio systems | Proc. ICC'85, Chicago 1985 | 466-470 | 1985 |
| E. H. L. Aarts & P. J. M. van Laarhoven | <i>E</i> | A new polynomial-time cooling schedule | Proc. ICCAD 85, Santa Clara, CA, 1985 | 206-208 | 1985 |
| M. Amato & V. Rumennik | <i>N</i> | Comparison of lateral and vertical DMOS specific on-resistance | Proc. IEDM 85, Washington, DC, 1985 | 736-739 | 1985 |
| G. Bruning | <i>N</i> | A comparative introduction of a new high voltage resonant oscillator | Proc. IEEE Applied Power Electronics Conf. & Exposition, New Orleans 1986 | 76-82 | 1986 |
| R. P. Tijburg, G. L. Dinghs & T. van Dongen | <i>E</i> | Criteria for the selection and treatment of contact metals to p-type III-V compounds | Proc. Int. Symp. GaAs and related compounds, Karuizawa 1985 | 355-360 | 1986 |
| F. P. M. Beenker | <i>E</i> | Systematic and structured methods for digital board testing | Proc. Int. Test. Conf. Philadelphia 1985 | 380-385 | 1985 |
| J. A. Clarke | <i>R</i> | Wide angle lenses for projection TV | Proc. SPIE 554 | 394-397 | 1986 |
| R. H. Coursant, J. M. Tellier & M. Naillon | <i>L</i> | Characterization of piezoelectric ceramics used in multielement ultrasonic transducers for medical imaging | Proc. Ultrasonics Int. '85, London 1985 | 563-568 | 1985 |
| M. Naillon & F. Besnier (<i>CISI-GRAPH, Rungis</i>) | <i>L</i> | Finite element analysis of piezoelectric structure electro-mechanical response | Proc. Ultrasonics Int. '85, London 1985 | 869-877 | 1985 |
| P. Pesqué & C. Méquio | <i>L</i> | A new and fast calibration method for ultrasonic hydrophones | Proc. Ultrasonics Symp., Dallas 1984 | 743-747 | 1984 |
| P. Pesqué & M. Fink | <i>L</i> | Effect of the planar baffle impedance in acoustic radiation of a phased array element theory and experimentation | Proc. Ultrasonics Symp., Dallas 1984 | 1034-1038 | 1984 |
| W. A. Smith, A. Shaulov & B. A. Auld (<i>Stanford Univ.</i>) | <i>N</i> | Tailoring the properties of composite piezoelectric materials for medical ultrasonic transducers | Proc. Ultrasonics Symp., San Francisco 1985 | 642-647 | 1985 |
| A. Shaulov & W. A. Smith | <i>N</i> | Ultrasonic transducer arrays made from composite piezoelectric materials | Proc. Ultrasonics Symp., San Francisco 1985 | 648-651 | 1985 |
| N. M. Marinovic & W. A. Smith | <i>N</i> | Suppression of noise to extract the intrinsic frequency variation from an ultrasonic echo | Proc. Ultrasonics Symp., San Francisco 1985 | 841-846 | 1985 |
| P. Bachmann | <i>A</i> | Review of plasma deposition applications: preparation of optical waveguides | Pure & Appl. Chem. 57 | 1299-1310 | 1985 |
| B. A. Joyce | <i>R</i> | Molecular beam epitaxy | Rep. Prog. Phys. 48 | 1637-1697 | 1985 |
| P. Blood | <i>R</i> | Capacitance-voltage profiling and the characterisation of III-V semiconductors using electrolyte barriers | Semicond. Sci. & Technol. 1 | 7-27 | 1986 |
| T. Carette*, M. Lannoo*, G. Allan* (* <i>ISEN, Lille</i>) & P. Friedel | <i>L</i> | Theoretical and experimental study of the hydrogenated (100) MBE grown surface of GaAs | Surf. Sci. 164 | 260-270 | 1985 |
| B. A. Joyce, P. J. Dobson, J. H. Neave, K. Woodbridge, J. Zhang (<i>Imp. College, London</i>), P. K. Larsen & B. Bölger | <i>R, E</i> | RHEED studies of heterojunction and quantum well formation during MBE growth — from multiple scattering to band offsets | Surf. Sci. 168 | 423-438 | 1986 |
| P. K. Larsen, P. J. Dobson, J. H. Neave, B. A. Joyce, B. Bölger & J. Zhang (<i>Imp. College, London</i>) | <i>E, R</i> | Dynamic effects in RHEED from MBE grown GaAs (001) surfaces | Surf. Sci. 169 | 176-196 | 1986 |
| P. Baudet, P. Rabinzohn, C. Venet, P. Gamand & C. Varin | <i>L</i> | A monolithic GaAs 2-18 GHz travelling wave amplifier | Workshop on Compound semiconductor integrated circuits, Visby 1986 | 4 pp. | 1986 |

Bright spots and bottlenecks in Europe's future industrial development

W. Dekker

It has long been our editorial policy to include occasional articles of a more general nature as well as the articles directly concerned with research and the results it provides. We are trying in this way to bring the development of science into a broader perspective and to paint in the backcloth against which our research is enacted.

In recent years economic developments have imposed much stricter boundary conditions for our research activities than ever before. One result of these economic circumstances has been that the management of our company has become more 'visible'. The problems that our company has to contend with have been brought to the notice of various kinds of public forum, especially since Dr W. Dekker became President.

We hope the publication of the speech 'Bright spots and bottlenecks in Europe's future industrial development' will make the reader more familiar with the outlook of the management on certain socio-economic and political problems. It gives, we think, an impression of the way in which the progress of technology and science are affected by external factors as well as internal ones.

The speech printed below was given by Dr Dekker at the conference 'A competitive future for Europe', organized by the Erasmus University of Rotterdam on the 12th and 13th December 1985, and has been made available to us by the Corporate External Relations Department. The text of the speech is given verbatim; the editors have only added the figures.

Anyone who has regularly taken part in conferences such as this over the last few years will have noticed that 'the future' is a subject which crops up again and again. For some of you 'the future' will stand for uncertainty and threats, while for others it will above all represent a challenge and new promises.

In my view, the future encompasses both these aspects. Its possibilities include positive as well as less desirable developments. The direction in which the future develops partly depends on our own efforts, our will and our understanding of events. The object

of a conference such as this is to establish what factors are important, what options we have and what strategies should be pursued.

'Forecasting is a difficult business, especially where the future is concerned', people can often be heard to sigh. And I heaved the same sigh when preparing this introduction on 'Bright spots and bottlenecks in Europe's future industrial development'.

It is not difficult to obtain a clear picture of which industries in Europe are performing well or badly at this moment. When it comes to finding the explanations for this, however, we are groping more in the dark. And in order to know which companies will have disappeared by the year 2000 and which ones will be flourishing, one would almost have to be a clairvoyant.

Dr W. Dekker, formerly President of N.V. Philips' Gloeilampenfabrieken, is now Chairman of the Supervisory Board. This speech is published with the permission of Philips International B.V., Corporate External Relations.

In seeking to answer questions like this one feels like Banquo in Shakespeare's *Macbeth*, who exclaims:
 'If you can look into the seeds of time
 And say which grain will grow and which will not,
 Speak then to me!'

Did Banquo perhaps himself foresee that he would not survive until the end of the drama?

I do not regard it as my task to offer a judgement here on which companies will and which will not survive beyond the year 2000. What I would like to do, however, is to consider what demands will be made on companies in the future and what changes will take place in the environment in which they operate.

Let me begin by briefly outlining the future business environment. One of the main conclusions here will be that we are living in an age of rapid change in which technological innovations are a crucial factor. Following on from this, I shall then look more closely at the significance of the new technology, or what is generally known as 'high tech'. On the basis of an example of a high-tech business within my own group of companies I shall seek to list a number of conditions for success. There are conditions which in my view can be generalized, in other words they are conditions which apply, *mutatis mutandis*, to many high-tech industries.

I have already stated that the future will be determined to a significant extent by our own action, our own objectives and understanding. Several actors are involved. Industrialists, national governments and the European organizations will all be called upon to make their own contribution and all of them have their own responsibility. I shall deal both with the tasks which face companies and the conditions which a European industrial policy should fulfil.

In conclusion I shall make a plea for what I call 'multiple reciprocity', an open form of cooperation both between the countries of Europe and between companies. In my view this type of cooperation is a precondition for further European integration and for the desired cooperation with regard to the new technologies.

It seems that people in Europe — the ordinary citizen, but also many politicians and parliamentarians and many businessmen — often scarcely realize that important shifts are taking place. It is true that books like Toffler's 'The third wave' and Naisbitt's 'Megatrends' are printed in large editions, but the ideas put forward scarcely seem to affect the reader's conception of the future^[1].

Thus, what is referred to in a recent SER report as a 'Eurocentric system of thought'^[2] often seems to be in evidence. In other words, there is insufficient awareness that Europe is no longer number one. In

the USA's future perspective, for instance, Japan is probably more important than Europe.

This shift in emphasis is already well under way. To take an example, trade between the USA and Western Europe in 1979 was still worth over one billion dollars more than trade between the USA and Japan plus South-East Asia. In 1983, however, trade between the USA and that part of Asia exceeded trade between the USA and Europe by almost 29 billion dollars. If we are not able to reverse this trend, a Pacific alliance will emerge in which Japan and the United States call the tune and Europe has to follow.

Some important trends which I see for the future of industrial enterprise in Europe are:

1. A transition towards a multipolar world economy in which, along with a further integrated Europe and the USA, Japan and the newly industrialized countries will constitute major centres of world production and world trade.
2. A globalization of industrial production, as companies increasingly produce for a world market (*fig. 1*), with modifications being made locally to the software and the design.
3. A transition to an information society in which the greater part of the working population will be involved in generating and processing information, software perhaps becoming more important than hardware.
4. An increase in the pace of technological innovation and the emergence of new technologies and in-



Fig. 1. Trento, Italy. Refrigerators being produced on an international scale.

dustries (in electronics, biotechnology, new materials) and of as yet unknown applications.

5. Significant socio-cultural and socio-political shifts resulting in new forms of organization (with the emphasis on autonomous groups, business units, holdings) and new political conditions with deregulation and a rethinking of the role of the state.

In describing the first trend as the emergence of a multipolar world system I placed Europe alongside the United States of America and Japan. Many people will regard this as wishful thinking, and they are right. I certainly wish Europe a place in the sun alongside the USA and Japan. But I also believe that this is possible. I am not one of those pessimists who now describe the situation of European industry in the same way as historians characterized the situation of the Netherlands in the year of calamity 1672: 'The rulers are helpless, the people senseless and the nation hopeless'.

If European industry shows common sense and goodwill its situation will by no means be hopeless. Allow me to support my optimism by reference to a number of indicators.

As far as technological innovation is concerned, Europe is certainly not lagging behind. According to the UN Statistical Year Book almost 96 000 patents were granted in 1980 to the ten countries of the European Community, compared with almost 62 000 to the United States and 46 000 to Japan.

Expenditure on research and development in Europe is roughly 20% less than in the USA, but greater than the amount spent by Japan.

In a number of sectors such as aviation, rail transport technology, nuclear energy, new materials and the chemicals industry, Europe still leads the way. And Europe's position in important growth areas like biotechnology and the space industry is far from weak.

According to figures published by the OECD, the increase in high-tech industry's share of total industrial production in the European Community does not lag behind that of the other major economic blocs, Japan and the USA. The share of the most advanced high-tech sectors (i.e. industries in which R&D expenditure was higher than 4% of the value of production) in total industrial production in the European Community rose between 1970 and 1980 from 11.4% to 11.7%. In the same period there was a fall from 14.6% to 10.8% in the United States and a fall from 14.1% to 13.4% in Japan.

Nor is the European Community lagging behind in average high-tech industry, where R&D expenditure is between 1 and 4% of the total value of production. On the contrary, with average high-tech industry

having a share of 32.2% of total industrial production (in 1980), the European Community was in fact ahead of Japan (30.1%) and the United States (31.6%).

However, Europe's relative share of trade in high-tech products is lower than that of the USA and Japan. On the one hand, this may be due to prices being too high or quality being too low (which might be the result of the fact that these industries emerged in Europe under protective conditions). On the other hand, trade barriers (often of a non-tariff nature, particularly in Japan) are also a factor. Over the last 20 years the export specialization coefficient — an indicator expressing the share of technically advanced products in total exports or total industrial exports — has shown that Europe's position has deteriorated in relation to the USA and Japan. Although this indicator is open to considerable criticism, because it does not take account, for instance, of such an important sector for Europe as chemicals and because intra-European trade has the effect of lowering this coefficient, the trend, namely the gap between Europe on the one hand and the USA and Japan on the other, should give us cause for concern.

Looking at all the indicators which I have mentioned, I come to the conclusion that European industry still has good prospects. There are, however, certain areas where Europe seems to be losing ground, and it is here that a specific policy on the part of government and joint efforts by European industry are necessary to prevent a further deterioration of Europe's position.

A European innovation policy will have to be supported both by individual governments and by the supranational organizations.

High-tech industry assumes an important place in the discussion about the kind of industrialization policy to be pursued and the need for industrial innovation. Although I in no way underestimate the importance of high tech — and you would not have expected anything else from someone representing the electronics industry, a high-tech industry *par excellence* — I would nevertheless like to put its importance in perspective.

• First of all, innovation not only has to be stimulated in the high-tech industries but should take place over a much broader field. New production technologies should also be introduced in the more traditional sec-

[1] A. Toffler, *The third wave*, PAN books, London 1981; J. Naisbitt, *Megatrends*, Macdonald, London 1984.

[2] The term 'Eurocentric system of thought' and the figures concerning the Pacific Basin are taken from the report 'Nederland in de wereld-economie — perspectieven en mogelijkheden' ('The Netherlands in the world economy — perspectives and possibilities') by the Dutch Social and Economic Council (SER) committee for development problems of companies, 1985.

tors such as agriculture, textiles or the building industry. If productivity growth lags behind in these sectors, this can indirectly have an adverse effect on the growth potential of the high-tech industries. Innovation in the more traditional sectors would also mean that a sales market for new production technologies would be created.

under the high-tech heading, even though they can no longer easily be made without very advanced technological knowledge. The oil industry is a case in point. Its end-product is not regarded as high tech, even though high tech plays an essential role in the overall production chain. The term 'high tech' is only relative, therefore, and we should not be blinded by it.

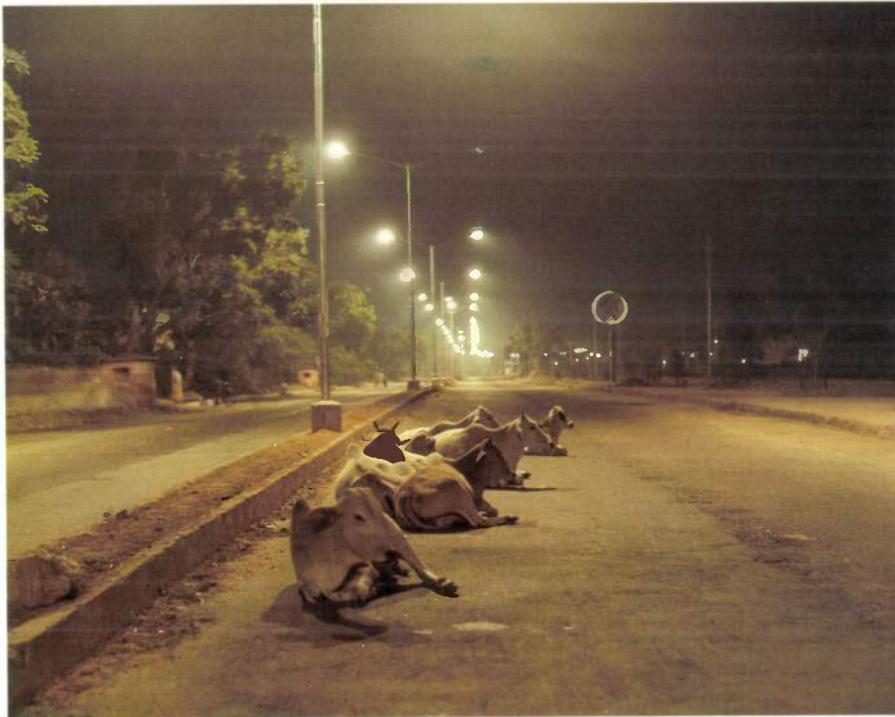


Fig. 2. Bokaro, India. Yet another picture of the contrast between the traditional East and the modern West. In the future this contrast will become less marked, since the manufacture of lamps and other high-value technical products will take place in the newly industrialized countries.

- Secondly, it is unlikely that high-tech companies will be able to provide sufficient employment to absorb the present surplus of labour. The cost of unemployment can exert severe pressure on the potential of the new industries for development. An all-round economic revival should therefore be the aim. At the same time we should beware that general measures aimed at reducing unemployment do not have an inhibiting effect on growth in the high-tech industries. These industries are often faced, for example, with a scarcity of labour. It would not make sense to conclude the same types of agreement on the reduction of working hours with this category as with the abundant supply of workers in the other sectors.

- My third observation concerns the term 'high-tech' industry itself. According to the OECD definition, high tech refers to industrial products which require an above-average intellectual input. However, a number of products or industrial activities do not fall

- Nor should our interest in high tech make us forget the importance of the service sector. The service sector already has a significant added value and is instrumental, moreover, in promoting high-tech industry. One need only think of the provision of venture capital and trade credit facilities, export assistance, marketing research, effective and reliable transport systems, etc.

Conditions for success

The business climate of the future will be different from that of today. Different demands will be made on high-tech companies. We must assume, for instance, that the newly industrialized countries will not be content to take over only the traditional industries of the West, such as the footwear and textile industries or shipbuilding. Their ambitions and potential go beyond this. In the future we shall also be con-

fronted with tough competition from the newly industrialized countries in the areas of high technology (fig. 2). As a result, the life cycles of the various new products will become shorter than we have been accustomed to until now and profit margins will become smaller. If we add to this the fact that there will be rapid increases in the cost of research, development and production, we can only conclude that industrial enterprise will increasingly become a high-risk enterprise in which only healthy companies with sufficient reserves will be able to engage.

Let me illustrate this process with a practical example from the manufacture of integrated circuits. If we put the costs of a building for IC manufacture in 1970 at 1, a further investment in installations and machinery of 2.5 was required. The depreciation periods for the building and the installations and machinery at that time were 50 and 10 years respectively.

In 1980 the investment costs for the building and the installations and machinery amounted to 2.5 and 4.5 and the depreciation period was 50 and 7 years respectively. For the near future — we are talking about 1987 — we can assume a drastic increase in the required investment up to an index value of 30 for buildings and 72 for installations and machinery, while the depreciation period will be further reduced to 30 and 6 years respectively.

In 17 years' time we therefore see an increase in the required investments by a factor of about 30, while the depreciation period is reduced by roughly 40%. The production capacity in this particular case increases by a factor of 20. Although the total market for integrated circuits will grow substantially in the future, we shall only be able to obtain a reasonable share of this growth market — and hence recoup the capital invested — if we can meet the specific requirements of the customer, if the products are of exceptionally good quality and reasonably priced, if our delivery reliability and service are unimpeachable and if we put the product on the market at the right moment.

I shall repeat these conditions once again because I consider that they are a *sine qua non* for any industry that wishes to be successful.

The conditions for success, together with their implications, are:

- *Customized production* — This means having a wide diversity of products, something which it will only be possible to achieve with a modular construction of the product and a flexible production process (fig. 3). The technological trend is moving in this direction in many production and assembly sectors. It also means having good relations with the customer

and in some cases supplying tailor-made software. We can also note here that for some industrial companies the delivery of software will become a more important source of income than the delivery of the original hardware.

- *Good quality* — The customer is becoming increasingly quality-conscious. He demands quality and will compare the quality of different products. As a result,

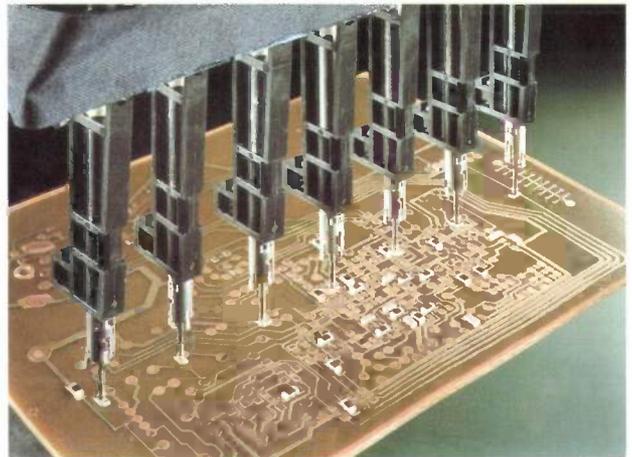


Fig. 3. The placement unit of an MCM machine (for Modular Chip Mounting). It directs about 100 000 electronic components an hour to positions specified on a printed-circuit board by computer (and therefore easily modified).

many technologically advanced consumer products are assuming a professional or semi-professional character.

- *Delivery reliability and service* — Operating on a world market means that one creates a worldwide distribution network and where necessary can offer direct service. This already applied in the past. For a number of products the trend is towards a lower service requirement. If, however, maintenance and repairs have to be carried out, this often requires very specific knowledge, and service will have to be provided by the manufacturer himself. Delivery reliability also means that one is able to make products which conform to the agreed performance characteristics — so we are again talking about quality requirements, which are becoming increasingly stringent.

- *Good timing* — In an innovative market with a lot of competition this means that the innovation process has to be shortened by the integration of research, development, production and marketing. This calls for integral design and a system-based approach to the whole of the business.

• *A good price* — In many cases it will only be possible to offer this if production is on a large scale and if there is a significant domestic market. A large domestic market often means more or less guaranteed sales, which in turn means that initial costs are covered and a much lower price can be offered abroad (a strategy much favoured by Japan). As far as the requirement of large-scale production is concerned, we have witnessed a change in the last few years, under the impact of new technological capabilities. The concept of 'economies of scale' is being replaced by that of 'economies of scope'. This implies that, while large volumes should be produced, there should also be a varied product mix, a family of products (*fig. 4*). In many branches of industry we can see that the learning curve is falling more sharply. This means that if one arrives on the market too late with one's product, the market price for this particular product has already fallen, as a result of which the high initial costs can no longer be recouped. — Closely related to the aspects of timing and price is the last condition for success, namely

• *The right place* — New products will have to be offered where the demand for them first arises. Only then will it be possible to offer them at the higher prices which apply to the beginning of the learning curve. This requirement with regard to location often means having a physical presence in regions where innovations first take place. This often means having a base in the United States (*fig. 5*) or Japan, where consumers and government think more in terms of technological innovations than is usually the case in Europe.

Industrial policy

I shall now turn to industrial policy — the objectives, the practice, and the conditions which from a European point of view it should fulfil.

Industrial policy can have different characteristics, such as:

- a. Creating or maintaining employment.
- b. Supporting or protecting new initiatives, or stimulating innovations.
- c. Safeguarding national independence and security.
- d. Restructuring a branch of industry; often bringing together a number of small companies.

These objectives apply at the national level but also at a European level. In practice, industrial policy as pursued by the various European countries has 'not always' — and this is an understatement — had the envisaged effect. Nationalistic and protectionist considerations have frequently led to what from a European point of view are the wrong measures.

Let me list some of the adverse consequences:

- Inefficient companies have been kept alive.
- Protectionism has reduced the incentive to innovate, both with regard to products and production processes.
- The same spearhead activities were often developed in different countries of the community, which inevitably resulted in duplication and over-capacity.
- The national market was often too small and the companies operating in it were often of less than optimum size.



Fig. 4. Family of stereo cassette recorders all based on two types of drive mechanism for the cassettes (foreground) and a circuit board (behind it). Depending on the type of recorder, all of the electronic components on the board are used, or only some of them.

- Nationally-oriented policies encouraged chauvinism and prevented integration between companies in different European countries.

A European industrial policy has to meet a number of conditions in order to prevent these problems. I will summarize them in five basic conditions:

1. An industrial revival must come from the companies themselves — national and community policy should support and stimulate initiatives from below; cooperation between companies across frontiers should be encouraged and measures taken to prevent overcapacity.
2. It should be innovative and future-oriented — with the emphasis on identifying and stimulating spearhead activities and new high-tech industries; a structural transformation of the market sector should be the aim, not the refurbishing and maintenance of outdated business.
3. The policy should be market-based — support measures should be of a temporary nature; companies

should ultimately be able to operate on a world market.

4. The policy should result in a single homogeneous European market. Key elements here are the harmonization of legislation and standardization.

5. The policy should be consistent, transparent and pragmatic — red tape should be cut; rules and regulations which have proved ineffective or have become outmoded should be withdrawn.

mation between the universities and the national and Community research institutes; prevent duplication.

e. Stimulate cooperation between these research institutes and industry; create a link between theoretical research and practical application; science transfer points for small companies should be created in the research institutes.

f. Prevent a brain drain (to the United States) by enhancing the status of research; start a number of



Fig. 5. In the United States the 'centre of gravity' of the electronics industry has for some time been moving away from the East Coast towards the area around San Francisco on the West Coast — including 'Silicon Valley', where Philips have established a base through the acquisition of Signetics.

Within the framework of these basic conditions and the problems referred to, a number of specific policy measures can be formulated which could give a particular boost to high-tech industry in Europe. These are measures to be taken both at the national and the Community level:

- a. Create a favourable tax climate for research and development.
- b. Stimulate cooperation between industries (not least between small companies) in the new technology and remove inhibiting regulations in this field.
- c. Stimulate/subsidize relevant research which could otherwise not be tackled because of the high risks involved.
- d. Stimulate cooperation and the exchange of infor-

inspiring, large-scale research projects; create confidence in Europe's capacity for innovation.

g. Raise the quality and quantity of university education; encourage more women to study in the areas of the new technology.

h. Create technology parks; abolish bureaucratic rules which stand in the way of the establishment of new companies; help small businesses to implement new initiatives.

i. Provide credit for financing and create venture capital.

j. Place orders for development work and pursue a procurement policy that stimulates innovation.

These 10 points for action should be tackled energetically if we in Europe do not want to be left behind

by the United States and Japan. To quote Mr Ponia-towski of the European Commission [3], we should prevent Europe from becoming the first 'colony of the third industrial revolution'. If we do 'too little, too late', Europe will become 'the continent of lost opportunities'.

But however energetically this industrial policy is pursued, no industrial policy, either at a Community or a national level, will be successful in the long term if it is not embedded in or supported by the further implementation and extension of the ideals of the European Community.

I regard as a central task facing the European Community — the member states and both sides of industry in the member states:

1. The creation of a truly common market in which European goods can move freely from one country to another, without papers, documents or other non-tariff trade barriers and with a uniform system of VAT.
2. The achievement of a convergent industrial policy through, for instance, the liberalization of the present nationalistic procurement policy; through the standardization of infrastructure facilities; through the stimulation of a number of major European projects; through acting as a single trading power in international trade.
3. The creation of a European capital market and the strengthening of the European monetary system.
4. Real European cooperation — through cooperative competition; through cooperation between European companies and through a broad interpretation of the much-quoted principle of a 'fair return'.

The importance of all these elements of European integration has been repeatedly emphasized by others and by myself over the past few years. It is remarkable how much consensus there is regarding what needs to be done. It is likewise remarkable how much consensus there is with regard to the urgency with which these things need to be done. We really are running out of time. Yet it is amazing how, despite this unanimity, we are often unable to achieve our objectives. In my opinion this can mainly be attributed to an incorrect conception of how cooperation between the parties involved should be organized.

Having discussed the changing face of the future, the significance of high-tech industry, the conditions for success and the industrial policy that should be pursued, permit me to conclude this introduction by saying a few words about the principle of a 'fair return', about which so much has been said recently in European circles. For the ancient Romans this was the *quid pro quo* principle, in other words the idea that if I do something for you, I can expect you to do

something for me in return — preferably directly and preferably a bit more.

This principle may work in bilateral relations, but not in more complex relations. It stands in the way of cooperation in a wider framework.

The Nobel Prize winner Prigogine provides an interesting theoretical illustration of the problem which I have referred to here. He uses the analogy of a game of billiards to illustrate how the principle of reciprocity fails. A good billiards player can work out precisely where a ball which is struck at a certain place with a certain force will hit the third cushion. And, conversely, another good billiards player can return the ball to its original position. In billiards it is possible to speak of perfect reciprocity — at least, the difference from the original position is so slight that we scarcely notice it. It is an 'exchange situation' in which what we give and what we get back are not exactly the same but so close to one another that it makes no difference. In an everyday bilateral exchange situation this is quite normal. But you know that things are much more complicated in a trilateral exchange relationship.

But let us return to Prigogine's example. He shows that in a game played with seven cushions, or a billiards table with seven sides, the principle of reciprocity no longer works. If, before hitting the ball, I arrange with someone behind the seventh cushion that he will hit the ball back with the same force as I hit the ball, so that I get back what I have given, I shall be cheated. Prigogine demonstrates that we cannot predict where the ball will come back. Making 'exchange agreements' thus becomes impossible.

Our example assumes an ideal billiards table with seven cushions. In the European Community, however — and I am now moving from physics to socio-political reality — we work with twelve partners (*fig. 6*) and the situation is far from ideal. The amounts contributed and received are difficult to measure; it is often a comparison of apples and pears and for someone who is thirsty (or for a poor country or region) an apple means more than for someone who has apples in abundance (a rich country). Nor will the priorities always be the same in the different countries.

Cooperation in Europe therefore calls for more than a narrow conception of a 'fair return' or the *quid pro quo* principle. If we restrict ourselves to direct bilateral or trilateral relations on the grounds that with multilateral relations it is not clear what we get back in return for what we put in, then a truly integrated Europe will never come about and we shall always fall short of our potential.

A rigidly conceived principle of a 'fair return' is based on a failure to recognize the principles of reci-



(Photo: European Community)

Fig. 6. The European Community building in Brussels. Are all the 'standard-bearers' prepared to band together in 'multiple reciprocity', unconstrained by the straitjacket of a 'fair return'?

procuity in a complex social context. If, by analogy with Prigogine's example, we assume that in a complex relationship with several actors you do not know whether you will get anything back directly in return for what you put in, we have to add that according to the same theory you also receive things without being able to show exactly where they come from.

A precondition for what I would like to call 'multiple reciprocity' is that all the partners take part in the game. If we — governments and politicians, as well as companies — remain the prisoners of a nar-

rowly defined 'fair return' strategy, this will be to everyone's disadvantage. 'Prisoner's dilemma' and 'fair return' straitjackets can only be broken by mutual trust, a joint objective and belief in the common cause.

Let us hope that our children will not be able to accuse us of having displayed a lack of trust, vision and faith in Europe.

^[3] See the interesting report No. A2-109/85, Europe's response to the modern technological challenge, drawn up in 1985 by M. C. Poniatowsky for the European Parliament.

A ceramic differential-pressure transducer

V. Graeger, R. Kobs and M. Liehr

In process control the accurate measurement of pressure differences is of great importance, since both volume flows and manometric heads can be derived from differential-pressure measurements. The differential-pressure transducers now widely used have metal diaphragms, but most metals eventually become corroded by aggressive process fluids. The Philips Hamburg research laboratories, Philips Forschungslaboratorium Hamburg, have developed a transducer with a ceramic aluminium-oxide diaphragm. The transducer is now in production at Process Control Instrumentation in Kassel, a German Philips company. The ceramic material is resistant to virtually all corrosive fluids. In the production of the new transducers good use is made of thick-film technology, originally developed for the manufacture of 'hybrid' electronic circuits.

Introduction

The volume flow in a pipeline can be measured with an orifice gauge, a plate containing an orifice of accurately defined dimensions. The difference Δp between the static pressure in front of the orifice plate and behind it is a measure of the volume flow Q_v , given by

$$Q_v = \mu A \sqrt{\frac{2\Delta p}{\rho}}$$

where μ is the flow coefficient, A the area of the orifice and ρ the density of the liquid. This expression is based on the well-known Bernoulli equation [1].

The head of liquid can be determined by measuring the difference between the pressure at the top of the liquid and at the bottom. The head H is proportional to the measured pressure difference:

$$H = \frac{\Delta p}{\rho g}$$

where g is the acceleration due to gravity. In this way the volume of liquid in a reservoir can be determined, provided the dimensions of the reservoir are known. Both types of measurement, of the volume flow in a pipeline and of the volume of liquid in a reservoir, are important in process control. It is essential in both

cases that the measurements remain accurate, that the measuring instrument can withstand the temperature, pressure and chemical aggressivity of the process fluid, and that the instrument cannot cause explosions by spark discharges.

Transducers for measuring pressure differences usually have a measuring element in the form of a stainless-steel diaphragm. The pressures are applied to both sides of the diaphragm, which is deflected by the difference in pressure. Small deflections are proportional to the difference in pressure. In general it is undesirable to expose such a diaphragm directly to liquids or gases, which are frequently corrosive. The differential-pressure transducers now in common use are therefore of the dual-chamber type [2]. The two chambers are filled with an inert fluid, such as silicone oil, and the chambers are separated from the actual process fluid by auxiliary diaphragms. Dual-chamber transducers therefore have three diaphragms, the outer one of stainless steel, tantalum, titanium or monel (a nickel-copper alloy). Transducers of this type tend to be expensive because of the difficulties in welding or soldering different metals.

The single-chamber differential-pressure transducer, which is the subject of this article, is simpler in design and consists of one kind of material: ceramic, see *fig. 1*. There are only two diaphragms, which are cir-

Dr V. Graeger, Dipl.-Ing. R. Kobs and M. Liehr are with Philips GmbH Forschungslaboratorium Hamburg, Hamburg, West Germany.

cular and separated by a chamber filled with silicone oil. Process fluids on both sides of the transducer exert pressures p_1 and p_2 on the diaphragms M_1 and M_2 . If the fluid in the transducer is incompressible and the diaphragms are identical, they will both have the same deflection δ . This is equal to the deflection of an imaginary diaphragm with a stiffness equal to twice the stiffness of the individual diaphragms.

The displacement of the diaphragms can be measured by strain gauges, or with inductive or capacitive displacement transducers. We used the capacitive type, by integrating two capacitors with the transducer. Electrodes for the capacitors are applied to the inside of the diaphragms and on opposite sides of the central ceramic plate Pl . The silicone oil that 'connects' the two diaphragms flows through an opening in the plate. The four electrodes form two capacitors with capacitances C_1 and C_2 . The deflection δ at the centre of each diaphragm is given by

$$\delta = \frac{1}{2}(d_1 - d_2),$$

where d_1 and d_2 are the distances between the electrodes at the centre of the two capacitors (with the same assumptions as before). Since the deflection δ of the diaphragms is proportional to the measured pressure difference, and d_1 and d_2 are inversely proportional to the respective capacitances, we have:

$$p_2 - p_1 \propto \frac{1}{C_1} - \frac{1}{C_2}. \quad (1)$$

This equation is fundamental to the operation of the new differential pressure transducer.

The ceramic material used for the various parts of the transducer is aluminium oxide. This material has several good features besides its resistance to practically all process fluids: the deformation of the diaphragms is accurately proportional to the force, and the material is not subject to hysteresis, creep, or plastic deformation. It therefore obeys Hooke's law. The deflection of the diaphragms in the laboratory models is no more than 10 μm for a diaphragm thickness of 0.24 mm and a differential pressure range of about 25 mbar for the most sensitive transducer, and 3.5 μm at a diaphragm thickness of 1.6 mm and a measurement range of about 3000 mbar for the least-sensitive transducer. (The deformations shown in fig. 1 and the following figures are shown greatly exaggerated.) It will be clear that the tolerances for the dimensions and shape of the diaphragms and the electrodes will be very tight. It has been found that the screen-printing techniques used in the thick-film technology for manufacturing 'hybrid' electronic circuits^[3] are also eminently suitable for the application of the electrodes. When this technology is used it is also possible

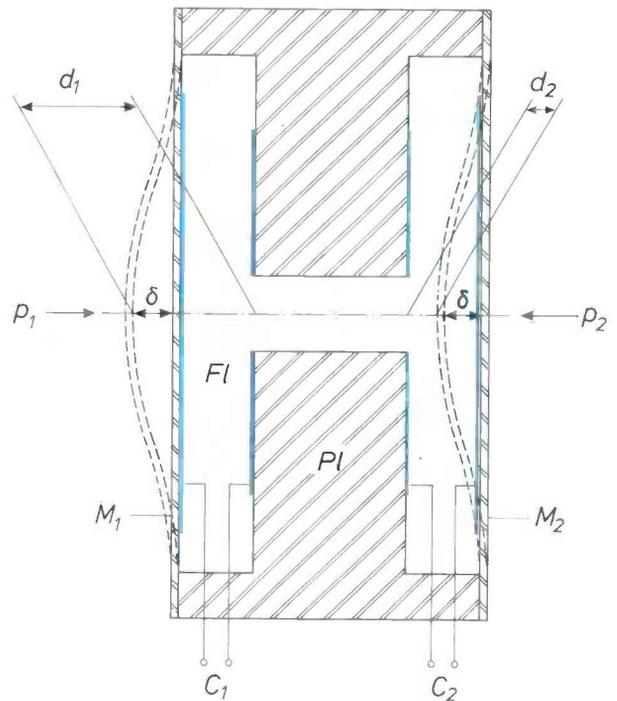


Fig. 1. Cross-section of the ceramic differential-pressure transducer. Pl central plate. Fl fluid (usually silicone oil). M_1 , M_2 diaphragms. C_1 , C_2 capacitances of the two capacitors with which the deflections of the diaphragms are measured; the capacitor electrodes are shown in blue. p_1 , p_2 pressures of the process fluid. The transducer measures the differential pressure $\Delta p = p_2 - p_1$. δ deflection in the middle of the diaphragms. d_1 , d_2 electrode spacings at the centre of the two capacitors. The measured differential pressure Δp is proportional to $1/C_1 - 1/C_2$.

to use 'glass solder' for hermetically sealing the diaphragms to the rest of the transducer. Fig. 2 shows that our ceramic differential-pressure transducer is much more compact than conventional types based on the dual-chamber principle.

When the transducer is filled with silicone oil, the oil pressure is made slightly higher than atmospheric pressure, so that the diaphragms are initially deflected outwards. When the transducer is in use, an increase in temperature or a drop in the mean pressure of the process fluid causes the sum $d_1 + d_2$ of the distances between the electrodes to increase. Although to a first approximation this does not introduce any error in the measured pressure difference given by (1), the accompanying change in the pressure and hence in the permittivity of the liquid in the transducer does cause a change in the measured result. It will be clear from the above that the quantity $d_1 + d_2$ can be used for deter-

- [1] J. Hengstenberg, B. Sturm and O. Winkler, Messen, Steuern und Regeln in der chemischen Technik, Band I, Springer, Berlin 1980 (in German).
 [2] See page 250 of H. N. Norton, Sensor and analyzer handbook, Prentice Hall, Englewood Cliffs, NJ, 1982.
 [3] W. Funk, Thick-film technology, Philips Tech. Rev. 35, 144-150, 1975.



a



b

Fig. 2. a) Conventional differential-pressure transducers, made of metal and operating on the dual-chamber principle. b) The new ceramic differential-pressure transducer, operating on the single-chamber principle; see fig. 1. Its diameter is 60 mm.

mining the change in the temperature and pressure of the liquid in the transducer. The effect of the temperature and mean pressure of the process fluid on the measured result is compensated by measuring the quantity $1/C_1 + 1/C_2$ as well as $1/C_1 - 1/C_2$ as in (1). In this way the effect of a temperature change of 10 °C on the measured result can be limited to less than 0.1% of the measurement range.

The new ceramic transducer is fitted in the type PD3 differential-pressure transmitter recently put into production at Process Control Instrumentation in Kassel, a German Philips company that forms part of the Philips Industrial and Electro-Acoustic Systems Division. There are six different types: the smallest adjustable measurement range has a maximum value of 10 mbar of pressure difference; the largest has a maximum value of 3000 mbar of pressure difference.

In the following we shall first consider the theoretical background of the differential-pressure transducer, and then we shall briefly discuss its technology and construction. Next we shall deal with the temperature compensation and discuss some measurement results. Finally we shall look briefly at a number of future developments.

Theoretical background

First of all we shall consider the magnitudes of the resultant pressures p_{M1} and p_{M2} acting on the two diaphragms (see fig. 3). As we have seen, the internal pressure p_1 of the fluid is such that initially both diaphragms are pushed outwards. If the pressure transducer were completely symmetrical, so that the diaphragms had exactly the same dimensions and material properties, each diaphragm would take up half the pressure difference $\Delta p = p_2 - p_1$. With the assumption that the internal pressure is a function of the transducer temperature t and the mean pressure $p = (p_2 + p_1)/2$ of the process fluid, then we have the following expressions for p_{M1} and p_{M2} :

$$p_{M1} = p_1(t, p) + \frac{1}{2}(p_2 - p_1), \tag{2a}$$

$$p_{M2} = p_1(t, p) - \frac{1}{2}(p_2 - p_1). \tag{2b}$$

The pressures acting on the surfaces of the diaphragms are taken as positive in the usual convention.

An increase dt in the temperature of the liquid will entail an increase in the liquid volume V_0 . The extra

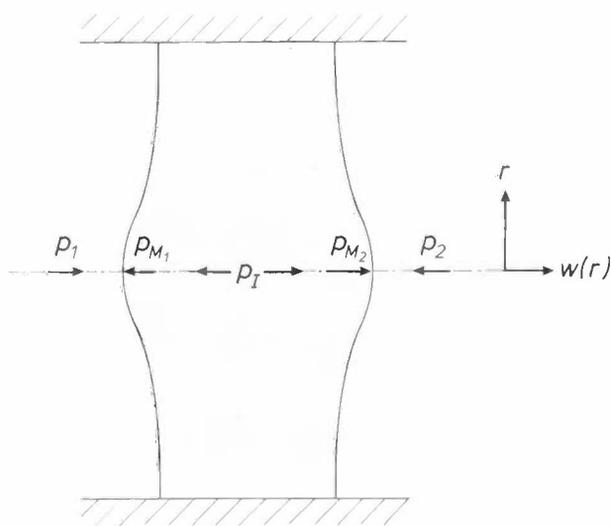


Fig. 3. Definition of some of the quantities used in calculating the output quantity $1/C_1 - 1/C_2$ as a function of the differential pressure $\Delta p = p_2 - p_1$ (equation 5). p_1 pressure of the liquid in the transducer, p_{M1} , p_{M2} resultant pressures on the diaphragms. r radius with respect to the axis of symmetry. $w(r)$ deflection of a diaphragm as a function of r , with $w = 0$ in the undeflected state.

volume $V_0 \alpha dt$, where α is the thermal expansion coefficient, will cause an additional deflection of the diaphragms and will increase the internal pressure p_1 . An increase dp in the mean pressure p of the process fluid will entail a decrease in the liquid volume by $V_0 dp/K$, where K is the bulk modulus of the liquid in the transducer. The effects of dp and dt show that it is desirable to keep the liquid volume V_0 as small as possible.

The shape of the diaphragms is determined by the function $w(r)$, which represents the deflection as a function of the radius, where $w = 0$ when there is no pressure difference across the diaphragms. This function follows from the theory of flat plates under a bending load [4]. The equations for the deflections of the diaphragms are:

$$w_1(r) = -p_{M_1} F \left(1 - \frac{r^2}{R_o^2}\right)^2, \quad (3a)$$

$$w_2(r) = p_{M_2} F \left(1 - \frac{r^2}{R_o^2}\right)^2, \quad (3b)$$

where

$$F = \frac{3(1 - \nu^2)R_o^4}{16Eh^3},$$

and R_o is the half-diameter and h the thickness of the diaphragms. E is Young's modulus and ν Poisson's ratio for the material of the diaphragms. Asymmetries can be taken into account by adding a subscript 1 or 2 to F , R_o and h for the left-hand or right-hand diaphragms.

For an imperfectly symmetrical transducer the following expressions apply for the capacitance of the capacitors C_1 and C_2 in fig. 1:

$$\frac{1}{C_1} = \frac{d_{0_1} + cF_1 p_{M_1}}{\epsilon_r \epsilon_0 \pi r_a^2}, \quad (4a)$$

$$\frac{1}{C_2} = \frac{d_{0_2} + cF_2 p_{M_2}}{\epsilon_r \epsilon_0 \pi r_a^2}, \quad (4b)$$

where d_{0_1} and d_{0_2} are the spacings of the capacitor plates for $p_{M_1} = p_{M_2} = 0$, r_a is the mean half-diameter of the capacitor electrodes, ϵ_r is the relative permittivity of the liquid in the transducer and ϵ_0 the permittivity of free space. The correction factor c accounts for the curvature of the diaphragm surface; for $r_a = R_o/2$ the correction is $c = 0.8$. Combining equations (2) and (4) gives an expression for the output quantity of the transducer:

$$\frac{1}{C_1} - \frac{1}{C_2} = \frac{d_{0_1} - d_{0_2} + c(F_1 - F_2)p_1(t,p)}{\epsilon_r \epsilon_0 \pi r_a^2} + \frac{c(F_1 + F_2)(p_2 - p_1)}{2\epsilon_r \epsilon_0 \pi r_a^2}. \quad (5)$$

The first term in this equation gives the zero error and the second term the sensitivity of the differential pressure transducer. For a symmetrical transducer we have $d_{0_1} = d_{0_2}$ and $F_1 = F_2 = F$, so that (5) becomes:

$$\frac{1}{C_1} - \frac{1}{C_2} = \frac{cF(p_2 - p_1)}{\epsilon_r \epsilon_0 \pi r_a^2}. \quad (6)$$

Equation (5) shows that the sensitivity is proportional to the sum $F_1 + F_2$ of the reciprocals of the stiffnesses of the diaphragms. Equation (5) therefore helps to explain the effect of the temperature t and the mean pressure p on the zero error and the sensitivity. The quantities p and t affect not only the internal pressure p_1 but also the relative permittivity ϵ_r , since ϵ_r is a function of the density of the liquid.

Technology and construction

As we have said, the transducer, apart from the electrodes, is made completely of a ceramic; see fig. 4. The material of the central plate Pl and the diaphragms M_1 and M_2 consists of alumina: 99.5% sintered aluminium oxide with a residue of vitreous constituents.

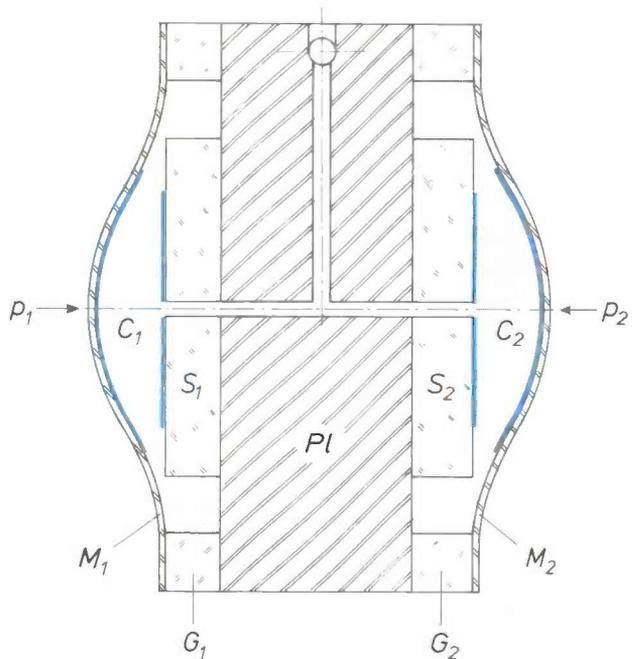


Fig. 4. Construction of the differential-pressure transducer. G_1 , G_2 glass-solder rings. S_1 , S_2 layers forming substrates for the inner electrodes. G_1 , G_2 , S_1 and S_2 are applied by screen printing (as in thick-film technology). See also the caption to fig. 1.

[4] S. P. Timoshenko and S. Woinowsky-Krieger, Theory of plates and shells, 2nd edition, McGraw-Hill, New York 1959; R. J. Roark and W. C. Young, Formulas for stress and strain, 5th edition, McGraw-Hill, New York 1975.

The permissible bending stress and the modulus of elasticity of this material are high. (Its Young's modulus is about 1.6 times that of steel.) The diaphragms do not therefore need reinforcing ribs. The central plate contains two holes, one axial and the other radial, both with a diameter of only 0.8 mm. The transducer is filled with silicone oil through the radial hole, which is then closed with a ball stopper.

The use of a ceramic makes it possible to use the thick-film technology^[3]. Coatings with an accurately determined thickness between 3 and 15 μm are deposited on a substrate by screen printing, and these coatings are then bonded to the substrate by sintering in a furnace. The experience gained in the electronics industry with the thick-film technology in the manufacture of hybrid circuits has been most useful in the design of our differential-pressure transducer. The process is inexpensive and there are many pastes available for making coatings with completely different physical properties. Specialities of the thick-film technology include:

- The metallization of ceramic with a minimum line-width of 100 μm .
- The insulation of surfaces by the application of a coating of glass ceramic.
- The coating and bonding of ceramic components with 'glass solder' (or 'melt glass'), which will give a glassy joint at 400 °C.

These features are turned to advantage in the following procedure. First the diaphragms are metallized locally. Next, glass-solder rings (G_1 and G_2) are applied to the central plate. The internal diameter of these rings determines the stiffness of the diaphragms, which means that these rings have to be very accurately made. Coatings of glass ceramic (S_1 and S_2) are then applied, which are only a little thinner than the glass-solder rings. The layers S_1 and S_2 limit the movement of the diaphragms and act as substrate for the inner electrodes that are now applied. At the same time electrical connections to the outside of the transducer are made for connecting the capacitors. Finally, the complete assembly is heated in a furnace; the glass solder melts to produce a hermetic and geometrically accurate seal.

Before the oil is introduced into a transducer the electrodes remain in contact, in pairs. The pressure of the silicone oil filling the transducer at the end of the manufacturing process determines the pre-stress of the diaphragms, and hence the electrode spacing at zero external pressure. In practice the oil pressure is set so that the electrodes of one of the capacitors just touch at an external differential pressure $p_2 - p_1$ of 1.3 times the range. As we noted earlier, the ampli-

tude of the diaphragm excursion is 3.5 to 10 μm . The outer diameter of the diaphragms is 60 mm and the oil volume is 60 mm^3 . The capacitance of each capacitor at zero external pressure is about 400 pF.

Temperature compensation and results of measurements

Fig. 5a shows the results of measurements on a differential-pressure transducer with a range of 500 mbar, in the positive and the negative directions. The output quantity $1/C_1 - 1/C_2$ is plotted as a function of the differential pressure $\Delta p = p_2 - p_1$, with the temperature as parameter. There is a zero error for the linear relation between $1/C_1 - 1/C_2$ and Δp corresponding to equation (5) and the sensitivity is temperature-dependent. For a temperature of 20 °C fig. 5b shows the relative deviation f of the measured points from the straight line of best fit; these results were ob-

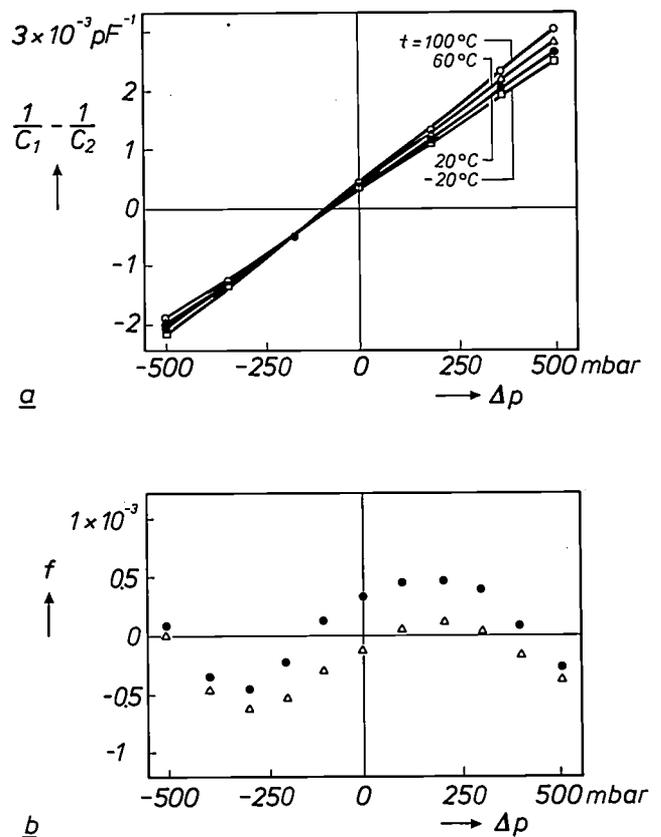
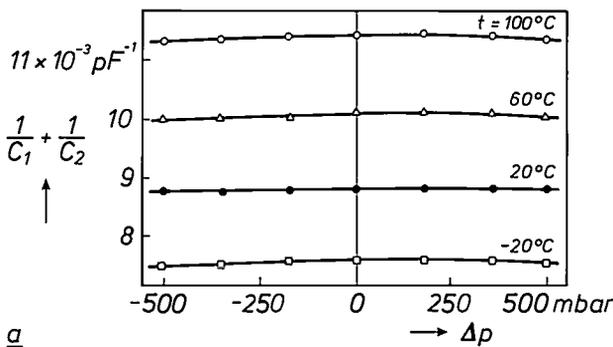


Fig. 5. a) Results of measurements on a transducer with a measuring range of 500 mbar in the positive and negative directions. The output quantity $1/C_1 - 1/C_2$ is plotted as a function of the differential pressure Δp , with the temperature t as parameter. b) The relative deviation f as a function of Δp at 20 °C. f is related to the measuring range and is the deviation of the measured points with respect to the straight line of best fit. The measured points are the results of a measurement cycle 0 \rightarrow 500 \rightarrow 0 \rightarrow -500 \rightarrow 0 mbar.

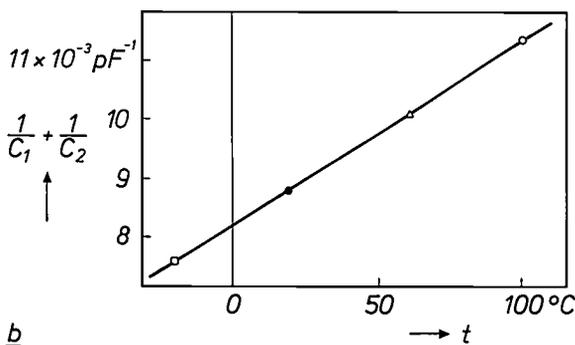
tained on going through the measurement range in the sequence $0 \rightarrow 500 \rightarrow 0 \rightarrow -500 \rightarrow 0$ mbar. The relative deviation f is related to the range, and the result of the measurements gives some idea of the hysteresis of the transducer. On varying the temperature the relative change in the zero error in fig. 5a ($1/C_1 - 1/C_2$ for $\Delta p = 0$) is 0.4% for every 10 °C and the relative change in the sensitivity is 1% for every 10 °C. The effect of increasing the mean pressure $p = (p_1 + p_2)/2$ by 100 bars corresponds to reducing the temperature by about 9 °C. The errors are too large for a differential-pressure transducer for industrial application, so that it is necessary to compensate for the effects of temperature and pressure.

When the temperature of the process fluid rises or the mean pressure drops the deflection of the diaphragms increases, so that the quantity $d_1 + d_2$ assumes a higher value. This quantity is proportional to $1/C_1 + 1/C_2$, and by combining equations (2) and (4) we see that

$$\frac{1}{C_1} + \frac{1}{C_2} = \frac{d_{01} + d_{02} + 2cFp_1(t,p)}{\epsilon_r \epsilon_0 \pi r_a^2} \quad (7)$$



a



b

Fig. 6. a) The quantity $1/C_1 + 1/C_2$ as a function of the differential pressure Δp , again with the temperature t as parameter. b) The quantity $1/C_1 + 1/C_2$ as a function of the temperature t at $\Delta p = 0$.

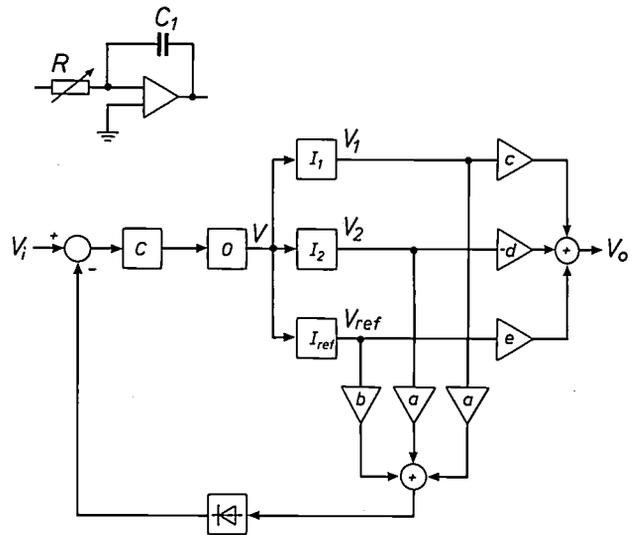


Fig. 7. Block diagram of a circuit for compensating linear temperature effects. V_1 direct voltage at the input. C controller. O oscillator with stable angular frequency. V amplitude of the alternating voltage at the oscillator output. I_1 , I_2 and I_{ref} integrators (see inset) with alternating output voltages that have amplitudes V_1 , V_2 and V_{ref} respectively. a , b , c , $-d$ and e voltage gains. V_o amplitude of the alternating output voltage. V_o is proportional to the differential pressure Δp and is to a first order independent of temperature.

It is assumed here that the transducer is perfectly symmetrical, so that $F_1 = F_2 = F$. The quantity $1/C_1 + 1/C_2$ can thus be used for measuring the temperature t and the mean pressure p . Fig. 6a shows that this quantity is in fact closely dependent on temperature and hardly affected at all by the differential pressure Δp . The temperature dependence can be approximated by a linear function:

$$\frac{1}{C_1} + \frac{1}{C_2} = a_0 + a_1 t, \quad (8)$$

where a_0 and a_1 are constants that can be determined by calibration; see fig. 6b. Similarly, in the relation between the output quantity $1/C_1 - 1/C_2$ and the differential pressure Δp to be measured the effect of the temperature t on the zero error and the sensitivity can be approximated by linear functions:

$$\frac{1}{C_1} - \frac{1}{C_2} = a_2 + a_3 t + (a_4 + a_5 t)(p_2 - p_1), \quad (9)$$

where a_2 , a_3 , a_4 and a_5 are constants.

Fig. 7 shows a circuit that we have designed to compensate for linear temperature effects. I_1 , I_2 and I_{ref} are integrators, which each consist of a resistor, a capacitor and an operational amplifier (see inset). I_1 contains the capacitance C_1 of the differential-pressure transducer and I_2 contains C_2 . If V denotes the amplitude of the output voltage of the oscillator O

with stable angular frequency ω , the amplitudes of the alternating voltages after the integrators are

$$V_1 = \frac{V}{\omega RC_1},$$

$$V_2 = \frac{V}{\omega RC_2},$$

$$V_{\text{ref}} = \frac{V}{\omega \tau_{\text{ref}}},$$

where R is the — variable — resistance in I_1 and I_2 , and τ_{ref} is the time constant of I_{ref} . The three alternating voltages after the integrators are multiplied by the factors a , a and b respectively, summed and then rectified. The result is fed back and compared with a direct voltage V_1 at the input. The factors b and a are chosen so as to satisfy the equation

$$b = \frac{a \tau_{\text{ref}}}{R} \left(\frac{a_1 a_4}{a_5} - a_0 \right),$$

which contains constants from (8) and (9). Calcula-

tions show that the feedback signal before rectification has an amplitude

$$\frac{V a a_1}{\omega R a_5} (a_4 + a_5 t).$$

The feedback signal is thus proportional to the factor in (9) that represents the sensitivity. The feedback and the action of the controller C therefore keep the value of $a_4 + a_5 t$ constant, so that the linear temperature-dependence of the sensitivity is eliminated.

The controller C affects the amplitude of the oscillator O ; we assume that the frequency of the oscillator is sufficiently constant. The output signal of the circuit is obtained by multiplying the signals after the integrators, by c , $-d$ and e respectively and then summing them. Calculations in which equation (8) is used show that the amplitude V_o of the output signal follows from the relation

$$V_o = \frac{V}{\omega R} \left\{ \frac{1}{C_1} - \frac{1}{C_2} - (a_2 + a_3 t) \right\}, \quad (10)$$

when the factors c , d and e are chosen so as to satisfy:

$$c = 1 - \frac{a_3}{a_1},$$

$$d = 1 + \frac{a_3}{a_1}$$

and

$$e = \frac{\tau_{\text{ref}}}{R} \left(\frac{a_0 a_3}{a_1} - a_2 \right).$$

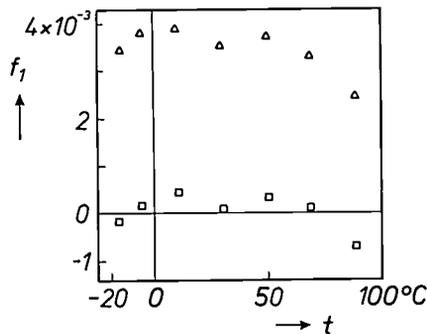
Using eq. (9) we can write eq. (10) as

$$V_o = \frac{V}{\omega R} (a_4 + a_5 t)(p_2 - p_1). \quad (11)$$

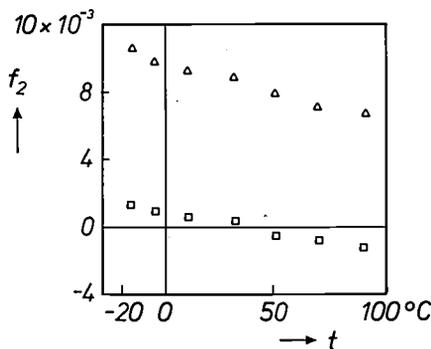
Since the sensitivity $a_4 + a_5 t$ is kept constant by feedback and control, the amplitude of the output signal is proportional to the differential pressure to be measured.

Fig. 8 shows the improvement achieved by introducing temperature compensation. For a temperature change of 10 °C fig. 8a gives the relative change f_1 in the zero error and fig. 8b the relative change f_2 in the sensitivity, both with and without temperature compensation. A roughly tenfold improvement is obtained.

Fig. 9 shows the Philips differential-pressure transmitter type PD3, which incorporates the transducer with the temperature-compensation circuit described here. The measuring range can be adjusted, for example from 20 to 100 mbar for the type with a maximum range of 100 mbar of differential pressure. The maximum error for the measured Δp value with this type is guaranteed to be less than 0.2% of the



a



b

Fig. 8. The improvement obtained with the circuit in fig. 7. a) The relative change f_1 in the zero error ($1/C_1 - 1/C_2$ at $\Delta p = 0$, see fig. 5a) for a temperature variation of 10 °C, as a function of temperature t . The upper and lower points are the results of measurements without and with the compensation circuit respectively. b) The relative change f_2 in the sensitivity for a temperature variation of 10 °C, again with and without compensation circuit.

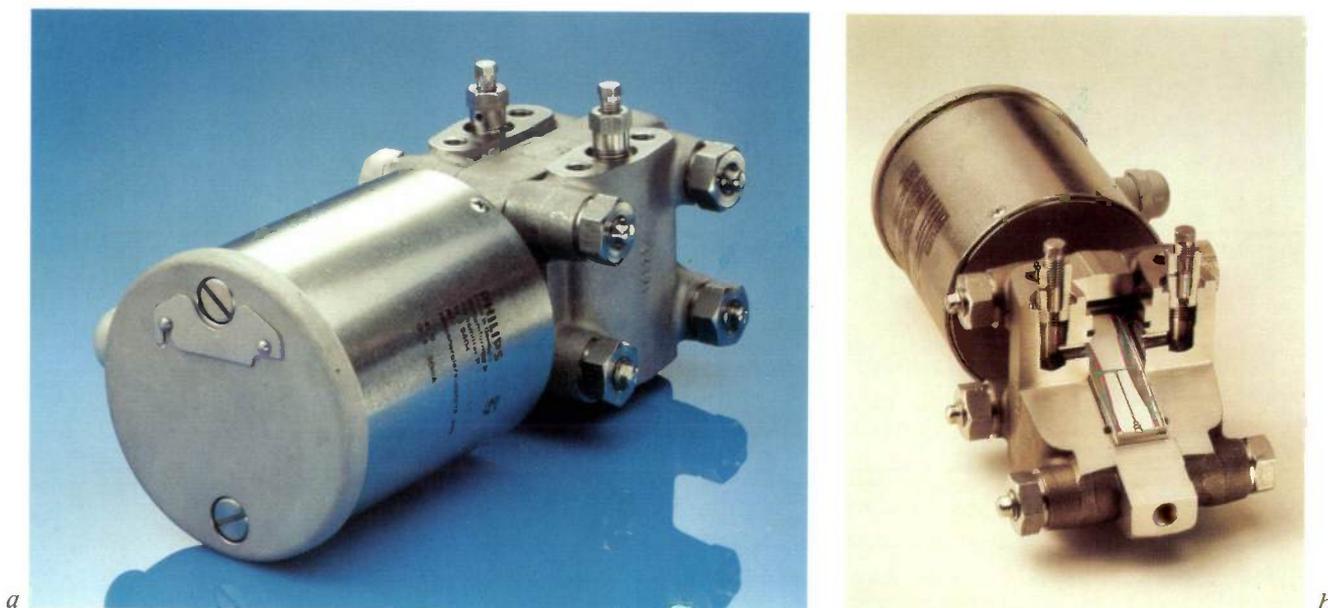


Fig. 9. *a*) The type PD3 differential-pressure transmitter, complete with transducer and temperature compensation circuit, produced by Process Control Instrumentation of Kassel, West Germany, a German Philips company. *b*) As (*a*), now partly cut away to show the transducer.

measuring range. This applies for a temperature range from $-30\text{ }^{\circ}\text{C}$ to $80\text{ }^{\circ}\text{C}$. The maximum direct-current output signal is 20 mA, which is proportional to the measured differential pressure or — for flow-rate measurements — to the root of the differential pressure.

Further developments

Higher-order temperature and pressure effects can be compensated with the aid of a microprocessor. This can also be used for calculating the mass flow from the measured differential pressure, the temperature and static pressure of the process fluid, and the geometry of the orifice plate. A device that combines a differential-pressure transducer and a microprocessor can also generate error signals if there is a malfunction. The data for the temperature compensation do not have to be entered by means of resistors but can be stored as numbers in a PROM (Programmable Read-Only Memory). The device can supply the results of the measurements to a standard data-transmission system (a 'bus'). Experimental versions of a differential-pressure transducer for laboratory use are now ready; the device has an exceptionally high accuracy — ten times better than the instrument developed for process control — and is suitable for a bus system complying with the standards IEEE 488 or IEC 625.

Another promising development is a differential-pressure transducer that delivers the results of its measurements along an optical-fibre link [6]. This is not sensitive to interference from electromagnetic

fields, which can be a problem with conventional copper wires. Such interference can be particularly serious when measurements are carried out in the vicinity of electric motors or high-power transformers. Since a transducer cannot generate an optical signal directly in a glass fibre, we use a method in which the transducer generates an optical signal via an electronic circuit containing LEDs (Light-Emitting Diodes). Our circuit is made from components in CMOS technology (Complementary Metal-Oxide Semiconductor), and takes only $50\text{ }\mu\text{A}$, so that the circuit can operate continuously for more than ten years from lithium batteries.

[6] J. Kordts, V. Graeger and G. Martens, *Hard & Soft (Mikroperipherik)* No. 4/86, 7, 1986.

Summary. The differential-pressure transducer developed at Philips Forschungslaboratorium Hamburg is made completely of alumina, and is therefore resistant to practically all corrosive fluids used in process engineering. The transducer has a single liquid-filled chamber, isolated by two diaphragms from the fluids whose differential pressure is to be measured. The deflections of these diaphragms are measured by a capacitive method. Four electrodes, applied to the two diaphragms and on opposite sides of a central plate, form two capacitors. The measured differential pressure is proportional to the difference of the reciprocals of their capacitances. This difference is affected, however, by the temperature of the liquid in the transducer. The effect is compensated by using the sum of the reciprocals of the capacitances, which is proportional to the temperature. The linear temperature compensation is produced by an electronic circuit that includes operational amplifiers and a feedback loop. This circuit and the ceramic transducer are incorporated in the Philips PD3 differential-pressure transmitter, which is now on the market in six models with different measuring ranges. Experimental versions are now ready of a differential-pressure transducer for laboratory use, in which nonlinear temperature effects can also be compensated. A method is being studied for transmitting the results of measurements along optical fibres.

1937

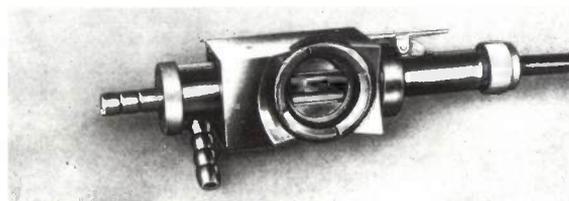
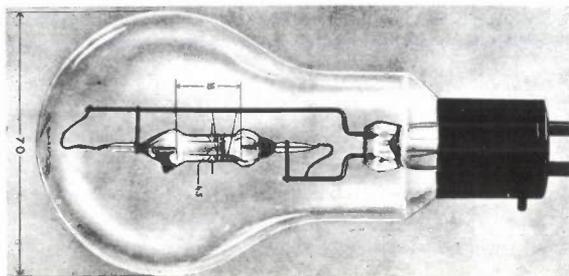
THEN AND NOW

1987

High-pressure mercury lamps

Some fifty years ago the first Philips high-pressure mercury lamp ^[*], the HP 300, with an electrical power of 75 watts and a luminous flux of 3000 lumens (upper right), was put on the market. It had about three times the luminous efficacy of a comparable incandescent lamp. Shortly after this a related lamp, the SP 500 W, with a power of 500 W and a luminous flux of 30 000 lm was introduced. This lamp ^[**] had to be water-cooled (centre right). Because of the small dimensions of the actual light source it gave a very high luminance, as required in projector lamps and searchlights. The SP 500 W was manufactured almost unmodified until a few years ago.

Today high-pressure mercury lamps are widely used in road lighting, factories and various other places where plenty of bright 'white' light with reasonably good colour rendering is required. Although the efficacy, life, reliability and colour rendering of these lamps are all much improved, they still have a close resemblance to their 'ancestor', the HP 300, especially internally (lower right). The inside of the glass envelope is now often coated with fluorescent powder. The power of the modern Philips high-pressure mercury lamps (no water cooling now) is between



50 and 2000 W, while the associated luminous flux is between 2000 and 125 000 lm.

For about 25 years Philips have also manufactured high-pressure mercury lamps containing accurately measured amounts of the iodides of sodium, thallium and indium ('metal-halide lamps'). These additives produce a 'whiter' white (large photo) and improve the efficacy by about 50%.

[*] All mercury lamps with a gas pressure higher than 1 bar are called high-pressure mercury lamps (or mercury-vapour lamps). For lighting, however, a pressure of 20 bars and above (sometimes as high as 200 bars) gives the best colour rendering and efficacy.

[**] From Philips Technical Review, June 1937.

Interactive MR image synthesis

M. H. Kuhn and W. Menhardt

In proton magnetic-resonance tomography (MR) the images can be obtained by the inversion-recovery method or by the spin-echo method. The contrast in these images is very dependent on the time parameters of the pulse trains used for obtaining the signals that ultimately produce the MR images. The optimum image contrast to use depends on the pathological abnormalities the observer is looking for. The image contrast required, and hence the value of the various parameters, is not always known beforehand, however. Researchers at Philips Forschungslaboratorium Hamburg, the Philips Research Laboratories in Hamburg, have now designed an image-processing procedure that permits the observer to synthesize, in real time, an image for a different setting of these parameters, without any need for a new MR exposure. The observer can thus experiment from the MR control console, to find the settings that give the contrast ratios he requires. It is also possible to use colour in the image. A further valuable aid is a facility for indicating two areas in the image that should give the maximum contrast.

Introduction

The most outstanding advantage of proton magnetic-resonance tomography (MR) is that it enables images to be obtained of cross-sections of the human body without causing any physiological damage. *Fig. 1* shows a Philips Gyroscan S5 installation that gives such images. Other advantages of magnetic-resonance tomography are that the difference between diseased and healthy tissue can often be made more visible than with other diagnostic methods, that in general no contrast agent is necessary and that the contrast in the image can be manipulated fairly easily. A disadvantage is that MR equipment is still expensive. This is mainly because of the need for a highly uniform and powerful magnetic field, now more likely than not produced by superconducting coils. This means that to reach and maintain the very low temperature required (near absolute zero) it is necessary to use liquid helium.

In proton MR tomography use is made of the precession of the spin axis of protons that form the nucleus of hydrogen atoms in the tissue. The fre-

quency of this 'Larmor precession' is proportional to the strength of the constant magnetic field. The precessing protons induce r.f. signals in a detector coil. When a small gradient is introduced into the constant field the location of the hydrogen atoms that are responsible for the detected signals can be determined from the frequency.

From the detected signals various parameters can be derived that when taken together are fairly characteristic of the relevant tissue. These parameters are the proton density ρ , the relaxation time T_1 of the longitudinal magnetization and the relaxation time T_2 of the transverse magnetization. Other parameters that can be derived relate to the rate of blood circulation and the chemical composition, but these will not be considered here. Later we shall return to the theoretical basis for the relaxation times T_1 and T_2 . All we need to say here is that it has been found that healthy and diseased tissue often show a remarkable difference in T_1 and T_2 , even when the proton density ρ is about the same for both kinds of tissue.

The magnitude of the detected r.f. signals is a function not only of the proton density ρ and the relaxa-

Dr M. H. Kuhn and W. Menhardt are with Philips GmbH Forschungslaboratorium Hamburg, Hamburg, West Germany.

tion times T_1 and T_2 but also of the 'pulse parameters'. The term pulse parameter refers to the time intervals between the successive pulses in a pulse train that is periodically applied to excite the Larmor precession of the protons. From equations giving the signal magnitude as a function of these quantities it follows that the pulse parameters greatly affect the magnitude of the detected signals and can even change their polarity in certain cases. The choice of the pulse parameters

be obtained without the need to make a further MR scan. The observer (normally a medical practitioner) can then experiment to determine the combination of pulse parameters that produces contrast ratios that give the best visualization of certain lesions. The difficulty with this procedure until now has been that calculating new grey levels for 256×256 pixels — in our case — required a great deal of computer time, which meant a relatively long wait for results.

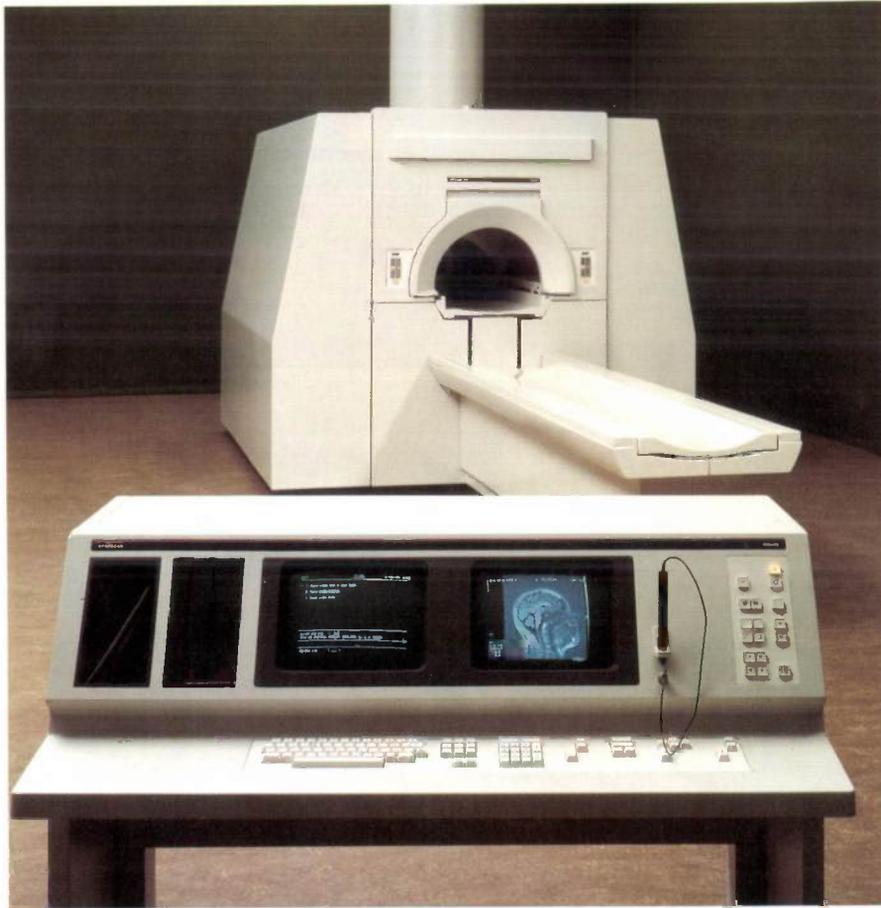


Fig. 1. Philips Gyroscan S5 MR installation, which produces images of 'slices' of the human body by proton magnetic resonance.

therefore has an important bearing on the grey levels in the resultant image and hence on the contrast between the tissues. Changing the pulse parameters, for example, has been found to reverse the contrast between grey and white brain matter.

When an MR image has been formed that is based on an exposure with a certain set of values for the pulse parameters, it is possible to use the equations mentioned above to convert the grey levels obtained for each picture element (pixel) into grey levels for other pulse parameters. The advantage of this procedure is that images with different contrast ratios can

Efforts have been made for some time to find methods in which calculations of images with different pulse parameters can be made virtually in 'real time', i.e. almost instantaneously. The most common method uses RAMs (random-access memories) organized in the form of a look-up table. Each binary number that corresponds to a pixel and is stored in an image memory with 256×256 locations is transformed with the aid of the RAM look-up table. If the binary number for each pixel has a length of n bits, the look-up table contains 2^n locations. In calculating a new image the binary number serves as the address for the look-

up table. A binary number for a new pixel is then stored at each address. Until now this procedure could only be used for one of the quantities T_1 , T_2 or ρ . Another possibility was to use a complicated procedure with three image memories for each of these quantities, including a network of RAM look-up tables and digital adders and multipliers.

For our MR system at Philips Forschungslaboratorium Hamburg we have developed a method in which numbers for T_1 , T_2 and ρ are stored together at each location in the image memory. Each location in this memory has a magnitude of ten bits. Since the original digital values for T_1 , T_2 and ρ have eight bits, this means that we can only use three, three and four bits respectively of these values. Additional quantization noise is caused by the rounding off made necessary by the availability of only $2^3 = 8$ levels for T_1 and T_2 and of $2^4 = 16$ levels for ρ . We have minimized the effects of this noise by using an algorithm proposed by R. Floyd and L. Steinberg^[1] in the conversion from eight to three or four bits. This algorithm distributes the quantization error of a pixel over the adjacent pixels such that the mean grey level of a group of pixels approximates to the grey level of the pixel before quantization.

The inclusion of this image-processing method enables the operator to adjust the pulse parameters from the console to obtain the optimum image contrast. An additional advantage is that the effect of varying the pulse parameters can easily be demonstrated to others, thus making the method particularly suitable for instructional use. The magnitude of the pulse parameters follows from the length of a number of bars on the monitor screen; see fig. 2. The computer program positions a cursor at the end of the bar for the pulse parameter that the operator wishes to modify. The operator can move the cursor by a 'track-ball' interactive device, incorporated in the operating console. After he has indicated the magnitude of the modified pulse parameter by moving the cursor to the left or right, the new image appears on the screen almost immediately — within 60 to 90 ms — so that the observer can evaluate the result at once. Once the required image contrast has been found, the observer can have the image calculated with the original 8-bit values of T_1 , T_2 and ρ . This image is virtually free from quantization artefacts. This computing procedure, however, requires about half a minute before the image appears. Finally, a photograph of the image can be taken with a 'matrix camera'.

The method described has recently been extended to include a number of new features. The first is contrast maximization. The observer uses a light pen to indicate two pixels between which the contrast should be

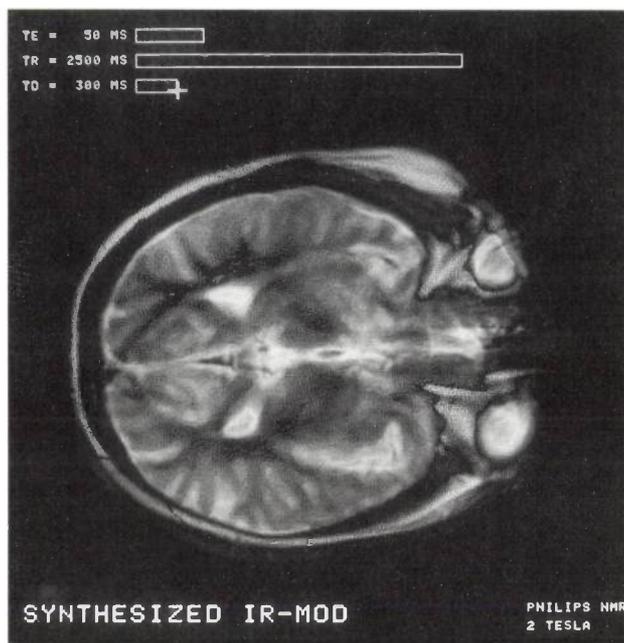


Fig. 2. Monitor screen of the MR installation at the Hamburg laboratories. The screen shows a cranial image. The magnitude of the pulse parameters used for making the image is given by the length of the bars at the top of the photograph. The observer is able to modify these parameters by moving the cursor along the bars.

as high as possible. The computer then calculates pulse-parameter values for which the maximum difference between the grey levels of the pixels is obtained. This gives the observer a better view of particular anatomical or pathological details. The second new feature is the use of colour. The MR image is then composed of differently coloured ρ , T_1 and T_2 components. The brightnesses correspond to the appropriate values of ρ , T_1 or T_2 . Since many biological tissues can be identified by a unique combination of ρ , T_1 and T_2 , the result is that colours appear in the composite image that are characteristic of specific tissues. After some practice an observer can recognize particular anatomical details at a glance and in many cases distinguish between malignant and benign tissue. In future, after the values of ρ , T_1 and T_2 that characterize specific tissues have been accurately charted, this method should evolve to become a kind of tissue-characterization system. The medical user will then be able to point at a particular region in an MR image with a light pen, and the computer will tell him the type of tissue likely to be located in that region.

In the following we shall first touch on the theoretical background of proton magnetic resonance. We shall then take a closer look at the procedures used in

[1] R. Floyd and L. Steinberg, An adaptive algorithm for spatial grey scale, SID 75 Digest, Washington, DC, 1975, pp. 36-37.

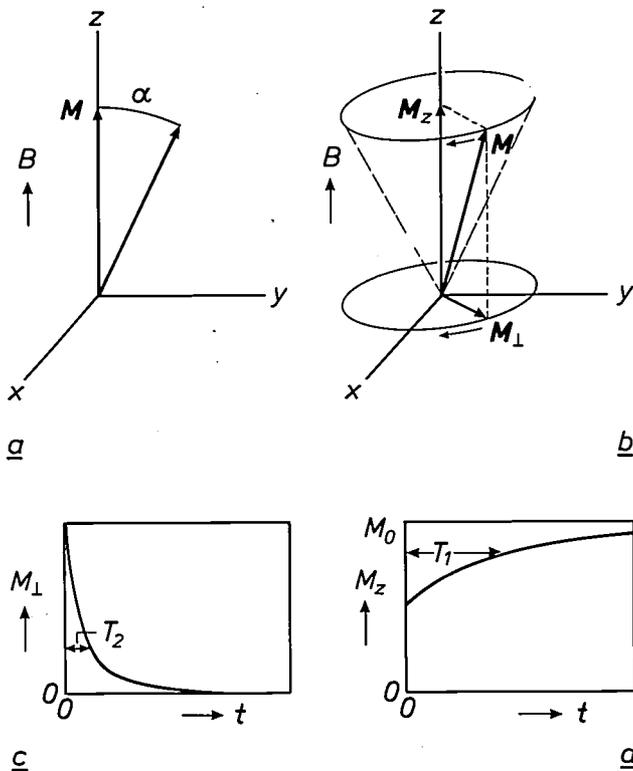


Fig. 3. The excitation of the Larmor precession. *a*) The brief application of an r.f. magnetic field in the *x*-direction at a frequency equal to that of the Larmor precession makes the magnetization **M** of the material rotate through an angle α in the (*y,z*)-plane (α pulse). The constant magnetic field of flux density *B* has the direction of the *z*-axis. *b*) The vector **M** then precesses about the *z*-axis. The longitudinal magnetization M_z does not change direction but the transverse magnetization M_1 does. *c*) The transverse magnetization M_1 decays exponentially with time constant T_2 ; *t* time. *d*) Recovery of the longitudinal magnetization. M_z increases until thermal equilibrium is reached with magnetization M_0 . The recovery follows an exponential curve with time constant T_1 , which is always greater than T_2 .

the new method of interactive MR image synthesis. Finally, some examples of applications will be discussed and a synthesized image in colour will be shown.

Theoretical principles of proton magnetic resonance

Proton magnetic resonance is based on the effect in which the axis about which protons spin describes a precessional motion about the direction of a constant magnetic field of flux density *B*. The angular frequency ω of this 'Larmor precession' is proportional to the flux density:

$$\omega = \gamma B. \tag{1}$$

The quantity γ is called the gyromagnetic ratio of the proton. An r.f. magnetic field perpendicular to the constant field excites the precession if its frequency is equal to the precessional frequency $\omega/2\pi$, which fol-

lows from eq. (1). This is the basis of proton magnetic resonance [2].

An object in thermal equilibrium in the constant magnetic field has a nuclear magnetization in the direction of the field, since there are more spins aligned with the field than spins aligned in the opposite sense. The magnetization, which is the sum of the magnetic moments per unit volume, is very small; at $B = 0.3\text{T}$ it is only 10^{-6} of the theoretical saturation value with all the spins in the same direction. The application of the r.f. field referred to above for a short time (an 'r.f. pulse') makes the magnetization **M** rotate through an angle α , see fig. 3*a*. Because of the Larmor precession of the individual protons, the vector **M** starts to precess as well (fig. 3*b*). After some time this precession decays away as the system returns to the state of thermal equilibrium. The longitudinal magnetization M_z , which does not change direction, has a slower relaxation to thermal equilibrium than the transverse relaxation M_1 , which rotates in the (*x,y*)-plane (fig. 3*c* and *d*). The more rapid return of

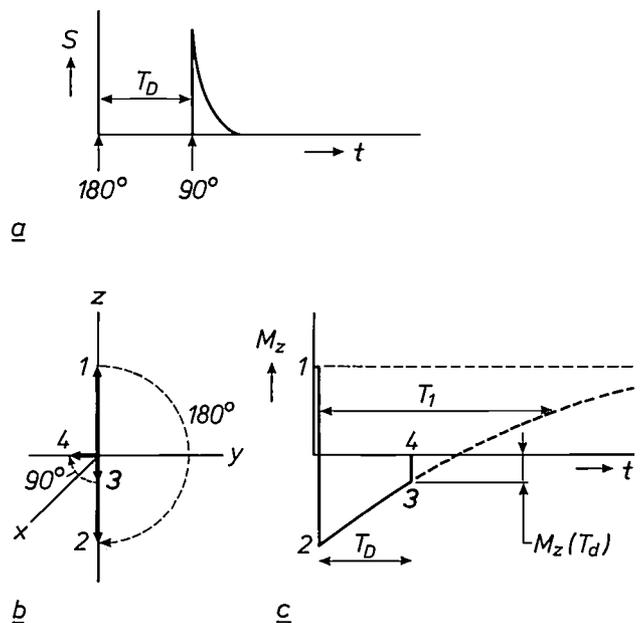


Fig. 4. Determination of the relaxation time T_1 by means of an inversion-recovery sequence. *a*) The amplitude *S* of the r.f. signal as a function of time *t* upon application of a 180° pulse at $t = 0$ (the pulse rotates the magnetization **M** through 180° , see fig. 3*a*) and a 90° pulse at $t = T_D$. The pulses of the r.f. magnetic field are indicated by the thin arrows. *b*) The magnetization vector (thick arrow) at different times: 1 immediately before the application of the 180° pulse, 2 immediately after this, 3 immediately before the application of the 90° pulse, 4 immediately after this. (The coordinate system should be considered as rotating at the frequency of the Larmor precession about the *z*-axis.) *c*) The longitudinal magnetization M_z as a function of time. The times 1, 2, 3 and 4 are indicated. After 4 the transverse magnetization induces an r.f. signal in the detector coil. The amplitude of this signal is a measure of the longitudinal magnetization at time 3. This gives one point — $M_z(T_D)$ — on the exponential curve.

the transverse magnetization to the thermal-equilibrium state is due to slight differences in the frequency of the Larmor precessions, caused by interaction between the spins. The Larmor precessions thus fall out of phase and eventually cancel one another out. The time constant T_2 associated with the relaxation of the transverse magnetization is therefore always smaller than the time constant T_1 associated with the relaxation of the longitudinal magnetization.

The rotating transverse magnetization M_1 induces an r.f. signal that decreases in amplitude in the detector coil. Its frequency, as we have seen, is determined by the magnitude of the constant magnetic field. It was the discovery that spatial information could be added to this signal by applying a small fixed gradient to the constant magnetic field that made proton MR tomography possible. This gradient has the direction of one of the space coordinates. The frequency of the detected signal is then a linear function of that space coordinate, so that the detected signal is the Fourier transform of a notional signal that is a linear function of location. Image-reconstruction methods related to those used in X-ray computer tomography can be used to translate the detected signals into grey-level values for pixels that together yield an (x, y) cross-section or 'slice' of the object.

On reconstruction, the signals obtained as described above give grey-level values that are approximately proportional to the local proton concentration, i.e. the local water content. Some identification of organs or tissues is possible from the images. Intensive investigations of proton magnetic resonance for medical applications have shown that local attribution of values for the relaxation times T_1 and T_2 considerably widens the useful range of application of MR tomography, particularly since in many cases the comparison of T_1 - and T_2 -dependent images makes it possible to distinguish between diseased and healthy tissue. It is not easy to see how information about T_1 can be derived from the detected r.f. signal, since the longitudinal magnetization M_z is constant in direction and does not induce any signal in the detector coil. M_z continues to increase along an exponential curve with a time constant T_1 (see fig. 3d) after the initiation of the Larmor precession.

Fig. 4 illustrates the principle of determining the relaxation time T_1 using the 'inversion-recovery method'. First, the r.f. field is applied for a time necessary to rotate the magnetization M through 180° ($\alpha = 180^\circ$ in fig. 3a). After this '180° pulse' the spin system starts to recover its thermal equilibrium ($2 \rightarrow 3$ in fig. 4b) along the exponential curve shown in fig. 4c. No r.f. signal has yet been induced in the detector coil, however, since the mean transverse magnetiza-

tion is zero. After a time T_D has elapsed, a 90° pulse is applied, which rotates the magnetization through 90° and does give transverse magnetization. The signal induced immediately after application of the 90° pulse is a measure of the longitudinal magnetization $M_z(T_D)$, as shown in fig. 4c. Repeating the experiment for other values of T_D would give more points on the exponential curve, but the procedure we have adopted is slightly different.

Fig. 5 shows the principle of determining the relaxation time T_2 by means of the spin-echo method. After a 90° pulse the longitudinal magnetization is zero and the transverse magnetization has the same magnitude

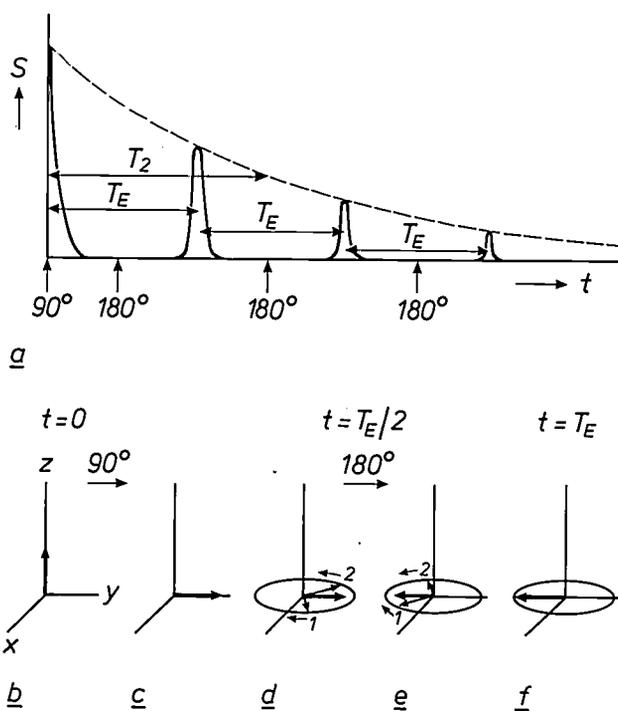


Fig. 5. Determination of the relaxation time T_2 by means of a spin-echo sequence. a) The amplitude S of the r.f. signal as a function of time t , with a 90° pulse applied at $t = 0$ and a 180° pulse at times $t = T_E/2$, $t = 3T_E/2$ and $t = 5T_E/2$. b)-f) The magnetization vector (thick arrow) at different times. (The coordinate system should be considered as rotating about the z -axis at the mean frequency of the Larmor precession.) Because of field inhomogeneities the magnetizations 1 and 2 of different domains start to go out of phase (c \rightarrow d) immediately after the 90° pulse (b \rightarrow c). The 180° pulse (d \rightarrow e) causes the submagnetizations to change sign. The submagnetizations 1 and 2 then converge again and, at a time T_E after the 90° pulse, they come back into phase again (f); they then give an echo signal. Owing to the decrease in the transverse magnetization, the recovery is not complete. The relaxation time T_2 can be determined from the decrease in the peak height of the different echoes.

[2] P. R. Locher, Proton NMR tomography, Philips Tech. Rev. 41, 73-88, 1983/84; L. Kaufman, L. E. Crooks and A. R. Margulis (ed.), Nuclear magnetic resonance imaging in medicine, Igaku-Shoin, New York 1981.

as the longitudinal magnetization immediately beforehand. The magnetization vectors of different domains now precess about the z-axis with an angular frequency approximately equal to ω , from eq. (1). However, owing to local field inhomogeneities, the angular velocities are slightly different, so that vector 1 of a submagnetization is slightly out of phase and rotates faster than the average angular velocity, while vector 2 rotates more slowly (see fig. 5d). To bring about convergence a 180° pulse is then applied, causing the y-component of the vectors of the submagnetizations to change sign and gradually eliminating the dephasing. With a time of $0.5 T_E$ between the 90° and 180° pulses the vectors of the submagnetizations are back in phase again a time T_E after the 90° pulse, causing a signal known as a spin echo to be induced in the detector coil. Periodic repetition of a 180° pulse produces a sequence of spin echoes. The peak height of the successive spin echoes decreases as an exponential curve with a time constant T_2 ; see fig. 5a.

The images of interest in clinical diagnostics are those that contain information about the relaxation time T_1 , the relaxation time T_2 or the proton density ρ . To obtain these images repeated pulse trains are applied, consisting of 180° and 90° pulses in a specific sequence. In our investigations we have used a spin-echo sequence SE consisting of one 90° pulse and four 180° pulses for measuring T_2 , and an inversion-recovery sequence IR consisting of one 180° pulse, one 90° pulse and four 180° pulses for measuring T_1 . This is illustrated in fig. 6, which also shows the variation in amplitude of the detected r.f. signal.

As we noted, after the application of a 90° pulse the longitudinal magnetization M_z becomes zero and the transverse magnetization becomes the same as the longitudinal magnetization just beforehand. Imme-

diately after the 90° pulse the longitudinal magnetization starts to recover again exponentially:

$$M_z = M_0(1 - e^{-t/T_1}),$$

where M_0 is the magnetization in thermal equilibrium. The transverse magnetization similarly decreases exponentially:

$$M_{\perp} = M_{\perp 0} e^{-t/T_2},$$

where $M_{\perp 0}$ is the transverse magnetization immediately after the 90° pulse. From the exponential functions given here for M_z and M_{\perp} we can derive approximate expressions (for $T_R^{SE} \gg T_E$ and $T_R^{IR} \gg T_E$, and apart from a constant factor) for the heights S_{A_1} to S_{A_4} of the four echoes in the signal due to the spin-echo sequence and the height S_B of the first echo in the signal due to the inversion-recovery sequence [3]; see fig. 6:

$$S_{A_i} \propto M_0 \rho (1 - e^{-T_R^{SE}/T_1}) e^{-i T_E/T_2}, \quad i = 1 \text{ to } 4 \quad (2)$$

$$S_B \propto M_0 \rho (1 - 2e^{-T_D/T_1} + e^{-T_R^{IR}/T_1}) e^{-T_E/T_2}, \quad (3)$$

where ρ represents the proton concentration and T_R^{SE} , T_R^{IR} , T_D and T_E are parameters that determine the time dependence of the patterns of pulse sequences. (For simplicity we shall not go into the significance of the second, third and fourth echoes in the signal produced by the inversion-recovery sequence.) The image-reconstruction methods mentioned earlier can be used to derive values for ρ , T_1 and T_2 for each pixel from the measured peak heights S_{A_1} to S_{A_4} and S_B . For this the ratio of S_{A_1} to S_B is calculated for each pixel. From the result, which no longer depends on ρ and T_2 , a value for T_1 can be obtained from a look-up table in the computer. A value for T_2 can then be obtained by using the method of least squares to calculate a line of best fit for $\ln S_{A_i}$ as a function of $i T_E$. Finally, with the

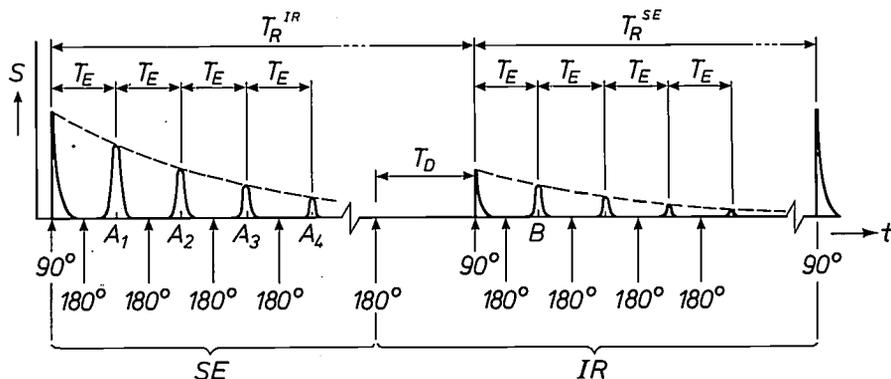


Fig. 6. The pulse sequence periodically applied to determine the proton density and the relaxation times T_1 and T_2 from the detected signals with the aid of equations (2) and (3). The 90° and 180° pulses are indicated by arrows and the amplitude S of the detected r.f. signal is shown diagrammatically as a function of time t . SE spin-echo sequence. IR inversion-recovery sequence. A_1 to A_4 echo signals of the spin-echo sequence. B first echo signal of the inversion-recovery sequence. T_R^{IR} , T_R^{SE} , T_E , T_D pulse parameters.

aid of T_1 and T_2 , S_{A_1} gives a value for the proton density ρ .

Image-synthesis procedure

We have already described how the quantities S_{A_1} , to S_{A_4} from the spin-echo sequences and S_B from the inversion-recovery sequences are obtained for each pixel and how the relaxation times T_1 and T_2 and the proton density ρ are calculated from them. We shall now see how the values of S_{A_1} and S_B are recalculated into other values for another setting of the pulse parameters T_E and T_R^{SE} , or T_E , T_D and T_R^{IR} . The use made of a video look-up table in this process will be illustrated with reference to *fig. 7*.

The results of an MR exposure are stored in the image memory *IM*. A distinction is to be made here between an SE image, built up from the signals of the spin-echo sequences, and an IR image, built up from the signals of the inversion-recovery sequences. For an SE image binary numbers are stored at each memory location for S_{A_1} (4 bits), T_1 and T_2 (3 bits each); for an IR image binary numbers are stored at each location for S_B (4 bits), T_1 and T_2 (3 bits each). In the periodic generation of an image on the monitor *M* each pixel in the image memory is 'translated' by means of the video look-up table in the memory *LTM*. Here the binary value of the combination S_{A_1} , T_1 , T_2 or S_B , T_1 , T_2 serves as the address *ADR* for *LTM*. The address thus has a length of 10 bits and the look-up table contains 2^{10} 8-bit numbers. The result of the translation with the look-up table is an 8-bit sample *VS* of the video signal. A sequence of $256 \times 256 = 2^{16}$ samples contains the information for a complete image. A video signal is constructed from the samples,

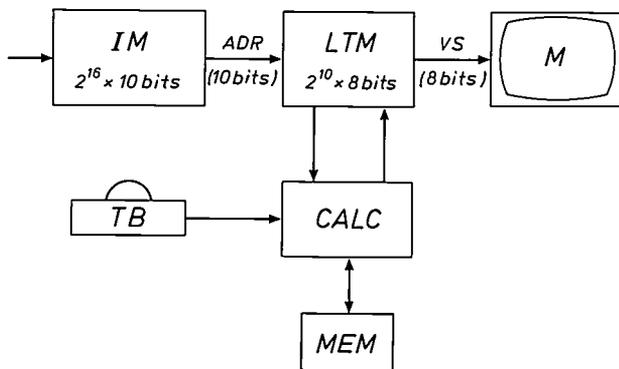


Fig. 7. Block diagram illustrating the image-synthesis procedure. *IM* image memory for 256×256 picture elements (pixels) with 10 bits per memory location. *ADR* address formed by a number in a memory location of *IM*. *LTM* look-up table in the form of a RAM with 2^{10} memory locations with 8 bits per location. *VS* video-signal sample for the monitor *M*. *CALC* part of the MR software that generates new contents for *LTM* when one of the pulse parameters is modified. *MEM* memory for the pulse parameters. *TB* track-ball, an interactive device for altering the pulse parameters.

and this is fed to the monitor. Arrangements are made to ensure that the video signal always contains the complete range of grey levels for the monitor, from white to black. We shall now consider how the contents of the look-up table ($2^{10} \times 8$ bits) are altered.

When the observer changes one of the pulse parameters, a new content for *LTM* is generated by the part of the MR software denoted by *CALC* in *fig. 7*. The relevant pulse parameters are stored in the memory *MEM*. The observer can change a pulse parameter in *MEM* by means of the interactive track-ball device *TB*. As soon as a new pulse-parameter value is entered, *CALC* generates new numbers for the look-up table. For an SE image the observer can alter T_E or T_R^{SE} and for an IR image T_E , T_D or T_R^{IR} . (T_E may be different for SE and IR images.)

The basis for the calculation of new numbers for the look-up table is formed by eqs. (2) and (3). From these equations other equations can be formulated for the new values S_{A_1}'' and S_B'' as a function of the old values S_{A_1}' and S_B' respectively. For an SE image, on changing the pulse parameter T_E' to T_E'' :

$$S_{A_1}'' = S_{A_1}' \frac{e^{-T_E''/T_2}}{e^{-T_E'/T_2}},$$

and on changing the pulse parameter $T_R^{SE'}$ to $T_R^{SE''}$:

$$S_{A_1}'' = S_{A_1}' \frac{1 - e^{-T_R^{SE''}/T_1}}{1 - e^{-T_R^{SE'}/T_1}}.$$

For an IR image on changing the pulse parameter T_E' to T_E'' :

$$S_B'' = S_B' \frac{e^{-T_E''/T_2}}{e^{-T_E'/T_2}},$$

on changing the pulse parameter T_D' to T_D'' :

$$S_B'' = S_B' \frac{1 - 2e^{-T_D''/T_1} + e^{-T_R^{IR}/T_1}}{1 - 2e^{-T_D'/T_1} + e^{-T_R^{IR}/T_1}},$$

and on changing the pulse parameter $T_R^{IR'}$ to $T_R^{IR''}$:

$$S_B'' = S_B' \frac{1 - 2e^{-T_D/T_1} + e^{-T_R^{IR''}/T_1}}{1 - 2e^{-T_D/T_1} + e^{-T_R^{IR'}/T_1}}.$$

The image-synthesis procedure described has been implemented in our MR laboratory configuration. The heart of the configuration is a superconducting magnet with a flux density of 2 T (Tesla). The configu-

[3] F. W. Wehrli, J. R. McFall, G. H. Glover, N. Grigsby, V. Haughton and J. Johanson, The dependence of nuclear magnetic resonance (NMR) image contrast on intrinsic and pulse sequence timing parameters, *Magn. Resonance Imaging* 2, 3-16, 1984; W. H. Perman, S. K. Hilal, H. E. Simon and A. A. Maudsley, Contrast manipulation in NMR imaging, *Magn. Resonance Imaging* 2, 23-32, 1984; M. H. Kuhn, W. Menhardt and I. C. Carlsen, Real-time interactive NMR image synthesis, *IEEE Trans. MI-4*, 160-164, 1985.



Fig. 8. Part of the RAMTEK 9460 graphic display, used in the MR configuration at the Philips Hamburg laboratories. The screen shows an SE cranial image. The track-ball, used in conjunction with the cursor for varying the pulse parameters, can be seen at the lower right.

ration contains a VAX 11-780 computer and a RAMTEK 9460 display system; see *fig. 8*. As mentioned, the observer can use the interactive track-ball device with display system to alter the pulse parameters to any required values. The slightly modified Floyd-Steinberg algorithm that we use for distributing the quantization error of a pixel over adjacent pixels prevents the occurrence of the 'interference patterns' that are characteristic of the 'ordered dither' technique^[41]. *Fig. 9* shows the effect of the quantization by

the Floyd-Steinberg algorithm. *Fig. 9a* is the original SE image of eight bits and *fig. 9b* shows the modified image with only four bits per pixel. Magnified details can be seen in *fig. 9c* and *d*. Although the resolution is not quite so good, since groups of adjacent pixels receive the same grey level, there are no interference patterns. The image in *fig. 9b* compared with that in *fig. 9a* seems to show no change in the grey-level transitions after the four-bit quantization. The eventual image is made with the original number of eight bits per pixel.

Some applications

As we noted, the method of interactive variation of the pulse parameters can be very useful for optimizing specific contrasts and for instructional use. In IR images, for example, the influence of the pulse parameter T_D is very marked. In the inversion-recovery sequence the relaxation time T_1 determines the rate at which the longitudinal magnetization M_z returns to its original value (see *fig. 4c*). It will be evident that the exponential curves for two tissues with different T_1 values may intersect. If T_D is approximately equal to the time interval between the 180° pulse and the point of intersection, a small change in T_D can cause a reversal of the contrast between the two tissues. This is illustrated in *fig. 10*: a change in T_D of only 40 ms causes a reversal of the contrast between the white

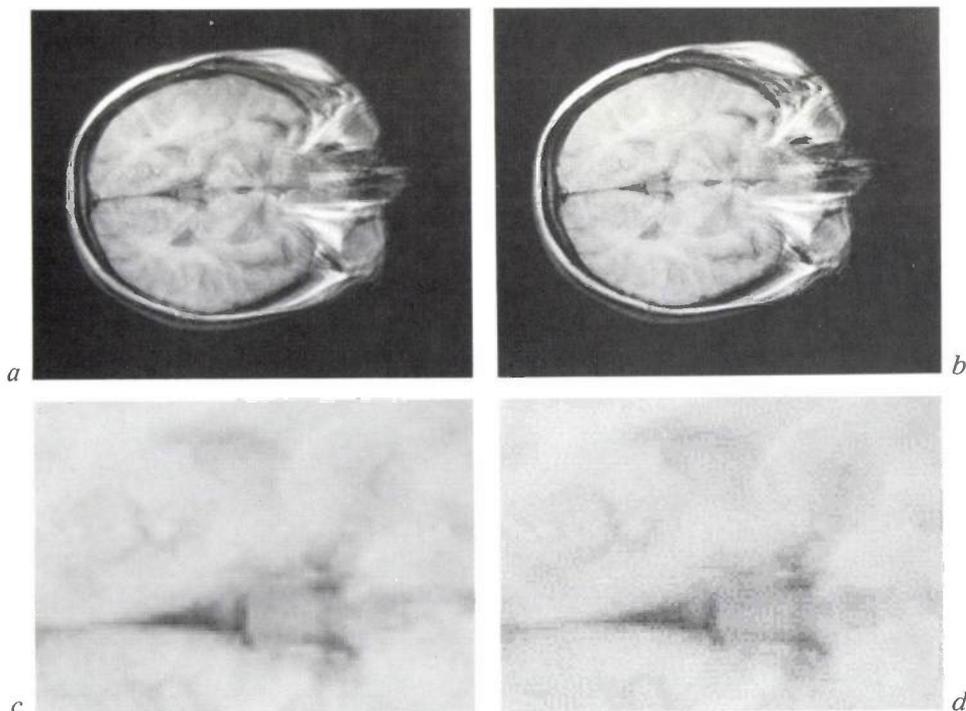


Fig. 9. Quantization effect with the Floyd-Steinberg algorithm^[1]. *a*) The original SE cranial image with 8 bits per pixel. *b*) The modified image with only 4 bits per pixel. *c*) and *d*) Magnified details of *a* and *b*.

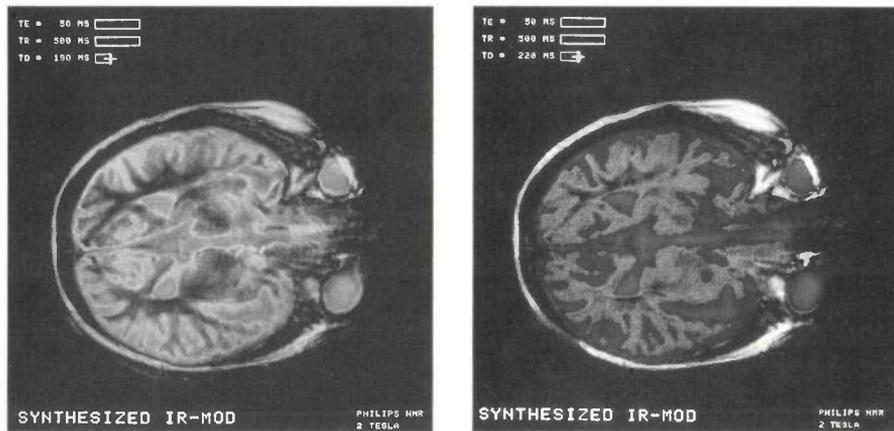


Fig. 10. Contrast reversal in an IR cranial image on changing the pulse parameter T_D . *a*) Exposure with $T_E = 50$ ms, $T_R = 500$ ms and $T_D = 180$ ms. The grey brain matter (cortex) is light grey, the white brain matter is shown almost black. *b*) Exposure with $T_D = 220$ ms and the same values for T_E and T_R . The grey brain matter now appears almost black, the white brain matter is shown dark grey.

and grey matter in an IR image of the main lobe of the brain. As we have said, our method can also be used for maximizing the contrast between two parts of an image indicated by the observer.

The use of colour in the image is a valuable addition. It is possible in principle to characterize biological tissue more or less unambiguously by the values of the proton density ρ and the relaxation times T_1 and T_2 . If tissues are represented by vectors in a three-dimensional (ρ, T_1, T_2) space, the end-points of vectors of identical tissue form clusters of points that do not as a rule overlap. A primary colour can be assigned to each of the quantities ρ , T_1 and T_2 , e.g. blue

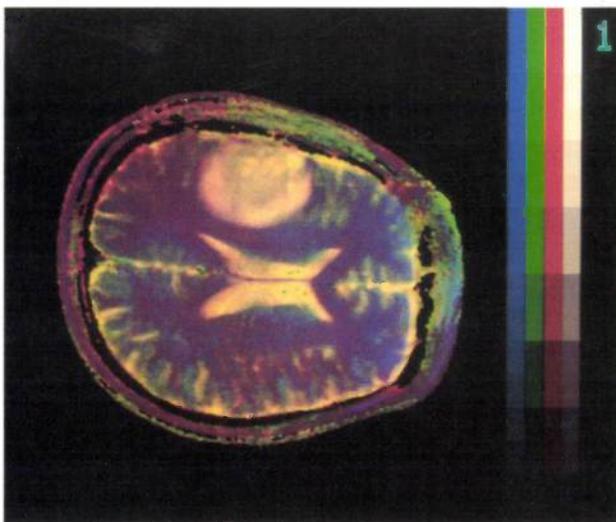


Fig. 11. Cranial image in colour^[6]. The image is a superposition of a blue image, a green image and a red image. In the blue image the brightness is proportional to the proton density ρ , in the green image it is proportional to the relaxation time T_1 and in the red image to the relaxation time T_2 . Each type of tissue has a characteristic additive colour. For example, the malignant tissue of the tumour at the top of the photograph can clearly be distinguished from the benign tissue around it by its orange-yellow colour.

to ρ , green to T_1 and red to T_2 . If the brightness of each of the primary colours is made equal to the values of the relevant quantities, each type of tissue is represented by a characteristic secondary additive colour. *Fig. 11* shows that with this procedure a tumour in the brain can be distinguished from the surrounding healthy tissue^[6].

Our laboratories are currently working on a database that will contain the characteristic quantities ρ , T_1 and T_2 of all biological tissues, with the associated spread in values. When this database has been completed, we hope that it will be possible to use it for identifying tissue in an MR image, perhaps with the aid of other characteristic data such as the age of the patient. The computer can then use look-up tables to establish the type of tissue present in any region of the image indicated by the medical user.

[4] C. N. Judice, J. F. Jarvis and W. H. Ninke, Using ordered dither to display continuous tone pictures on an AC plasma panel, *Proc. SID* 15, 161-169, 1974.

[6] This image was made available to us by Dr Maas of the Eppendorf University Hospital, Hamburg.

Summary. In MR tomography the contrast in images made with the spin-echo method or the inversion-recovery method depends greatly on the time parameters of the r.f. pulses that are applied to obtain signals. The user of MR equipment cannot, however, always indicate the contrast values — the parameters required — beforehand. Once a single MR exposure has been made, the grey levels of points in the MR image can be recalculated to give grey levels for other pulse parameters. The use of a RAM look-up table enables the computer to perform this calculation in real time. The user can thus see the result of a different choice of pulse parameters on the monitor screen immediately. Because small numbers of bits have to be used for the digital values of the proton density and the two relaxation times the quantization noise in the image is increased, but the application of a modified Floyd-Steinberg algorithm prevents the occurrence of 'interference patterns' in the image. The definitive image, which takes rather longer to calculate, is free from quantization noise, however, since it is obtained by using the original numbers of bits. Refinements of the method are the addition of colour to the image and a procedure for maximizing the contrast between areas in the image indicated by the user.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

	Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	<i>A</i>			
	Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	<i>B</i>			
	Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	<i>E</i>			
	Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>			
	Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>			
	Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>			
	Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	<i>R</i>			
J. W. D. Martens	<i>E</i> A magneto-optic investigation of polished manganese-zinc ferrite surfaces		Adv. Ceram. 16	231-235	1986
P. K. Bachmann, W. Hermann, H. Wehr & D. U. Wiechert	<i>A</i> Stress in optical waveguides. 1: Preforms		Appl. Opt. 25	1093-1098	1986
J. J. P. Bruines, R. P. M. van Hal, H. M. J. Boots, W. Sinke* & F. W. Saris* (* FOM, Amsterdam)	<i>E</i> Direct observation of resolidification from the surface upon pulsed-laser melting of amorphous silicon		Appl. Phys. Lett. 48	1252-1254	1986
G. W. 't Hooft & S. Colak	<i>E</i> Influence of the donor depth on the determination of the band discontinuity of isotype heterojunctions by the capacitance-voltage technique		Appl. Phys. Lett. 48	1525-1527	1986
J. Hasker, J. van Esdonk & J. E. Crombeen	<i>E</i> Properties and manufacture of top-layer scandate cathodes		Appl. Surf. Sci. 26	173-195	1986
C. Ronse	<i>B</i> Criteria for approximation of linear and affine functions		Arch. Math. 46	371-384	1986
F. L. van Nes (<i>Inst. Perception Res., Eindhoven</i>)	<i>E</i> Space, colour and typography on visual display terminals		Behav. & Inf. Technol. 5	99-118	1986
G. Gärtner, P. Janiel, H. Rau, H. A. M. van Hal & H. J. P. Nabben	<i>A, E</i> Flow method vapour pressure determination and characterization of tetrakis(trifluoropentanedionato)-thorium(IV) and tetrakis(heptafluorodimethyl-octanedionato)-thorium(IV)		Ber. Bunsenges. Phys. Chem. 90	459-463	1986
B. Chantepie, M. Rocchi, C. Rocher & M. Fichelson	<i>L</i> Packaged 7 mW, 1.2 GHz dynamic 60/61 GaAs prescaler		Electron. Lett. 22	355-356	1986
W. Albers	<i>E</i> Composite materials: nonmechanical properties		Encyclopedia of materials science and engineering, M. B. Bever (ed.), Pergamon, Oxford	763-767	1986
J. M. Robertson	<i>E</i> Crystal growth of ceramics		Encyclopedia of materials science and engineering, M. B. Bever (ed.), Pergamon, Oxford	961-981	1986
J. T. Klomp	<i>E</i> Glass bonding in advanced ceramics		Encyclopedia of materials science and engineering, M. B. Bever (ed.), Pergamon, Oxford	1958-1966	1986
J. T. Klomp	<i>E</i> Joining of ceramic-metal systems: procedures and microstructures		Encyclopedia of materials science and engineering, M. B. Bever (ed.), Pergamon, Oxford	2467-2475	1986

- R. Woltjer, R. Eppenga, J. Mooren, C. E. Timmering & J. P. André *E, L* A new approach to the quantum Hall effect *Europhys. Lett.* 2 149-155 1986
- P. W. J. M. Boumans *E* A century of spectral interferences in atomic emission spectroscopy — can we master them with modern apparatus and approaches? *Fresenius Z. Anal. Chem.* 324 397-425 1986
- G. Söllner *A* Zur Entstaubung der Luft in Wohnräumen *Gesundh. Ing. Haustechn. Bauphys. Umwelttech.* 107 147-155 1986
- J. J. Kelly, J. E. A. M. van den Meerakker & P. H. L. Notten *E* Electrochemistry of photoetching and defect-revealing in III-V materials *Grundlagen von Elektrodenreaktionen*, J. W. Schultze (ed.), VCH, Weinheim 453-464 1986
- D. Meissner*, C. Sinn*, R. Memming* (* *Univ. Hamburg*), P. H. L. Notten & J. J. Kelly *E* On the nature of the inhibition of electron transfer at illuminated p-type semiconductor electrodes *Homogeneous and heterogeneous photocatalysis*, E. Pelizzetti & N. Serpone (eds), Reidel, Dordrecht 317-333 1986
- C. B. Marshall *R* Low-power integrable paging receiver architecture *IEE Proc. F* 133 449-455 1986
- A. D. Sayers *R* Radiometric sky temperature measurements at 35 and 89 GHz *IEE Proc. H* 133 233-237 1986
- L. C. M. G. Pfennings, W. G. J. Mol, J. J. J. Bastiaens & J. M. F. van Dijk *E* Differential split-level CMOS logic for subnanosecond speeds *IEEE J. SC-20* 1050-1055 1985
- M. Rocchi, B. Chantepie, B. Gabilard, G. Haldemann, J. P. Debost*, J. Andrieux* & F. Robert* (* *TRT, Le Plessis-Robinson*) *L* A 1.2-GHz frequency synthesizer using a custom-design divide-by-20/21/22/23/24 GaAs circuit *IEEE J. SC-20* 1194-1199 1985
- T. Ducourant, J.-C. Baelde, M. Binet & C. Rocher *L* 1-GHz, 16-mW, 2-bit analog-to-digital GaAs converter *IEEE J. SC-21* 453-456 1986
- D. H. Evans *R* A millimetre-wave self-oscillating mixer using a GaAs FET harmonic-mode oscillator *IEEE MTT-S Int. Microwave Symp. Digest, Baltimore 1986* 601-604 1986
- P. E. Wierenga (*PDM, Oosterhout*), J. A. van Winsum & J. H. M. van de Linden *E* Roughness and recording properties of particulate tapes: a quantitative study *IEEE Trans. MAG-21* 1383-1385 1985
- S. B. Luitjens, C. P. G. Schrauwen, J. P. C. Bernards & V. Zieren *E* Magnetostatic and recording analysis of RF-sputtered double layer media for perpendicular recording *IEEE Trans. MAG-21* 1438-1440 1985
- A. H. van Ommen & F. H. P. M. Habraken (*Univ. Utrecht*) *E* Bombardment induced redistribution of hydrogen in LPCVD-Si₃N₄ films *Insulating films on semiconductors*, J. J. Simonne & J. Buxo (eds), Elsevier Science, Amsterdam 237-240 1986
- J. P. Charlier & F. Crowet *B* Wave equations in linear viscoelastic materials *J. Acoust. Soc. Am.* 79 895-900 1986
- T. D. Rossing (*Northern Illinois Univ., DeKalb, IL*) & A. J. M. Houtsma *E* Effects of signal envelope on the pitch of short sinusoidal tones *J. Acoust. Soc. Am.* 79 1926-1933 1986
- M. Erman, J. P. André & J. LeBris *L* Spectroscopic ellipsometry study of InP, GaInAs, and GaInAs/InP heterostructures *J. Appl. Phys.* 59 2019-2025 1986
- A. E. T. Kuiper, M. F. C. Willemsen, A. M. L. Theunissen, W. M. van de Wijgert, F. H. P. M. Habraken*, R. H. G. Tijhaar*, W. F. van der Weg* (* *Univ. Utrecht*) & J. T. Chen *E, S* Hydrogenation during thermal nitridation of silicon dioxide *J. Appl. Phys.* 59 2765-2772 1986
- H. A. van Sprang, R. G. Aartsen & A. J. S. M. de Vaan *E* Fast switching in a bistable 270° twist display *J. Appl. Phys.* 59 3087-3090 1986
- J. W. D. Martens *E* Magneto-optical properties of substituted cobalt ferrites: CoFe_{2-x}Me_xO₄ (Me = Rh³⁺, Mn³⁺, Ti⁴⁺ + Co²⁺) *J. Appl. Phys.* 59 3820-3823 1986
- M. H. L. M. van den Broek *E* Aberrations in diverging electron beams caused by an inhomogeneous current density *J. Appl. Phys.* 59 3923-3925 1986

W. Reints Bok	E	Spectroscopic study of the wall corrosion in super high-pressure sodium discharge lamps	J. Appl. Phys. 59	3996-4003	1986
G. J. van der Kolk & M. J. Verkerk	E	Microstructural studies of the growth of aluminum films with water contamination	J. Appl. Phys. 59	4062-4067	1986
J. W. de Haan*, H. M. van den Bogaert, J. J. Ponjeé & L. J. M. van de Ven* (* Eindhoven Univ. Technol.)	E	Characterization of modified silica powders by Fourier transform infrared spectroscopy and cross-polarization magic angle spinning NMR	J. Colloid & Interface Sci. 110	591-600	1986
P. J. A. Thijs, W. Nijman & R. Metselaar (Eindhoven Univ. Technol.)	E	LPE of InP and InGaAsP on InP substrates; a verification of the diffusion limited growth model	J. Cryst. Growth 74	625-634	1986
J. C. Brice, C. Papper & C. L. Jones (Mullard, Southampton)	R	The phase diagram of the pseudo-binary system CdTe-HgTe and the segregation of CdTe	J. Cryst. Growth 75	395-399	1986
J. van de Ven*, J. L. Weyher* (* Catholic Univ. Nijmegen), J. E. A. M. van den Meerakker & J. J. Kelly	E	Kinetics and morphology of GaAs etching in aqueous CrO ₃ -HF solutions	J. Electrochem. Soc. 133	799-806	1986
J. P. André, E. P. Menu, M. Erman, M. H. Meynadier & T. Ngo	L	Ga _{1-x} In _x As-InP abrupt heterostructures grown by MOVPE at atmospheric pressure	J. Electron. Mater. 15	71-74	1986
H. K. Kuiken	E	On the influence of longitudinal diffusion in time-dependent convective-diffusive systems	J. Fluid Mech. 165	147-162	1986
R. Grössinger*, R. Krewenka*, R. Eibler*, H. R. Kirchmayr* (* Tech. Univ. Vienna), J. Ormerod (Mullard Magn. Components, Southport) & K. H. J. Buschow	E	The hard magnetic properties of sintered Nd-Fe-B permanent magnets	J. Less-Common Met. 118	167-172	1986
K. H. J. Buschow	E	Invar effect in R ₂ Fe ₁₄ B compounds (R ≡ La, Ce, Nd, Sm, Gd, Er)	J. Less-Common Met. 118	349-353	1986
R. A. de Groot (Univ. Nijmegen) & K. H. J. Buschow	E	Recent developments in half-metallic magnetism	J. Magn. & Magn. Mater. 54-57	1377-1380	1986
J. W. D. Martens & W. F. Godlieb	E	The surface properties of mechanically polished (100) manganese zinc ferrites	J. Magn. & Magn. Mater. 59	325-332	1986
J. W. M. Jacobs & J. F. C. M. Verhoeven	E	Specimen preparation technique for high resolution transmission electron microscopy studies on model supported metal catalysts	J. Microsc. 143	103-116	1986
J. Droho*, G. Bruning & W. Kepka* (* Advance Transformer, Chicago, IL)	N	Feasibility of a high frequency power converter for the microwave magnetron	J. Microwave Power 21	91-92	1986
C. A. M. Mulder, J. G. van Lierop & G. Frens	E	Densification of SiO ₂ -xerogels to glass by Ostwald ripening	J. Non-Cryst. Solids 82	92-96	1986
D. M. Krol & E. M. Rabinovich (AT&T Bell Labs, Murray Hill, NJ)	E	Raman study of fluorine-doped silica gels	J. Non-Cryst. Solids 82	143-147	1986
C. A. M. Mulder, G. van Leeuwen-Stienstra, J. G. van Lierop & J. P. Woerdman (Univ. Leiden)	E	Chain-like structure of ultra-low density SiO ₂ sol-gel glass observed by TEM	J. Non-Cryst. Solids 82	148-153	1986
J. G. van Lierop, A. Huizing, W. C. P. M. Meerman & C. A. M. Mulder	E	Preparation of dried monolithic SiO ₂ gel bodies by an autoclave process	J. Non-Cryst. Solids 82	265-270	1986
P. Dolizy & F. Grolière	L	Dissociation energies of alkali antimonides as thin layers	J. Phys. D 19	687-698	1986
J. C. M. Henning, J. P. M. Ansems & P. J. Roksnoer	E	Spectroscopic determination of L ₆ conduction band minima in Al _x Ga _{1-x} As	J. Phys. C 19	L335-L338	1986
H. Schink*, G. Packeiser* (* Siemens, München), J. Maluenda & G. M. Martin	L	GaAs substrate material assessment using a high lateral resolution MESFET test pattern	Jap. J. Appl. Phys. 25	L369-L372	1986
C. C. Paige (McGill Univ., Montreal) & P. van Dooren	B	On the quadratic convergence of Kogbetliantz's algorithm for computing the singular value decomposition	Linear Algebra & Appl. 77	301-313	1986
H. J. W. van Houtum & I. J. M. M. Raaijmakers	E	First phase nucleation and growth of titanium disilicide with an emphasis on the influence of oxygen	Mater. Res. Soc. Symp. Proc. 54	37-42	1986

A. G. Dirks, T. Tien, J. M. Towner <i>E, S</i>	Homogeneous Al-Ti and Al-Ti-Si thin alloy films	Mater. Res. Soc. Symp. Proc. 54	181-186	1986
A. H. van Ommen, M. F. C. Willemssen, P. R. Boudewijn & A. H. Reader <i>E</i>	Bombardment-induced mixing of Mo films on Si	Mater. Res. Soc. Symp. Proc. 54	221-226	1986
H. A. van Sprang & G. Vertogen (<i>Univ. Nijmegen</i>) <i>E</i>	Vloeibare kristallen	Natuur & Tech. 54	384-397	1986
P. America <i>E</i>	Design issues in a parallel object-oriented language	Parallel computing 85, M. Feilmeier <i>et al.</i> (eds), Elsevier Science, Amsterdam	325-330	1986
E. A. M. Odijk <i>E</i>	Overview of the Philips object-oriented parallel computer	Parallel computing 85, M. Feilmeier <i>et al.</i> (eds), Elsevier Science, Amsterdam	527-532	1986
P. A. Devijver <i>B</i>	On the editing rate of the MULTIEDIT algorithm	Pattern Recognition Lett. 4	9-12	1986
A. Browne <i>R</i>	Vision and the robot	Philips J. Res. 41	232-246	1986
G. Clark & R. F. Milsom <i>R</i>	Analysis of surface acoustic wave filter by expansion guided modes	Philips J. Res. 41	247-267	1986
P. F. Fewster <i>R</i>	X-ray diffraction from multiple quantum well structures	Philips J. Res. 41	268-289	1986
A. J. Fox <i>R</i>	Calculation for the diffraction efficiency of acousto-optic modulators using an integral equation approach	Philips J. Res. 41	290-312	1986
C. T. Foxon & J. J. Harris <i>R</i>	The growth of high purity III-V structures by molecular beam epitaxy	Philips J. Res. 41	313-324	1986
A. W. Woodhead, D. Washington, D. L. Lamport, A. G. Knapp, J. R. Mansell & K. G. Freeman <i>R</i>	The channel multiplier cathode-ray tube	Philips J. Res. 41	325-342	1986
G. G. P. van Gorkom & A. M. E. Hoeberechts <i>E</i>	An efficient silicon cold cathode for high current densities, II. Comparison with theory and discussion	Philips J. Res. 41	343-384	1986
J. M. M. Pasmans & J. Haisma <i>E</i>	An environmental-nonpolluting antireflective coating for ZnSe optics at CO ₂ -laser wavelengths	Philips J. Res. 41	385-390	1986
W. J. van Gils <i>E</i>	An error-control coding system for storage of 16-bit words in memory arrays composed of three 9-bit wide units	Philips J. Res. 41	391-399	1986
D. B. de Mooij & K. H. J. Buschow <i>E</i>	A novel ternary Nd-Fe-B compound	Philips J. Res. 41	400-409	1986
K. A. Schouhamer Immink <i>E</i>	Coding methods for high-density optical recording	Philips J. Res. 41	410-430	1986
A. Osseyran (<i>Philips Elcoma Div., Eindhoven</i>) <i>R</i>	Computer aided design of magnetic deflection systems	Philips J. Res. 41, Suppl. No. 1	151 pp.	1986
E. Kauer & E. Schnedler <i>A</i>	Möglichkeiten und Grenzen der Lichterzeugung	Phys. Bl. 42	128-133	1986
A. R. Hepburn*, J. M. Marshall*, C. Main* (* <i>Dundee College Technol.</i>), M. J. Powell & C. van Berkel <i>R</i>	Metastable defects in amorphous-silicon thin-film transistors	Phys. Rev. Lett. 56	2215-2218	1986
J. A. Geurst <i>E</i>	Variational principles and two-fluid hydrodynamics of bubbly liquid/gas mixtures	Physica 135A	455-486	1986
A. M. van der Kraan (<i>Interuniv. Reactor Inst., Delft</i>) & K. H. J. Buschow <i>E</i>	The ⁵⁷ Fe Mössbauer isomer shift in intermetallic compounds of iron	Physica 138B	55-62	1986
E. Delhaye, C. Rocher, J.-M. Gibereau & M. Rocchi <i>L</i>	A 2.5 ns, 40 mW, 4 × 4 GaAs multiplier in 2's complement mode	Proc. ESSCIRC'86, Delft 1986	56-58	1986
G. Gärtner, P. Janiel & H. Lydtin <i>A</i>	Chemical vapour deposition of tungsten/thorium structures	Proc. Eur. Conf. on Chemical vapour deposition, Uppsala 1985	52-59	1985
P. Gamand, R. Leblanc, M. Levent-Villegas & P. Rabinzohn <i>L</i>	New design tools for high fabrication yield low cost monolithic GaAs circuits using passive and active components statistical models	Proc. Eur. Microwave Conf., Dublin 1986	607-612	1986

V. Pauker, P. Dauriche, A. Giakoumis (RTC Limeil, Limeil-Brévannes) & A. Kazeminejad (Paris Univ.)	L	Normally-off MESFET analogue circuits	Proc. Eur. Microwave Conf., Dublin 1986	631-635	1986
H. Sari, L. Desperben & S. Moridi	L	Optimum timing recovery for digital equalizers	Proc. GLOBECOM'85, New Orleans 1985	1460-1465	1985
M. A. Shaik & P. Datsaris (Univ. Rhode Island, Kingston, RI)	N	A workspace optimization approach to manipulator linkage design	Proc. IEEE Int. Conf. on Robotics & Automation, San Francisco, CA, 1986	75-81	1986
N. M. Marinovic & W. A. Smith	N	Application of joint time-frequency distributions to ultrasonic transducers	Proc. IEEE Int. Symp. on Circuits & systems, San Jose, CA, 1986	50-54	1986
M. P. A. Vieggers, B. H. Koek & A. H. van Ommen	E	Precipitation of coesite in oxygen implanted silicon	Proc. Int. Cong. on Electron Microscopy, Kyoto 1986	121-122	1986
M. P. A. Vieggers, B. H. Koek, J. J. P. Bruines, R. P. M. van Hal & H. M. J. Boots	E	Subsurface extension of explosive crystallization	Proc. Int. Cong. on Electron Microscopy, Kyoto 1986	1521-1522	1986
J. P. André, M. Wolny & M. Rocchi	L	MO VPE versus MBE for III/V heterostructure devices and ICs	Proc. Int. Symp. GaAs & Related Compounds, Karuizawa 1985	379-384	1986
J. B. Hughes, N. C. Bird & R. S. Soin	R	A novel digitally self-tuned continuous-time filter technique	Proc. ISCAS '86, San Jose 1986	1177-1180	1986
M. G. Collet, J. G. C. Bakker, L. J. M. Esser, H. L. Peek, M. J. H. van de Steeg, A. J. P. Theuwissen & C. H. L. Weijtens	E	High density frame transfer image sensors with vertical anti-blooming	Proc. SPIE 570	27-34	1985
R. H. Coursant, P. Eyraud* & L. Eyraud* (* INSA, Villeurbanne)	L	Influence of material texture and ionic substitutions in piezoceramics on the resonance frequency dispersion diagrams of narrow-strip elements; application to multielement ultrasonic transducers	Proc. Ultrasonics Symp., San Francisco 1985	614-619	1985
J. P. Arragon	L	Utilisation des mémoires d'images dans les récepteurs grand public	Rev. Radiodif.-Télév. No. 91	1-13	1986
B. Bölger & P. K. Larsen	E	Video system for quantitative measurements of RHEED patterns	Rev. Sci. Instrum. 57	1363-1367	1986
M. Mijiyawa	L	Hybrid motion estimation in successive television pictures	Signal processing III, I. T. Young <i>et al.</i> (eds), Elsevier Science, Amsterdam	837-840	1986
J. J. H. van den Biesen	E	A simple regional analysis of transit times in bipolar transistors	Solid-State Electron. 29	529-534	1986
P. Friedel, S. Gourrier, J. B. Theeten, D. Arnoult, M. Taillepie & M. Erman	L	Study of the interaction of plasmas with III-V semiconductor surfaces, application to passivation	Surf. Sci. 168	635-644	1986
C. Ronse	B	A topological characterization of thinning	Theor. Comput. Sci. 43	31-41	1986
M. J. Verkerk & W. A. M. C. Brankaert	E	Effects of water on the growth of aluminium films deposited by vacuum evaporation	Thin Solid Films 139	77-88	1986
F. Pasqualini	L	Influence des effets piezoelectriques sur la protection passivation des TEC's AsGa n-off	Vide/Couches Minces 41	197-199	1986
J. P. Chané, B. G. Martin, J. N. Patillon & J. L. Gentner	L	Dépôts diélectriques appliqués à la passivation des photodiodes PIN InGaAs/InP	Vide/Couches Minces 41	203-204	1986
M. Taillepie & S. Gourrier	L	Mise en évidence du rôle de l'oxyde natif sur la stabilité des MISFETs en InGaAs	Vide/Couches Minces 41	209-210	1986
J. L. Gentner, C. Guedon & J. P. Chané	L	Préparation d'hétérostructures (In,Ga)As/InP par épitaxie en phase vapeur aux chlorures	Vide/Couches Minces 41	253-254	1986

Layered semiconductor structures

In research on materials for the electronics industry we have recently seen a growing interest in layered epitaxial structures. If the layers are very thin, these structures often have very different physical properties from the individual materials, and different kinds of applications. Typical examples are the combinations of the III-V semiconductors GaAs and $Al_xGa_{1-x}As$ in field-effect transistors with a high electron mobility and in semiconductor lasers with a short emission wavelength.

Although these new structures do not seem to be compatible with the laws of ordinary chemical thermodynamics, they can be fabricated if careful attention is paid to the facilities for preparation. The layers or films are deposited on a substrate by evaporation in ultra-high vacuum (molecular beam epitaxy, or MBE) or by deposition from a gas mixture containing metal-organic compounds (metal-organic vapour-phase epitaxy, or MO-VPE). When these technologies are combined with sophisticated methods of characterization and analysis, monolayers of atoms or molecules can be applied, a layer at a time, and abrupt transitions can be introduced into the composition.

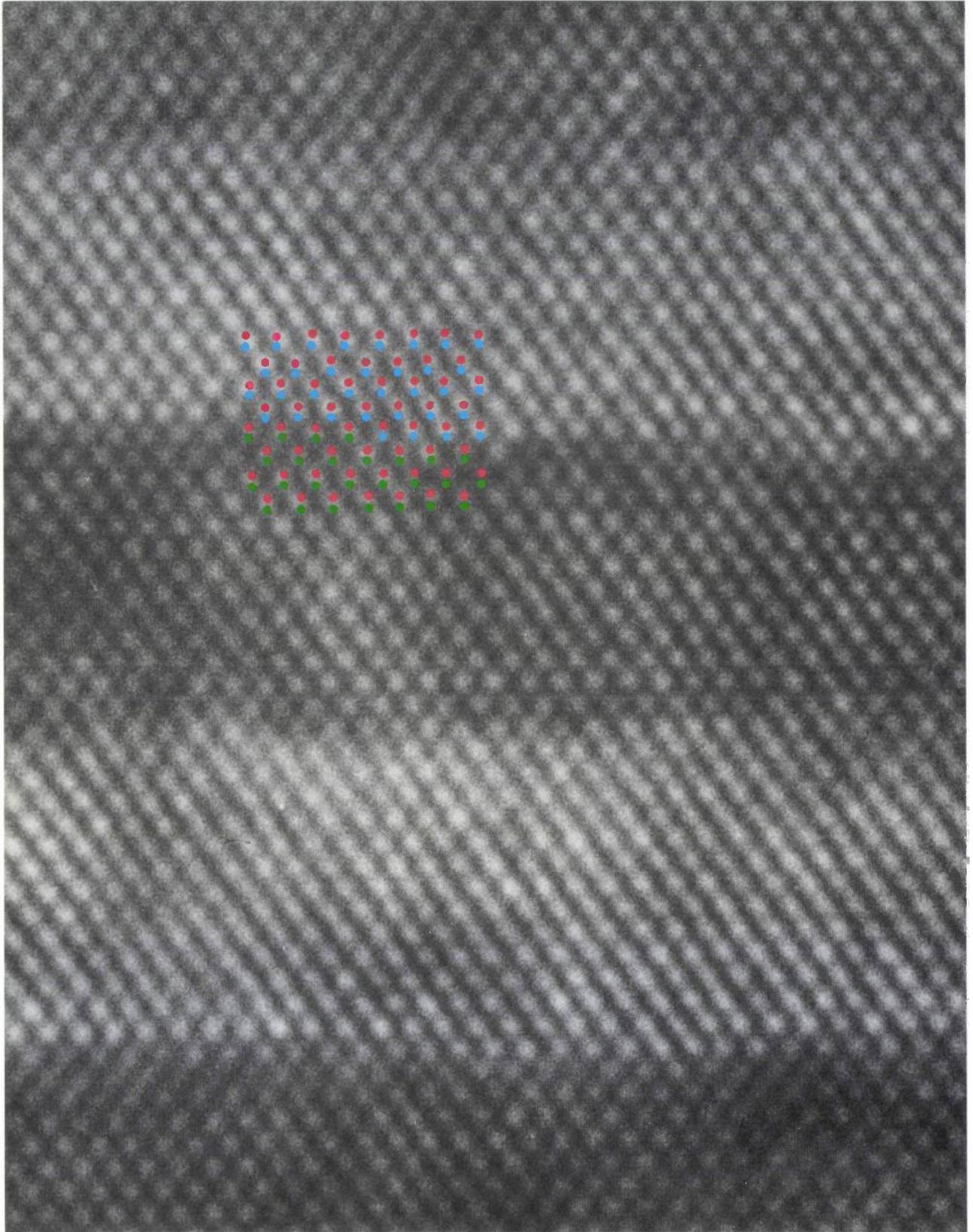
The main centre for the Philips activities in the MBE of III-V semiconductors is the Research Laboratories at Redhill, in Britain. The leading position of these Laboratories can be judged by the record established there for electron mobility in semiconductor structures. Work on the MBE of II-VI semiconductors has recently been started at the Research Laboratories at Briarcliff, in the U.S.A., while the MBE of silicon and combinations of metals is being investigated at the Research Laboratories in Eindhoven.

The MO-VPE technology has perhaps won greater industrial acceptance than MBE. Contributions from our Laboratories at Limeil-Brévannes (France) and Eindhoven have helped to bring this technology to a high degree of perfection in recent years.

Combinations of very different materials, such as silicon on GaAs (or GaP), GaAs on silicon, semiconductors on oxide substrates or combinations of metals and semiconductors also receive attention. There are good prospects for the discovery of new physical effects and unusual properties of materials that can be used in new applications.

The time was clearly ripe for a special issue of Philips Technical Review on these structures. We now present a general survey of the subject and four articles that describe some of the activities of Philips Research on layered semiconductor structures. However, our special issue can really give no more than a glimpse of a field still forging strongly ahead.

A. R. Miedema



The crystal perfection that can be achieved in layered semiconductor structures can be seen from images made by transmission electron microscopy of the crystal lattice. This image shows alternate 3.5-nm layers of GaAs (dark) and AlAs (light); the arrangement of the atoms is indicated by blue for Al, green for Ga and red for As. Changes in composition are completed within a thickness of about one monolayer.

Research on layered semiconductor structures

J. Wolter

Introduction

The special properties of semiconductor single crystals have regularly led to the introduction of devices based on new principles of application. This has resulted for example in silicon integrated circuits with high packing densities and complex functions, and in optoelectronic devices based on III-V compounds such as gallium arsenide (GaAs). Materials scientists have contributed extensively to these developments by making it possible to grow high-purity crystals free from lattice defects.

Most semiconductor devices manufactured today contain only one basic material. We now see the emergence of a completely new class of devices with structures consisting of 'stacked' thin layers of dissimilar semiconductors. In these 'heterostructures' the electric potential in the direction perpendicular to the layers can be modified to produce interesting physical effects, sometimes with novel applications.

During the last ten years or so, research on these structures has expanded rapidly. Their physical background is now much more clearly understood and various discoveries have been evaluated for practical application. This has only been possible through the interplay of ideas, potential device applications and advanced growth techniques, with multi-disciplinary contributions from large numbers of workers in industrial and university laboratories. Efforts in this field are expected to intensify in the years ahead.

Among the structures investigated, probably the most fascinating are the 'superlattices'. These are built up from alternate ultra-thin layers, typically 1 nm thick, of two dissimilar semiconductor materials. Then there is the 'quantum-well structure', consisting of a semiconductor layer less than 30 nm thick between two layers of another semiconductor with a larger band gap. This also has peculiarly interesting

properties. And even a simple heterojunction between two dissimilar semiconductors offers interesting potential applications.

This article gives a brief review of the research on these multilayer structures. It includes a general discussion of some of the physical aspects and practical implications. An indication is also given of the methods of preparation and the semiconductor materials. Finally there is a summary of the Philips activities in this field.

Superlattices

The basic idea of the formation of superlattices was put forward in about 1970 by L. Esaki and R. Tsu^[1]. They considered an array of alternating ultra-thin layers of two dissimilar semiconductors with different band gaps; see *fig. 1*. The alternation gives a periodic modification of the electric potential perpendicular to the interfaces. Extra potential wells are therefore created for the electrons, in addition to the 'ordinary' potential wells around each atom in the crystal lattice. When the periodicity of the superlattice becomes less than about 10 nm, there is a marked change in the energy-level diagram for the electrons: the valence and conduction bands are split into 'minibands' and new forbidden zones ('minigaps') are created. A similar modification applies to the energy of the quantized lattice vibrations (phonons) as a function of their wave number; this has already been demonstrated experimentally by Raman spectroscopy^[2].

A curious effect can occur when an electric field is applied to a superlattice^[1]. The presence of minigaps gives rise to 'Bloch oscillations' and to a negative dif-

Prof. Dr J. Wolter, Professor in Semiconductor Physics at Eindhoven University of Technology, was formerly with Philips Research Laboratories, Eindhoven.

^[1] L. Esaki and R. Tsu, Superlattice and negative differential conductivity in semiconductors, *IBM J. Res. & Dev.* **14**, 61-65, 1970.

^[2] J. L. Merz, A. S. Barker, Jr., and A. C. Gossard, Raman scattering and zone-folding effects for alternating monolayers of GaAs-AlAs, *Appl. Phys. Lett.* **31**, 117-119, 1977.

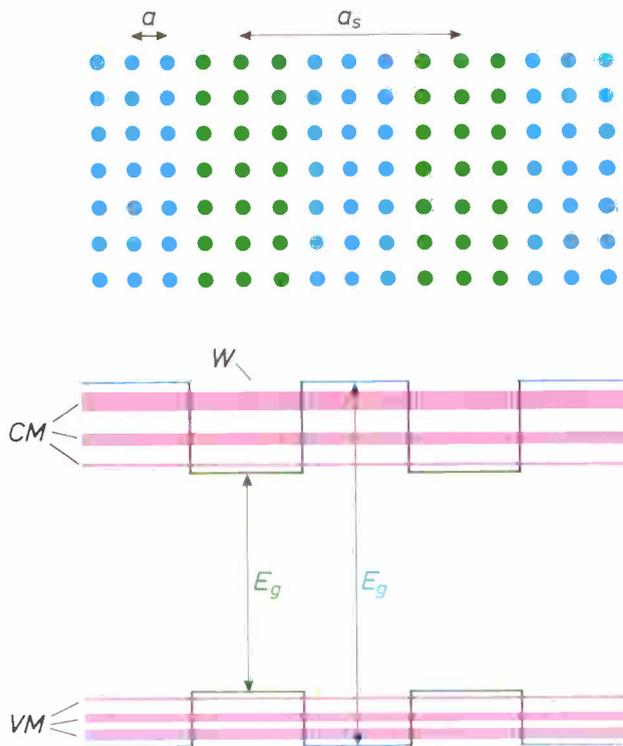


Fig. 1. Atomic arrangement and energy-band diagram of a superlattice consisting of alternate ultra-thin layers of two semiconductors with different band gaps E_g . In addition to the normal lattice period a , there is also a superlattice period a_s . At the interfaces the band-gap difference creates potential wells W , which split the valence and conduction bands into valence minibands VM and conduction minibands CM .

ferential conductivity, as illustrated in *fig. 2*. The energy bands are tilted and the resulting slope is proportional to the voltage. Electrons are driven towards the upper edge of the conduction band, but in a conventional semiconductor they never arrive there, since the distance they have to travel is much farther than the mean distance between two phonon emissions. In a superlattice, however, the minibands may be so narrow that the electrons do have a good chance of reaching the upper edge. Once they arrive there, they turn back towards the bottom edge because they cannot pass the forbidden zone. The electrons therefore oscillate between the two edges many times before emitting a phonon. The average electron position is shifted by a phonon emission, and the shift decreases as the angle of tilt becomes larger. A higher voltage therefore gives a lower current.

Another type of superlattice has also been proposed and investigated^[3]. This consists of layers of the same semiconductor with alternate p-type and n-type doping. The donor atoms in the n-type layers deliver electrons and the acceptor atoms in the p-type layers deliver holes. The resulting charge distribution creates a new set of potential wells, again producing energy-

band structures with minibands and minigaps. Important features of these superlattices are the nearly perfect lattice matching and the spatial separation of electrons and holes, which can increase the recombination lifetimes. The effective band gap can also be changed by applying an external voltage.

Superlattices have unusual physical properties, which can be varied in a controlled way by modifying the band structure via the composition, the doping or

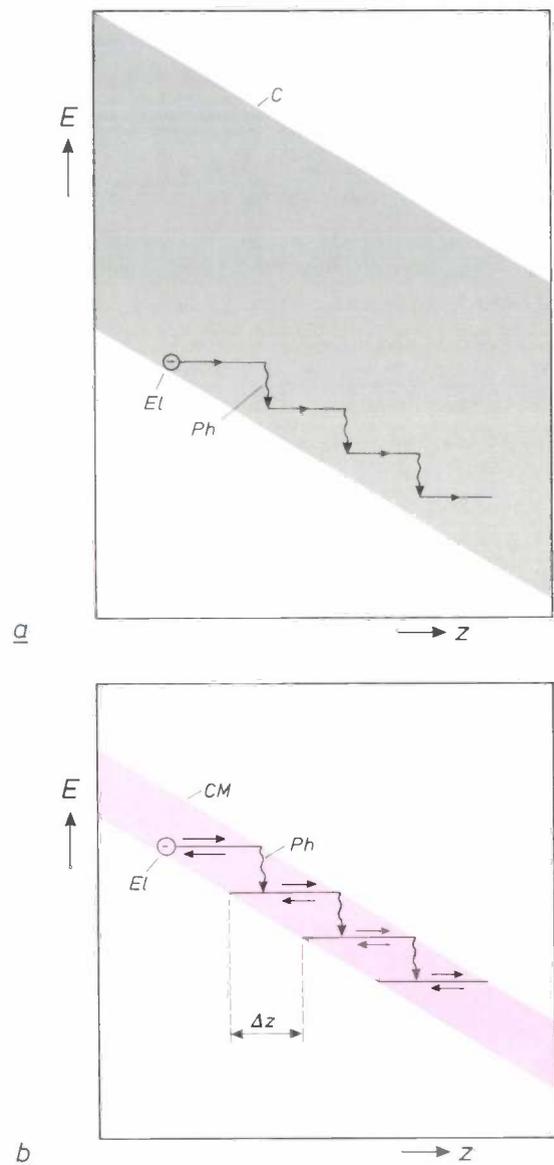


Fig. 2. The occurrence of Bloch oscillations and a negative differential conductivity in a superlattice in an electric field^[3]. The energy diagrams (energy E as a function of position z) are shown for the conduction electrons in a conventional semiconductor (a) and in a superlattice (b). In a conventional semiconductor the emission of phonons Ph prevents electrons El from reaching the upper edge of the conduction band C . In a superlattice the miniband CM is so narrow that electrons do have a chance of arriving at the upper edge. They are reflected there, and may oscillate repeatedly between the band edges before emitting a phonon. The shift Δz in the mean position of the electrons at a phonon emission decreases as the slope of the tilted band increases. This means that the current through the lattice decreases with increasing voltage.

the layer thickness. A practical example of this 'band-gap engineering' is the 'staircase' modification of the energy-band structure for an avalanche photodiode^[4]. Here the band gap is increased in steps in such a way that the band-edge discontinuities provide the entire ionization energy; see *fig. 3*. The multiplication of electrons then only occurs at the well-defined steps. The statistical fluctuations in the internal gain will therefore be much smaller than in conventional avalanche photodiodes, resulting in a better noise performance.

Another example is a superlattice of two semiconductors with an indirect band gap. In the separate materials the minimum of the conduction band does not correspond to the same wave number as the maximum of the valence band, and therefore the coupling to a phonon is required for an electron-hole recombination. Consequently the recombination probability is much smaller (e.g. by a factor of 1000) than in a direct-gap semiconductor. In a superlattice, however, the band structure can be tailored in such a way that a material with a direct band gap is obtained^[5].

This principle might be applied to gallium phosphide (GaP) and aluminium phosphide (AlP), which have nearly the same lattice constant and have band-gap transitions in the green and the blue, respectively. If these two indirect-gap semiconductors could be used to form a superlattice with a direct band gap, a blue-emitting laser material could perhaps be obtained. It might even be possible to obtain light emission from a superlattice of silicon and germanium. This would be

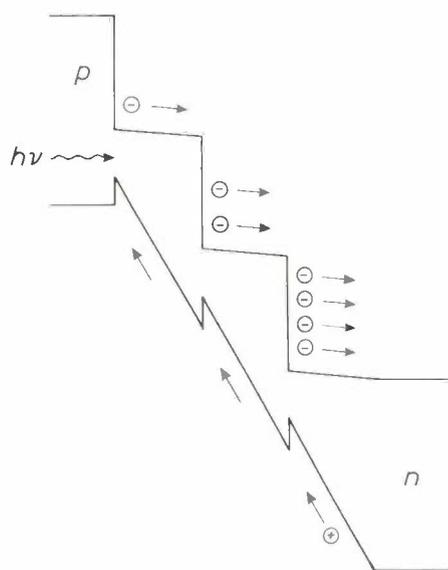


Fig. 3. Energy-band diagram of a 'staircase' superlattice for an avalanche photodiode^[4]. The electrons generated by the incident photons can only ionize the atoms in the crystal at the steps in the conduction band. The creation of new electrons therefore occurs at well-defined positions in the lattice. This reduces the statistical fluctuations in the electron multiplication.

a revolutionary development in optoelectronics, and the idea of growing these structures and studying them experimentally and theoretically presents a challenge.

Quantum-well structures

A semiconductor layer sandwiched between two other semiconductor layers with a larger band gap is called a quantum well when it is so thin (< 30 nm) that the states for electrons and holes moving perpendicular to the interfaces are quantized. The confinement of charge carriers then gives rise to discrete energy levels; see *fig. 4*. The energies corresponding to movements parallel to the interfaces have to be added to these levels. The difference between the lowest conduction level and the highest valence level is greater than in the bulk material. Luminescence due to electron-hole recombinations will therefore occur at shorter wavelengths. The changes in the densities of states associated with quantization^[6] also help to increase the probability of recombination.

In recent years quantum wells consisting of III-V semiconductors have been studied in many research laboratories. Most attention is given to structures of

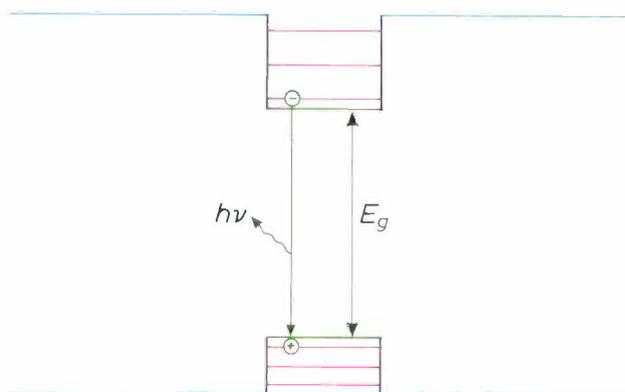


Fig. 4. Band diagram for a quantum well consisting of a semiconductor layer (green) sandwiched between two other layers of semiconductors with a larger band gap (blue). Only the discrete levels (red) are shown that correspond to states in which the electrons and holes only move at right angles to the interfaces. The difference in energy between the lowest conduction level and the highest valence level is larger than the band gap E_g in the bulk material, so that the luminescence due to electron-hole recombinations occurs at shorter wavelengths.

^[3] G. H. Döhler, Electron states in crystals with 'nipi-superstructure', *Phys. Stat. Sol. B* **52**, 79-92, 1972. This author also gives a clear description of superlattices in: *Sci. Am.* **249**, No. 5 (November), 118-126, 1983.

^[4] F. Capasso, W. T. Tsang and G. F. Williams, Staircase solid-state photomultipliers and avalanche photodiodes with enhanced ionization rates ratio, *IEEE Trans. ED-30*, 381-390, 1983.

^[5] U. Gnuzmann and K. Clausecker, Theory of direct optical transitions in an optical indirect semiconductor with a superlattice structure, *Appl. Phys.* **3**, 9-14, 1974.

^[6] R. Dingle, Confined carrier quantum states in ultrathin semiconductor heterostructures, *Festkörperprobleme* **15**, 21-48, 1975.

GaAs sandwiched between $\text{Al}_x\text{Ga}_{1-x}\text{As}$, two materials that have nearly the same lattice constant. Similar structures, but with a thicker active layer with no quantization effects, are applied in semiconductor lasers^[7], e.g. for Compact Disc players^[8] and digital optical recording^[9]. With a quantum-well version of this laser^[10] the wavelength of the laser emission can be reduced by making the well thinner. Similar wavelength tuning can of course be applied in other devices such as light-emitting diodes and photodetectors.

In addition to the wavelength tuning, quantum-well devices have other attractive features, such as a high luminescence efficiency and a lower and less temperature-dependent threshold current for laser operation^[11]. They may also offer stability and long life. Recognition of these features has already led to an industrial application.

High electron mobilities

It is well known that III-V semiconductors have good transport properties. Electrons in GaAs, for instance, have a smaller effective mass and therefore a higher mobility than in silicon. The electron mobility is generally limited, however, by a number of scattering processes, especially at built-in impurities. Collisions between electrons and their donor atoms are often unavoidable because they are present in the same medium. These processes, which are more important at low temperatures (< 77 K) than the interactions with phonons, were until recently the obstacles in the race to achieve higher electron mobilities by improving methods for growing crystals.

The introduction of modulation doping in heterostructures containing GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ provided a breakthrough^[12]. The principle, based on ideas put forward earlier by Esaki and Tsu^[11], is simple and ingenious: only the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ is doped, e.g. by the incorporation of donor atoms (n-type doping). Because of the band-gap discontinuity the electrons move from the donor atoms to the undoped GaAs. This results in a spatial separation between the donor atoms and the electrons, thus increasing the mobility by several orders of magnitude^[13]. A further increase can be obtained by introducing a thin spacer layer of undoped $\text{Al}_x\text{Ga}_{1-x}\text{As}$, which increases the distance between the donor atoms and the electrons still further. The increase in mobility is particularly large at low temperatures. At the Philips laboratories at Redhill in England a value of $3.1 \times 10^6 \text{ cm}^2\text{V}^{-1}\text{s}^{-1}$ has recently been measured^[14].

This principle can be applied in super-fast transistors and integrated circuits. Most of them are based on a structure of the MOSFET type (metal-oxide-

semiconductor field-effect transistor). This consists of a p-type silicon layer with connections to source and drain, an insulating SiO_2 layer and a metal layer connected to the gate; see fig. 5a. When a negative voltage is applied to the gate the conduction electrons are pushed away from the interface between semiconductor and insulator and no current can flow between source and drain. When a positive voltage is applied to the gate the conduction electrons are attracted to the interface, forming a very thin inversion layer with a 'two-dimensional electron gas'. The electrons can now only move freely in the directions parallel to the

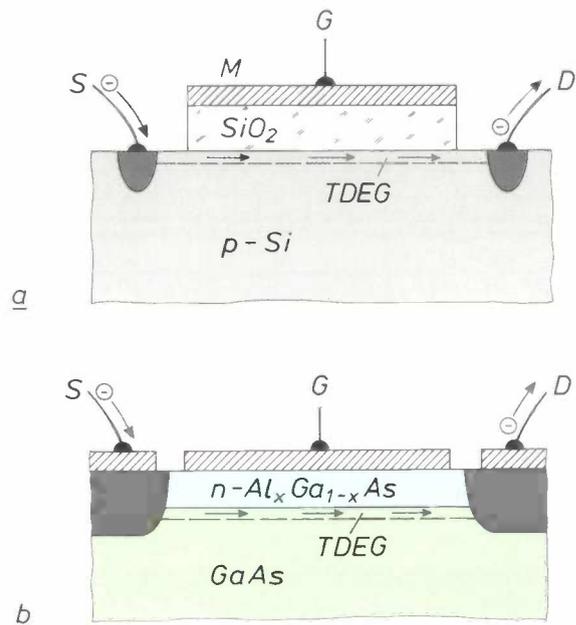


Fig. 5. *a*) Schematic cross-section of a MOSFET (metal-oxide-semiconductor field-effect transistor). The metal layer M is connected to the gate G and is separated from the p-type silicon by an insulating layer of SiO_2 . Contact regions connect the silicon to the source S and the drain D . At a positive gate voltage the electrons flow from S to D in a very thin inversion layer containing a two-dimensional electron gas (TDEG) near the interface with the SiO_2 layer. *b*) Schematic cross-section of a field-effect transistor with a modulation-doped heterojunction. The donor atoms of the $n\text{-Al}_x\text{Ga}_{1-x}\text{As}$ layer deliver electrons to the undoped GaAs layer. Near the interface a two-dimensional electron gas is formed at a positive gate voltage. Owing to the spatial separation from the parent donor atoms, the electrons here have a very high mobility.

interface. Partly because of their good switching characteristics MOSFETs are widely used in computers today.

If the p-type silicon layer is replaced by GaAs and the SiO_2 layer is replaced by n-type $\text{Al}_x\text{Ga}_{1-x}\text{As}$ (fig. 5b), structures with higher electron mobility can be produced. The two-dimensional electron gas is then formed in the undoped GaAs layer near the interface with the $n\text{-Al}_x\text{Ga}_{1-x}\text{As}$ layer. Since this layer is fully depleted, it behaves as an insulator. Whereas a MOSFET contains amorphous SiO_2 , the heterojunction here is formed by single-crystal semiconductors

with nearly perfect lattice matching. The absence of dislocations and asperities at the interface gives a further reduction in electron scattering.

This type of field-effect transistor is commonly referred to as a HEMT (high-electron-mobility transistor) [16]. Nowadays many industrial laboratories are fabricating and studying HEMTs. Integrated circuits containing these transistors have already been made. An important feature of HEMTs is the high speed at which they can process signals. Considerable amounts of charge are rapidly accepted and released, with relatively little energy dissipation [16]. Switching times of the order of 10^{-11} s have already been observed, which makes them faster, for the same energy dissipation, than any other semiconductor device. It seems highly likely that they will be widely used in advanced applications such as high-frequency amplifiers, supercomputers and super-fast signal-processing systems.

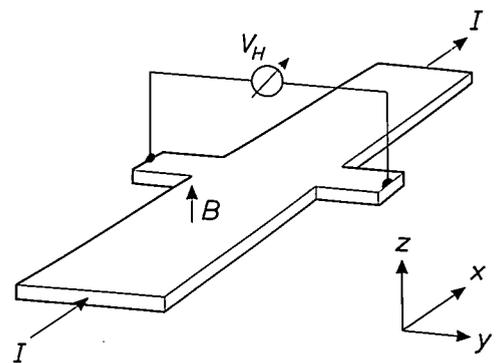
Modulation doping is also applicable to superlattices and quantum wells, of course. The combination with the minibands in superlattices offers much scope for speculation, but we shall leave the matter here.

Quantized Hall effect

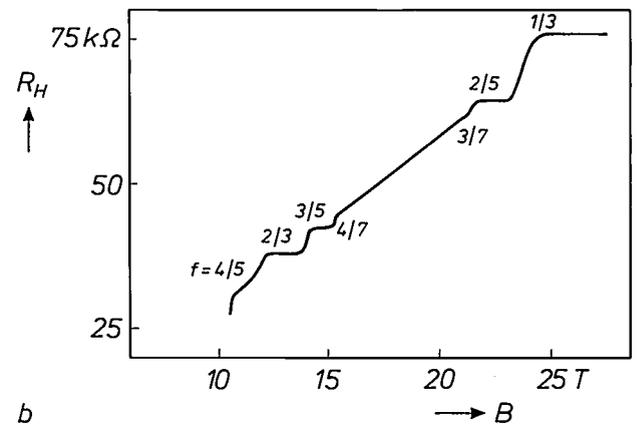
In a measurement of the Hall effect the quantity determined is the voltage induced in the y -direction (for example) of a sample by the combination of a current in the x -direction and a magnetic field in the z -direction:

see *fig. 6a*. The ratio of this voltage to the current is called the Hall coefficient (or Hall resistance). The structures described here, with a two-dimensional electron gas, can exhibit an intriguing effect, the quantized Hall effect, which was first observed in a silicon MOSFET [17]. At very low temperatures and strong magnetic fields there are plateaus in the curve of the Hall coefficient as a function of the magnetic field. The conductivity at these plateaus is independent of the experimental parameters and is given very accurately by h/ne^2 , where h is Planck's constant, n is an integer and e is the electronic charge. Klaus von Klitzing received the Nobel Prize for Physics in 1985 for his discovery of this effect.

There is as yet no full explanation of this effect. It is however generally accepted that it is related to the properties of the two-dimensional electron gas, whose energy levels for the motion perpendicular to the



a



b

Fig. 6. a) Arrangement for Hall-coefficient measurements. I current in the x -direction. B magnetic flux density in the z -direction. V_H induced Hall voltage in the y -direction. The Hall coefficient (or Hall resistance) R_H is given by V_H/I . b) The fractional quantized Hall effect measured for a heterojunction of GaAs and $Al_xGa_{1-x}As$ at a very low temperature and a very strong magnetic field [20]. The Hall-coefficient curve has plateaus with a value of h/fe^2 , where h/e^2 is equal to 25813Ω and f is a fraction of 1.

- [7] J. C. J. Finck, H. J. M. van der Laak and J. T. Schrama, A semiconductor laser for information read-out, *Philips Tech. Rev.* **39**, 37-47, 1980.
- [8] Special issue 'Compact Disc Digital Audio', *Philips Tech. Rev.* **40**, 149-180, 1982.
- [9] K. Bulthuis, M. G. Carasso, J. P. J. Heemskerk, P. J. Kivits, W. J. Kleuters and P. Zalm, Ten billion bits on a disk, *IEEE Spectrum* **16**, No. 8 (August), 26-33, 1979.
- [10] N. Holonyak, Jr., R. M. Kolbas, R. D. Dupuis and P. D. Dapkus, Quantum-well heterostructure lasers, *IEEE J. QE-16*, 170-186, 1980.
- [11] W. T. Tsang, Extremely low threshold (AlGa)As modified multiquantum well heterostructure lasers grown by molecular-beam epitaxy, *Appl. Phys. Lett.* **39**, 786-788, 1981.
- [12] R. Dingle, H. L. Störmer, A. C. Gossard and W. Wiegmann, Electron mobilities in modulation-doped semiconductor heterojunction superlattices, *Appl. Phys. Lett.* **33**, 665-667, 1978.
- [13] A concise survey of the electron mobilities in modulation-doped heterostructures of GaAs and $Al_xGa_{1-x}As$ is given in: H. L. Störmer, *Surf. Sci.* **132**, 519-526, 1983.
- [14] Measurement by C. T. Foxon and J. J. Harris, to be published shortly.
- [15] T. Mimura, S. Hiyamizu, T. Fujii and K. Nanbu, A new field-effect transistor with selectively doped GaAs/n- $Al_xGa_{1-x}As$ heterojunctions, *Jap. J. Appl. Phys.* **19**, L225-L227, 1980.
- [16] T. Mimura, Why HEMT are necessary and how they are made, *J. Electron. Eng.* **20**, No. 200, 60-62, 1983; H. Morkoc and P. M. Solomon, The HEMT: a superfast transistor, *IEEE Spectrum* **21**, No. 2 (February), 28-35, 1984.
- [17] K. von Klitzing, G. Dorda and M. Pepper, New method for high-accuracy determination of the fine-structure constant based on quantized Hall resistance, *Phys. Rev. Lett.* **45**, 494-497, 1980.

interface are quantized. Because of the magnetic field, quantization also occurs parallel to the interface, giving rise to the formation of 'Landau levels'. These levels are successively depleted as the magnetic field is increased. Whenever a level becomes fully depleted, there is a plateau in the curve of the Hall coefficient. To explain this it is generally assumed that between the Landau levels there are localized energy states in which the electrons do not contribute to the conductivity. It may also be significant that the electrons are not homogeneously distributed throughout the sample [18].

Soon after the discovery of the quantized Hall effect, it was also found in heterostructures of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ [19]. Further investigations of these structures revealed another unexpected effect, the fractional quantized Hall effect [20]. It was found that plateaus also occur where the Hall coefficient is given by h/fe^2 , with $f = \frac{1}{3}, \frac{2}{5}, \frac{2}{3}, \dots$; see fig. 6b. This means that plateaus must also occur when the Landau levels are only partially occupied. This effect cannot be explained from the single-electron model. This provides further motivation for much more theoretical and experimental investigation of these structures.

The observation of the quantized Hall effect also has practical consequences in view of the high accuracy with which the Hall coefficient at the plateaus can be experimentally determined. Attempts are being made in various standards institutions to use this effect for creating a new reference resistance with an accuracy better than 1 in 10^8 .

Growth technologies

The progress made in research on the structures described here has only been possible because of the development of advanced methods for the epitaxial growth of thin-film semiconductor layers. The potential applications that have emerged have provided a tremendous stimulus. At the same time the growth technologies have greatly stimulated the development of sophisticated semiconductor devices.

Two growth technologies are pre-eminent for the ability to control compositions and layer thicknesses and for obtaining virtually defect-free surfaces and interfaces. One of them is metal-organic vapour-phase epitaxy (MO-VPE). In this technology layers of a material such as $\text{Al}_x\text{Ga}_{1-x}\text{As}$ are deposited on a heated GaAs substrate from a reactive gas mixture containing AsH_3 and metal-organic compounds of aluminium and gallium [21] [22]. The other technology is molecular beam epitaxy (MBE), which is essentially a special form of evaporation in ultra-high vacuum [23] [24]. Here the layers are deposited by the re-

action of molecular (or atomic) beams incident on the heated substrate.

The success of these technologies is largely due to the facilities available for controlling the growth and analysing the structures. They include the use of photoluminescence, transmission electron microscopy and various methods of surface analysis [25]. With these technologies structures can be made that are virtually defect-free, with transitions that are abrupt on an atomic scale.

Materials

The materials that have received most attention are GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$. They have very similar lattice constants, which means that epitaxial structures on a GaAs substrate are nearly free of stress. These structures have proved suitable for testing new ideas and for practical evaluation. Some work has also been done with other materials, mainly III-V semiconductors. However, the considerable scope offered by modification of the band structure discussed here will be a considerable stimulus to the use of other materials.

Many materials cannot easily be combined because their lattice constants are too different. This inhibits perfect epitaxial growth and induces undesired dislocations at the interface. Some lattice mismatch can nevertheless be accommodated by an elastic strain in the crystal lattice, without the generation of misfit dislocations at the interface. Interesting combinations of semiconductors that have very different properties can be made in this way.

An example of such a combination is silicon on GaAs [24], with a 4% difference between the lattice constants. With such a structure integrated circuits in silicon could be combined with optoelectronic

[18] R. Woltjer, R. Eppenga, J. Mooren, C. E. Timmering and J. P. André, A new approach to the quantum Hall effect, *Europhys. Lett.* 2, 149-155, 1986.

[19] D. C. Tsui and A. C. Gossard, Resistance standard using quantization of the Hall resistance of $\text{GaAs-Al}_x\text{Ga}_{1-x}\text{As}$ heterostructures, *Appl. Phys. Lett.* 38, 550-552, 1981.

[20] D. C. Tsui, H. L. Störmer and A. C. Gossard, Two-dimensional magnetotransport in the extreme quantum limit, *Phys. Rev. Lett.* 48, 1559-1562, 1982.

A survey of the fractional quantized Hall effect is given by H. L. Störmer in: *Festkörperprobleme* 24, 25-44, 1984.

[21] P. M. Frijlink, J. P. André and M. Erman, this issue, pp. 118-132.

[22] M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft, this issue, pp. 133-142.

[23] B. A. Joyce and C. T. Foxon, this issue, pp. 143-153.

[24] P. C. Zalm, C. W. T. Bulle-Lieuwma and P. M. J. Marée, this issue, pp. 154-165.

[25] See the other articles in this issue. A general description of various methods of surface analysis is given in: H. H. Brongersma, F. Meijer and H. W. Werner, *Philips Tech. Rev.* 34, 357-369, 1974.

[26] J. C. Bean, J. C. Feldman, A. T. Fiory, S. Nakahara and I. K. Robinson, $\text{Ge}_x\text{Si}_{1-x}/\text{Si}$ strained-layer superlattice grown by molecular beam epitaxy, *J. Vac. Sci. & Technol. A* 2, 436-440, 1984.

components, such as lasers and photodiodes, in GaAs, on the same wafer. Another example is a strained-layer structure of silicon and SiGe. It has been demonstrated that dislocation-free SiGe layers as thick as 75 nm can be grown on silicon [26]. The ideas of modulation doping and band-structure modification can also be applied in this way for silicon. In layer structures that have an elastic strain degeneracies in the conduction band can also be eliminated, so that higher electron mobilities are possible in principle.

Structures of materials with a moderate lattice mismatch are now being widely studied by many materials scientists. It is hard to say whether this will lead to novel devices. The only certainty is that there will be many problems to overcome in the fabrication.

Research at Philips

In recent years Philips have done a great deal of work on thin-film epitaxial semiconductor structures. The study of MO-VPE and MBE processes has resulted in a continuous improvement in the control of the growth of multilayer structures. A variety of interesting structures have been fabricated and their properties have been investigated, with particular attention to the fundamental physics of the structures. Potential device applications, such as the quantum-well laser and the HEMT, are also topics of research.

There is work on MO-VPE growth at the research laboratories in Limeil-Brévannes, in France, and in Eindhoven. Structures of various III-V semiconduc-

tors are being investigated, but the emphasis is on combinations of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$. The same applies to the work on MBE at the Redhill laboratories, in England. Another MBE activity has been the growth of silicon layers on GaP and GaAs, a joint project undertaken by the FOM Institute for Atomic and Molecular Physics in Amsterdam and the Philips Research Laboratories in Eindhoven. At the Briarcliff laboratories, in the U.S.A., work has recently started on MBE of II-VI semiconductors. MBE of silicon and silicon-germanium structures is also being studied at the Eindhoven laboratories.

Some of these activities are reviewed in the following articles in this special issue. The first two deal with various aspects of the growth of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ structures by MO-VPE. The growth of similar structures by MBE is then discussed. The final article deals with MBE of silicon films on GaP and GaAs.

Summary. Research on layered epitaxial structures of semiconductor materials is a field of rapidly growing significance. Structures of special interest are the superlattice consisting of alternate ultra-thin layers of two dissimilar semiconductors, the quantum well formed by a thin semiconductor layer sandwiched between two semiconductor layers with a larger band gap, and the modulation-doped heterojunction. The controlled modification ('tailoring') of the energy-band structure gives some fascinating physical effects. There are also potential new applications, such as quantum-well lasers and high-electron-mobility transistors. Suitable growth technologies are metal-organic vapour-phase epitaxy (MO-VPE) and molecular beam epitaxy (MBE). Various aspects of fundamental physics, growth procedures and device applications are topics of research at several Philips laboratories.

Metal-organic vapour-phase epitaxy of multilayer structures with III-V semiconductors

P. M. Frijlink, J. P. André and M. Erman

Introduction

Thin single-crystal films of III-V semiconductor materials such as GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ are now widely used in devices such as solid-state lasers^[1] and microwave field-effect transistors^[2]. The films are obtained by epitaxial growth, in which they adopt the crystal structure of the substrate (usually GaAs). Conventionally, the films are grown by deposition from a melt or a solution (liquid-phase epitaxy, LPE) or from a reactive gas mixture (vapour-phase epitaxy, VPE). In recent years the control over these growth processes has been improved considerably to reduce the dimensions of the semiconductor devices and circuits. There is however only a limited capability for growing multilayer structures with very abrupt interfaces, required for particular applications as mentioned in the opening article^[3].

In the last ten years or so two growth techniques have been developed which can meet the special requirements for interface abruptness. One is a modified version of VPE, in which the reactive gas mixture contains metal-organic compounds (metal-organic vapour-phase epitaxy, MO-VPE), the other is molecular beam epitaxy (MBE), which will be discussed later in this issue^[4]. An important feature of these techniques is that the chemical composition near the growing surface can be changed in a time interval an order of magnitude smaller than the time necessary for growing a single atomic layer. This has made it possible to grow multilayer structures with sharp interfaces on the scale of one atomic monolayer. Consequently, new devices requiring such abrupt interfaces can now be made. An example of such a device is the $\text{Al}_x\text{Ga}_{1-x}\text{As}/\text{GaAs}$ quantum-well laser: a

semiconductor laser with a very thin active GaAs layer between $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers, which gives a narrow emission peak at a relatively short wavelength with a low threshold current^[5]. Another example is the high-electron-mobility transistor (HEMT): a microwave field-effect transistor with an $\text{Al}_x\text{Ga}_{1-x}\text{As}/\text{GaAs}$ heterojunction, which has an excellent high-frequency performance^[6].

At LEP the MO-VPE growth of III-V layer structures has been extensively investigated^[7]. Several aspects have been considered, such as the development of a suitable growth-reactor system and the assessment of the relevant growth parameters for obtaining multilayer structures with the required properties. The layers and their interfaces have been characterized by spectroscopic ellipsometry and photoluminescence measurements. In addition, the applicability of various structures for practical devices has been studied.

It has been found that MO-VPE is a reliable and versatile method for growing a large variety of device-quality III-V semiconductor materials and multilayer structures. It has proved possible to grow successive layers with precise control of composition and lattice match to the substrate, giving high electrical and optical quality. Precise compositional control is possible not only with ternary materials such as $\text{Al}_x\text{Ga}_{1-x}\text{As}$ but also with quaternary materials such as $\text{In}_x\text{Ga}_{1-x}\text{P}_y\text{As}_{1-y}$. Doping profiles are controlled with a resolution of less than 10 nm. MO-VPE growth is particularly successful for the preparation of quantum-well lasers and HEMT devices based on GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$.

A large number of results have shown that MO-VPE and the characterization methods have great potential, and that practical devices can be made.

Ir P. M. Frijlink, Dr J. P. André and Dr M. Erman are with Laboratoires d'Electronique et de Physique Appliquée (LEP), Limeil-Brévannes, France.

Some examples are:

- The photoluminescence of a quantum well in which the active layer is only about 0.6 nm thick.
- A strong visible luminescence of a quantum well with $\text{Al}_{0.2}\text{Ga}_{0.8}\text{As}$ in the active layer.
- A quantum-well laser emitting at 775 nm with a maximum power of 230 mW per facet in continuous-wave (CW) operation at 300 K and a CW threshold current density of 660 A/cm^2 .
- A two-dimensional electron gas at the interface of a GaAs/ $\text{Al}_x\text{Ga}_{1-x}\text{As}$ heterojunction giving an elec-

Finally, modulation-doped heterostructures will be discussed, with emphasis on the special properties related to the formation of a two-dimensional electron gas.

Preparation of multilayer structures with abrupt interfaces

A suitable reactor system for the MO-VPE growth of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ is shown in *fig. 1*. A heated GaAs substrate is exposed to a gaseous mixture containing

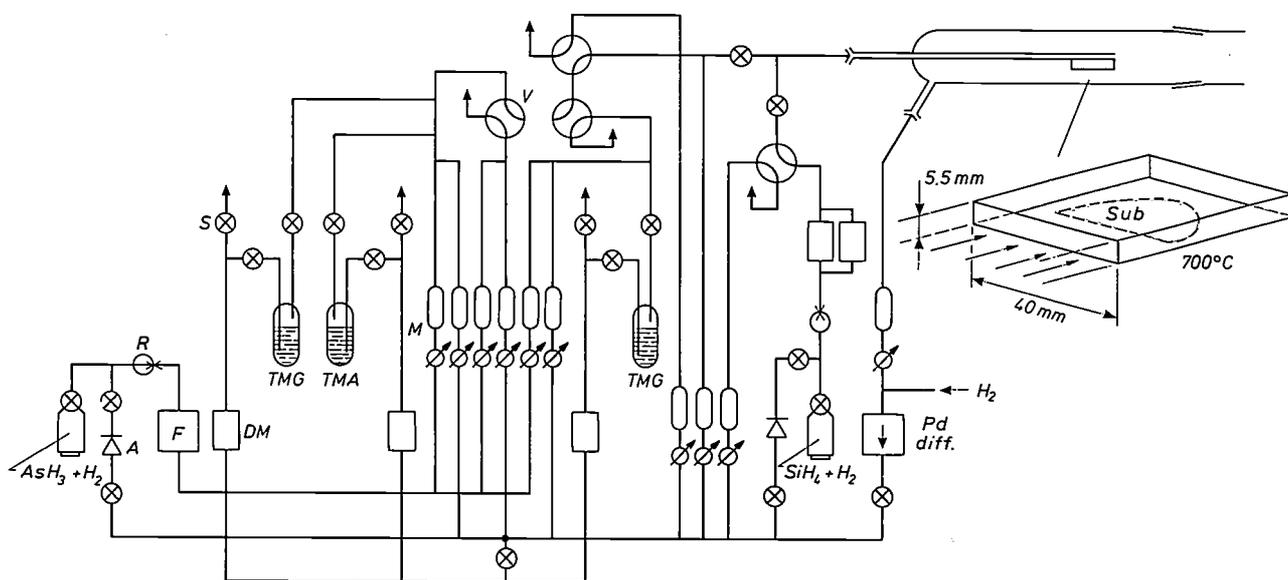


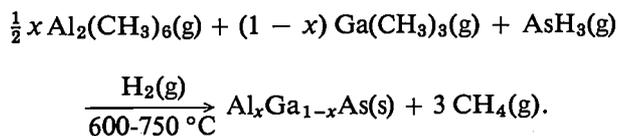
Fig. 1. Schematic diagram of the MO-VPE equipment used for growing $\text{Al}_x\text{Ga}_{1-x}\text{As}$ multilayer structures. In a quartz reactor layers are deposited from the gas phase on to a GaAs substrate *Sub*. This is kept at a temperature of about 700°C by inductive r.f. heating of the substrate holder. The gas phase contains hydrogen (used as the carrier gas) and the reactive components trimethyl aluminium (TMA), trimethyl gallium (TMG) and arsine (AsH_3), supplemented by doping reactants, e.g. silane (SiH_4) for n-doping. The hydrogen is purified by diffusion through a membrane of palladium (*Pd diff.*). The supply of the several gases can be controlled and mixed in such a way that any desired gas composition can be obtained very rapidly. *S* standard valve; *R* pressure regulator; *M* flow meter with regulating valve; *DM* mass-flow meter; *F* filter for gas purification; *A* non-return valve; *V* four-way valve.

tron Hall mobility of $7.5 \times 10^3 \text{ cm}^2\text{V}^{-1}\text{s}^{-1}$ at 300 K and of $2.7 \times 10^5 \text{ cm}^2\text{V}^{-1}\text{s}^{-1}$ at 4 K, yielding field-effect transistors with a transconductance as high as 250 S/m for a gate length of $1 \mu\text{m}$.

In this article, the discussion of our MO-VPE investigations will be restricted to multilayer structures of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$, the most widely investigated of the III-V structures and the most promising for immediate practical application. After a general description of the MO-VPE method, it will be shown how structures with abrupt interfaces can be obtained. It will be demonstrated that spectroscopic ellipsometry can be a useful tool for depth profiling and interface characterization. Next, the composition and luminescence properties of quantum wells will be dealt with, as well as the application in quantum-well lasers.

- [1] See for example J. C. J. Finck, H. J. M. van der Laak and J. T. Schrama, A semiconductor laser for information read-out, *Philips Tech. Rev.* **39**, 37-47, 1980.
- [2] See for example P. Baudet, M. Binet and D. Boccon-Gibod, Low-noise microwave GaAs field-effect transistor, *Philips Tech. Rev.* **39**, 269-276, 1980.
- [3] J. Wolter, Research on layered semiconductor structures, this issue, pp. 111-117.
- [4] B. A. Joyce and C. T. Foxon, Molecular beam epitaxy of multilayer structures with GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$, this issue, pp. 143-153.
- [5] One of the first descriptions of a quantum-well laser, operating at room temperature, has been given by R. D. Dupuis, P. D. Dapkus, N. Holonyak, Jr., E. A. Rezek and R. Chin, *Appl. Phys. Lett.* **31**, 295-297, 1978.
- [6] Successful fabrication of a HEMT device was reported for the first time by T. Mimura, S. Hiyamizu, T. Fujii and K. Nanbu, *Jap. J. Appl. Phys.* **19**, L225-L227, 1980.
- [7] Many of our colleagues at LEP contributed to the investigations described here. The MO-VPE growth of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ multilayer structures has also been investigated at the Philips Research Laboratories in Eindhoven, see M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft, this issue, pp. 133-142.

the metal-organic compounds trimethyl aluminium ($\text{Al}_2(\text{CH}_3)_6$) and trimethyl gallium ($\text{Ga}(\text{CH}_3)_3$) with arsine (AsH_3) and hydrogen (H_2), which is used as carrier gas. The pyrolysis of the reactive compounds at the substrate surface leads to the epitaxial deposition of a single-crystal $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layer in the overall reaction:



The composition of the deposited layer depends on the partial pressures of $\text{Al}_2(\text{CH}_3)_6$ and $\text{Ga}(\text{CH}_3)_3$ in the gas phase. Addition of substances such as silane (SiH_4) or diethyl zinc ($\text{Zn}(\text{C}_2\text{H}_5)_2$) to the gas phase leads to n- or p-doping of the layer.

Compared with the more conventional VPE and LPE methods, the MO-VPE method has some advantages for growing multilayer structures with abrupt interfaces: it is essentially a 'far-from-equilibrium' or 'one-way-deposition' process. Since the growth rate is essentially proportional to the supply of the reactants that provide the group III elements, decreasing the supply of these reactants will in principle make the growth rate arbitrarily small. By changing the gas composition very thin layers with abrupt compositional changes can be obtained. However, for high electrical and optical quality, these layers must be pure. If we assume a constant rate of contamination, e.g. from outgassing of the reactor walls, it is clear that an extremely small growth rate will result in a high contamination level in the solid grown. This means that there are essentially two ways of obtaining very pure heterostructures with very abrupt interfaces: either we use a reactor system in which we make sure that the overall contamination level is very low, and reduce the growth rate sufficiently to avoid transient control problems on changing the gas composition to grow an interface, or we aim for effective control of very rapid changes in the gas composition and use a high growth rate, thus reducing the requirements for contamination control. In view of the present 'state of the art' in MO-VPE reactors, we chose the second option.

In our investigations, the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers are generally deposited on substrates consisting of chromium-doped semi-insulating GaAs wafers. These wafers are sawn from a single-crystal ingot, grown by the Czochralski method. The mechanically polished substrate surface is disoriented with respect to the (001) crystal plane by a few degrees. The substrate is supported by a graphite susceptor. The appropriate substrate temperature, between 600 and 750 °C, is

obtained by inductive r.f. heating of the susceptor. A thermocouple in the susceptor measures the substrate temperature. The carrier gas is purified by diffusion through a palladium membrane. The supply of the various reactive compounds is controlled by mass-flow meters.

As pointed out, the MO-VPE growth of multilayer structures with abrupt interfaces requires an instantaneous and precise control of the partial pressures in the reactor system. In the system of fig. 1 the gas transport is arranged so that the gas composition over the wafer can be changed in a controlled way within 0.1 s, neglecting adsorption and desorption at the reactor walls. At a total gas pressure of 10^5 Pa (1 atm) and a suitably low growth temperature of 650 °C, the growth rate of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers can be as low as 0.5 nm/s. This means that the time required for changing the chemical composition at the crystal surface is much less than the growth time of one monomolecular layer of $\text{Al}_x\text{Ga}_{1-x}\text{As}$. This permits the growth of well-defined multilayer structures, with a reproducible control over composition, thickness and doping profile in the direction of growth on an atomic scale.

As will be discussed later, depth profiling and interface quality can be demonstrated for quantum wells by analysing their luminescence properties and for modulation-doped heterostructures by considering the increase in the electron mobility near the interface. However, we shall first discuss a more general characterization method: spectroscopic ellipsometry. This method has been found to be very useful in investigating the chemical and structural quality of the interfaces.

Depth profiling and interface characterization by spectroscopic ellipsometry

Ellipsometry is an optical technique that makes use of the change in polarization of light at the surface to be investigated [8]. The state of polarization of polarized light can be characterized by the amplitude A and the phase δ of two perpendicular components into which the total electric vector of the light can be resolved. The component parallel to the plane of incidence (E_p) and the component perpendicular to it (E_s) are given by:

$$\begin{aligned} E_p &= A_p \exp j(\delta_p + \omega t), \\ E_s &= A_s \exp j(\delta_s + \omega t), \end{aligned} \quad (1)$$

where j is the imaginary unit and ω is the angular frequency of the light wave. At arbitrary values of the amplitudes and phases, the tip of the total electric vector will describe an ellipse; the light is then said to be elliptically polarized. On reflection the change in

the parallel component is different from the change in the perpendicular component, so that the shape of the ellipse will change; see *fig. 2*.

In an ellipsometer the complex ratio ρ of the reflectances of the parallel and perpendicular components is determined; this can be expressed as:

$$\rho = \frac{(E_p/E_s)_r}{(E_p/E_s)_i}, \quad (2)$$

where the indices *r* and *i* refer to the reflected and incident beams. Substitution from eq. (1) gives the relation with the changes in amplitude and phase on reflection, which can be described by two parameters ψ and Δ :

$$\rho = \tan \psi \exp j \Delta, \quad (3)$$

$$\tan \psi = \frac{(A_p/A_s)_r}{(A_p/A_s)_i}, \quad (4)$$

$$\Delta = (\delta_p - \delta_s)_r - (\delta_p - \delta_s)_i. \quad (5)$$

For a reflection at a clean substrate, the quantities ψ and Δ can be correlated with its refractive index and absorption coefficient. If the substrate is coated with a thin layer, information about the thickness of the layer and its optical properties can also be obtained.

The layers we have deposited and their interfaces have been investigated by a special form of ellipsometry, in which ρ is determined as a function of the energy of the incident light (spectroscopic ellipsometry). Since ρ is directly related to the dielectric function of the reflective material, such measurements enable us to make a detailed structural and chemical analysis of the surface. The spectroscopic ellipsometer used in our experiments has the following optical components: a xenon lamp, a rotating polarizer, the sample, a fixed analyser, a double monochromator

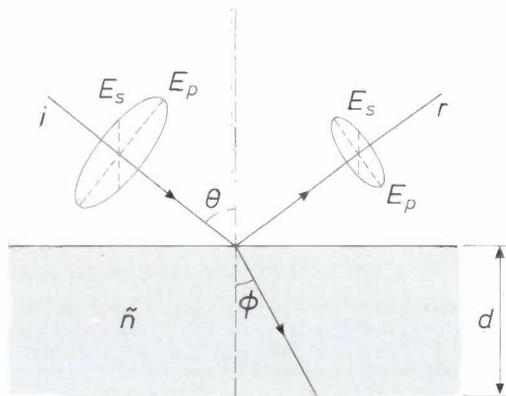


Fig. 2. Schematic representation of the reflection of elliptically polarized light at the surface of a layer. The parallel and perpendicular components E_p and E_s of the electric vector of the reflected light beam (*r*) differ from those of the incident beam (*i*). The ratio of the reflectances of the two components depends on the angle of incidence θ , the angle of refraction ϕ and the complex index of refraction \tilde{n} and the thickness d of the layer.

and a photomultiplier. The optical range is from 1.6 eV (because of the limited sensitivity at longer wavelengths of the S20 extended-UV photocathodes) to 5.4 eV (because of the use of calcite prisms for the polarizer and analyser). The output signal of the photomultiplier is sampled as a function of the angular position of the polarizer and a Fourier analysis is made by a minicomputer to determine the value of ρ .

The dielectric function ϵ can be derived at each energy by using the Fresnel expressions for the reflection coefficients of the parallel and perpendicular components and Snell's law. For a clean wafer ϵ is given by:

$$\epsilon = \epsilon_1 - j\epsilon_2 = \frac{(1 - \rho)^2}{(1 + \rho)^2} \sin^2 \theta \tan^2 \theta + \sin^2 \theta, \quad (6)$$

where θ is the angle of incidence. In *fig. 3a* the continuous lines give the dielectric function for a carefully cleaned GaAs wafer under vacuum at room temperature. A characteristic feature is the presence of a double-peak structure in the ϵ_2 curve, corresponding to different optical transitions (at 2.93 and 3.16 eV) between the valence and conduction bands^[9]. The spectrum of the dielectric function of GaAs is very sensitive to crystalline imperfections induced for example by high temperatures (lattice vibrations), doping or Al incorporation. Generally speaking, any type of imperfection tends to reduce the ϵ_2 peak height and especially the contrast in the double-peak structure. Evaluation of the spectrum obtained can therefore give useful information about the ordering of the GaAs lattice and the doping level or Al incorporation^[10].

The dielectric functions of various $\text{Al}_x\text{Ga}_{1-x}\text{As}$ alloys have also been measured. As an example, the continuous lines in *fig. 3b* give the spectrum for an undoped $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ sample. Compared with the

^[8] See for example D. E. Aspnes, Spectroscopic ellipsometry of solids, in: B. O. Seraphin (ed.), Optical properties of solids, New developments, North-Holland, Amsterdam 1976, pp. 799-846.

Surface analysis by ellipsometry has also been described in this journal: K. H. Beckmann, Optical investigations of semiconductor surfaces, Philips Tech. Rev. 29, 129-142, 1968; F. Meijer and G. A. Bootsma, Investigation of the chemical behaviour of clean silicon and germanium surfaces, Philips Tech. Rev. 32, 131-140, 1971.

^[9] A survey of the various valence-conduction band transitions in GaAs has been given by Y. Petroff in: M. Balkanski (ed.), Optical properties of solids, North-Holland, Amsterdam 1980, Chapter 1.

^[10] M. Erman, J. B. Theeten, N. Vodjdani and Y. Demay, Chemical and structural analysis of the GaAs/AlGaAs heterojunctions by spectroscopic ellipsometry, J. Vac. Sci. & Technol. B 1, 328-333, 1983; M. Erman, J. B. Theeten, P. Frijlink, S. Gaillard, Fan Jia Hia and C. Alibert, Electronic states and thicknesses of GaAs/GaAlAs quantum wells as measured by electroreflectance and spectroscopic ellipsometry, J. Appl. Phys. 56, 3241-3249, 1984.

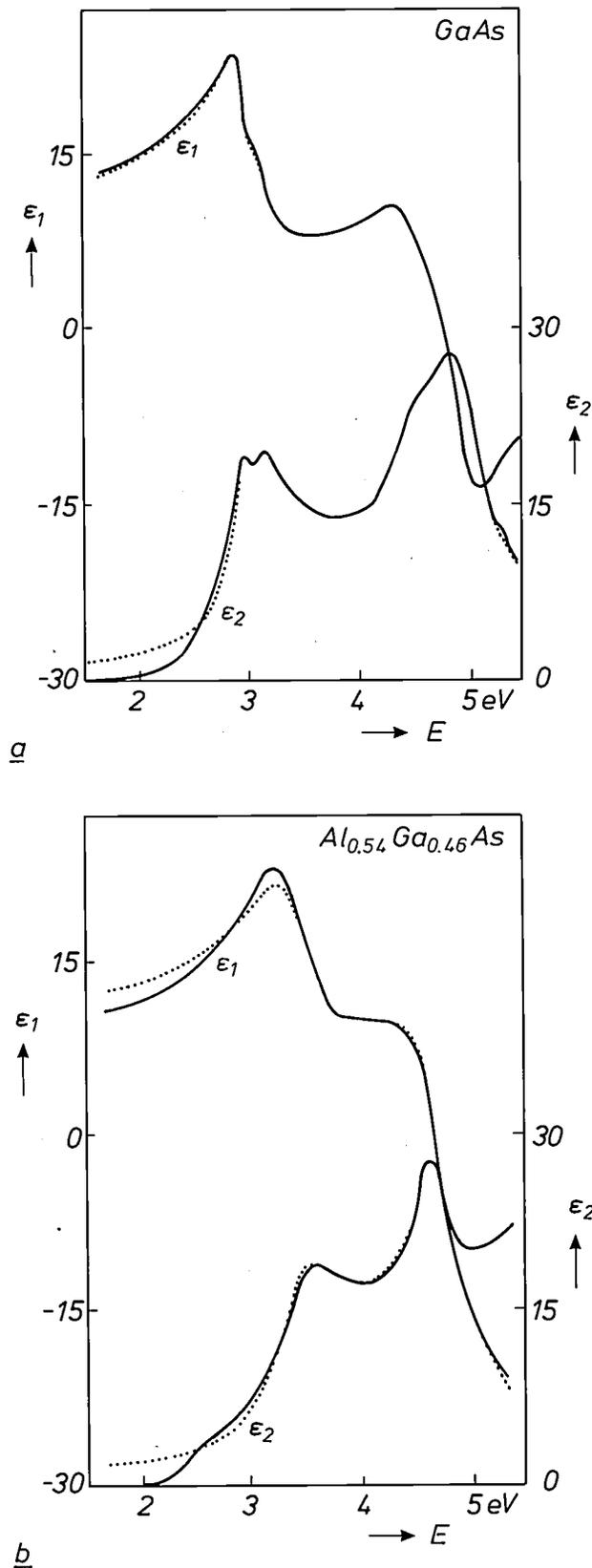


Fig. 3. Real and imaginary parts ϵ_1 and ϵ_2 of the dielectric functions of GaAs (a) and $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ (b), plotted against the energy E of the incident light. The continuous curves were derived from the measured complex ratio ρ as given by eq. (6). The ϵ_2 curve for GaAs has a characteristic double-peak structure at about 3 eV. The dotted curves refer to simulations generated by using a set of seven harmonic oscillators (see text). Above 2.5 eV good agreement with the measurements is obtained.

pure GaAs case, the spectrum shows no double-peak structure in the neighbourhood of 3 eV, but otherwise it has the same general features. It has been demonstrated that the spectrum of an arbitrary $\text{Al}_x\text{Ga}_{1-x}\text{As}$ alloy can be derived by interpolation. This can be done accurately by describing the variation of ϵ with the energy E as a sum of harmonic oscillators, each characterized by its amplitude A_i , its centre position E_i and its half-width Γ_i :

$$\epsilon = \sum_i A_i \{ (E_i - E - j\Gamma_i)^{-1} + (E_i + E + j\Gamma_i)^{-1} \}. \quad (7)$$

The GaAs dielectric function can be simulated by a set of seven oscillators as shown in fig. 3a by the dotted lines. The description is satisfactory above 2.5 eV. The basic idea for simulating the dielectric function of any $\text{Al}_x\text{Ga}_{1-x}\text{As}$ alloy is to use the *same* set of seven oscillators with a simple variation (quadratic in x) of A_i , E_i and Γ_i . The law of variation can be determined from measurements at a few x -values. This simulation method is illustrated by the dotted lines in fig. 3b for $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$. Good agreement with the measurements is again obtained above 2.5 eV.

The results for pure wafers can be used for simulations on actual samples of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers deposited on a GaAs substrate. These simulations are based on the 'effective-medium approximation' [11] for thicknesses less than the wavelength of the incident light. This means that the dielectric function ϵ of a mixture of two materials A and B, with known dielectric functions ϵ_A and ϵ_B , can be evaluated from the relation:

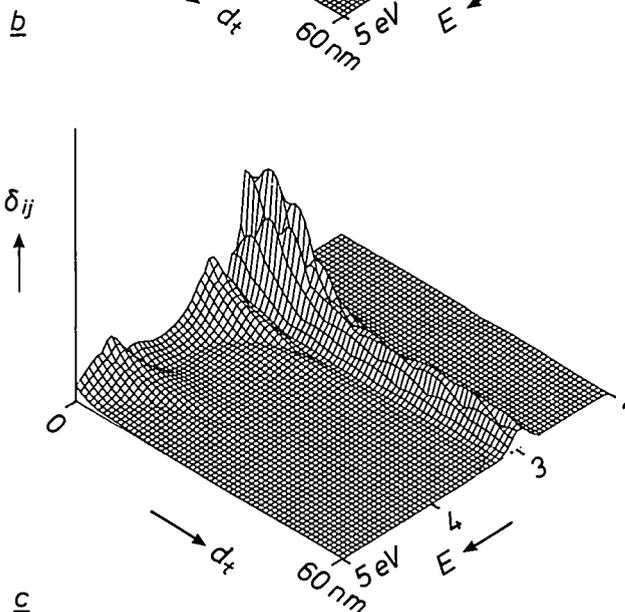
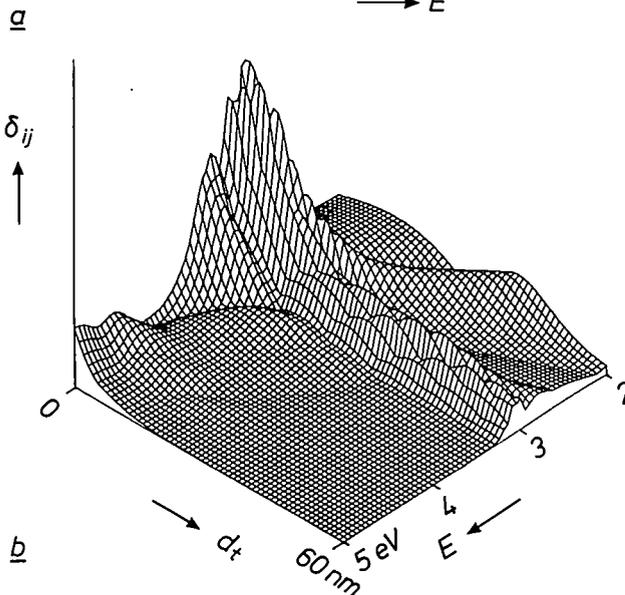
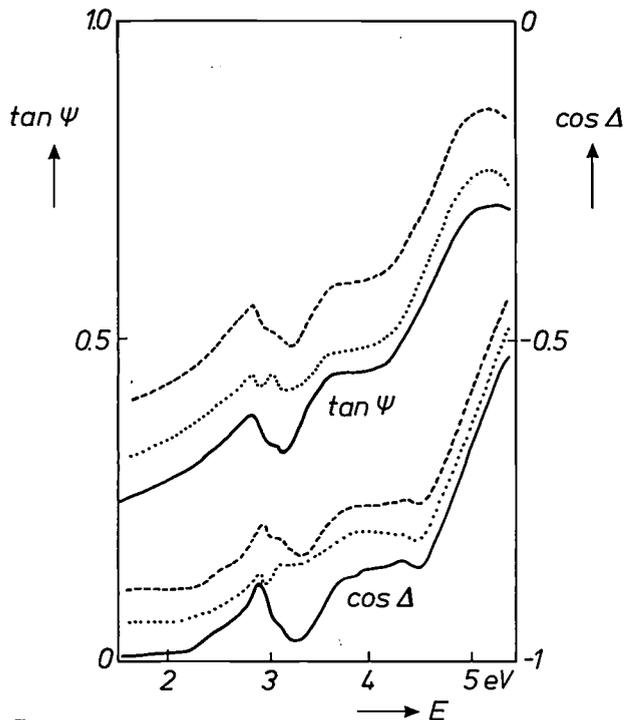
$$v(\epsilon_A - \epsilon)(\epsilon_A + 2\epsilon)^{-1} + (1 - v)(\epsilon_B - \epsilon)(\epsilon_B + 2\epsilon)^{-1} = 0, \quad (8)$$

where v is the volume fraction of A.

In evaluating and understanding the specific properties of heterojunctions and multilayer structures, an important aspect is the structure, on the atomic scale, of the transition regions between the various layers. In the ideal case of an abrupt heterojunction, the top of the semiconductor substrate A has the same composition and crystalline structure as the bulk of A, while the bottom of the deposited layer of the other semiconductor B has the same composition and structure as the bulk of B. Two types of deviation can occur from this ideal case:

- Interdiffusion of A and B and formation of an alloy with a chemical nature distinct from both materials: this will be referred to as a 'chemical' interface.
- Interpenetration of A and B, without changing their chemical nature, so that their separation is not planar

[11] C. G. Granqvist and O. Hunderi, Optical properties of ultra-fine gold particles, *Phys. Rev. B* 16, 3513-3554, 1977.



but rough on the atomic scale; this will be referred to as a 'physical' interface.

Investigations of different types of structures by spectroscopic ellipsometry have demonstrated the sensitivity of this method to the presence and nature of interface regions. As an example, *fig. 4a* gives the simulated curves of $\tan \psi$ and $\cos \Delta$ plotted against E for a structure with an ideal interface and for structures with chemical and physical interfaces. At the critical energy of about 3 eV (*fig. 3*) the curves for the ideal and physical interfaces are quite similar, while the curves for the chemical interface differ markedly. This is due to the presence of the intermediate composition at the chemical interface, giving rise to a shifted double-peak structure. This indicates that spectroscopic ellipsometry on $\text{Al}_x\text{Ga}_{1-x}\text{As}$ structures will be especially sensitive to the presence of a chemical interface.

A more extensive comparison is obtained by considering the mean-square difference δ_{ij} between the cases i and j , given by

$$\delta_{ij} = (\tan \psi_i - \tan \psi_j)^2 + (\cos \Delta_i - \cos \Delta_j)^2. \quad (9)$$

In *fig. 4b* and *4c* the calculated differences are given as a function of the incident light energy E and the total thickness d_t . Because of the optical absorption of the top layer, the interface region cannot be detected when $E \geq 3.5$ eV and $d_t \geq 15$ nm. The major contrast at about 3 eV between chemical and ideal is a factor of two greater than between physical and chemical. Even at this energy the contrast is strongly reduced when d_t becomes larger than 30 nm. This means that d_t should be of the order of 10 nm to achieve a satisfactory interface analysis. This limitation can be taken as an advantage when successive heterojunctions on the same substrate have to be investigated. An *in situ* analysis will be possible for each individual heterojunction, provided that sufficient material is deposited between two successive transitions.

The comparison between modelling and experiment will be illustrated by the results obtained

4
Fig. 4. Theoretical comparison between the ellipsometry properties of three types of interface. *a*) Calculated curves of $\tan \psi$ and $\cos \Delta$ as a function of the energy E . The continuous curves refer to an ideal interface between GaAs and 10 nm $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$, the dotted curves to a chemical interface region of 5 nm with the intermediate composition $\text{Al}_{0.27}\text{Ga}_{0.73}\text{As}$ and the dashed curve to a physical interface region that is a mixture of GaAs and $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$. For clarity the curves have been shifted with respect to one another in the vertical direction. At about 3 eV, the deviation from the ideal interface is larger for the chemical interface than for the physical one. *b*) Mean-square difference δ_{ij} , given by eq. (9), between the chemical and ideal interfaces as a function of the total thickness d_t and the energy E . The greatest contrast is obtained at $d_t \approx 10$ nm and $E \approx 3$ eV. *c*) Mean-square difference δ_{ij} between the chemical and physical interfaces plotted against d_t and E . The maximum contrast is about half that in (*b*).

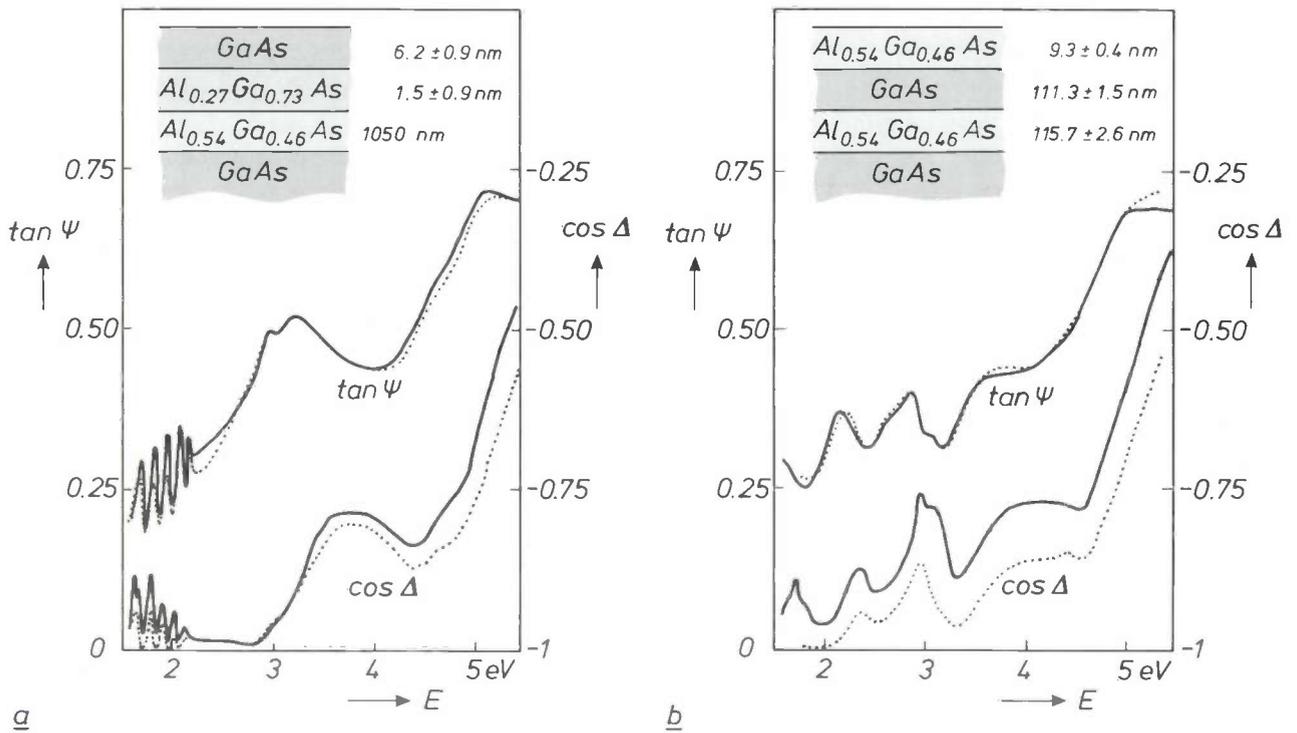
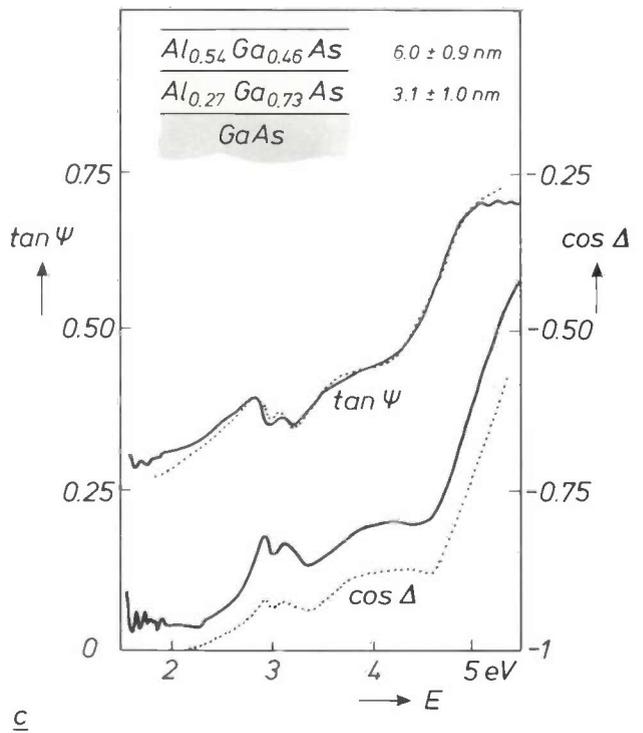


Fig. 5. Results of ellipsometry experiments and modelling for three types of heterostructure. The experimentally determined values of $\tan \psi$ and $\cos \Delta$ are plotted against the energy E (continuous curves). The modelling was performed for the $\tan \psi$ curves. The sequences of layer compositions and thicknesses shown refer to the best simulations (dotted curves). A chemical interface region is detected for a thin GaAs layer on $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ (a). No interface region is detected for a thin $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ layer on GaAs (b). When there is a rapid alternation between these compositions, a chemical interface region is also formed in this case (c).

on three samples of multilayer structures. Fig. 5a gives the measured data and best simulations for a thin GaAs layer on $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$. The oscillations between 1.6 and 2.2 eV are due to interferences in the $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ layer and can be used to evaluate the thickness of this layer (1050 ± 10 nm). Because of the presence of a thin native oxide layer (measurements were made in air), the comparison between modelling and experiment is made for the $\tan \psi$ curve only, since the oxide layer induces a small shift in the $\cos \Delta$ curve. The comparison indicates the presence of a chemical interface region, with a thickness of 1.5 ± 0.9 nm.

Fig. 5b shows the results for a thin $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ layer on GaAs. The shift in the $\cos \Delta$ curve is larger than for the first sample, and this can be attributed to a thicker native oxide layer due to the presence of Al in the top layer. The best modelling is obtained by assuming no interface region (ideal interface). The conclusion can therefore be drawn that the $\text{Al}_x\text{Ga}_{1-x}\text{As}/\text{GaAs}$ heterojunction (the 'normal' situation) is essentially abrupt and is better defined than the converse situation discussed earlier. A careful error analysis of the results shows that the maximum transition thickness is about 0.5 nm,



compared with about 1.5 nm in the converse situation.

To test the ability of spectroscopic ellipsometry to detect an interface region we have also examined a structure in which a 3-nm interface region had been deliberately grown. This was done by rapid alternation of the gas flows corresponding to GaAs and $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ in such a way that the interface could be considered as a sequence of twenty layers of

alternating composition, with each individual layer about 0.15 nm thick. Since this thickness is below that of one monolayer a chemical mixture would be expected: the interface region should look like an $\text{Al}_{0.27}\text{Ga}_{0.73}\text{As}$ layer about 3 nm thick. This is confirmed experimentally; see fig. 5c. The modelling assuming a chemical interface gives the best fit with the experimental data, with fewer uncertainties in the layer thicknesses. The best simulation is obtained with an interface of 3.1 ± 1.0 nm.

Quantum wells

When a GaAs layer is sandwiched between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers that have a larger band gap, there are steps in the energy-level diagram at the interfaces [11]. The existence of energy bands is due to the periodicity of the crystals, but their configuration is essentially determined by interactions over a range of about one unit cell of the crystal. In fact, the magnitude of the band gap of each material appears to be maintained right up to the interface. The band-gap difference leads to discontinuities in the conduction and valence bands at the interfaces. The conduction-band discontinuity is related to the difference in electron affinity between the two materials and is now assumed to be about 60% of the band-gap difference [12].

An interesting quantization effect is observed when the GaAs layer becomes thinner than the mean free path of the electrons between two collisions with lattice vibrations (phonons), lattice imperfections such as point defects, or impurity atoms [13]. The electrons trapped in the GaAs layer can then be reflected many times between the two interfaces. If the half-wavelength of the electrons in the direction perpendicular to the interfaces does not fit into the quantum well an integer number of times, the state will only have a very short lifetime, because of extinction by interference. Only the particular states that do fit in this way can exist permanently. Thus for the movement perpendicular to the interfaces only discrete energy values are permissible.

In directions parallel to the interfaces the electrons can move freely, and to a first-order approximation this does not affect the wavelength of the perpendicular movement. This implies that the electrons may have some extra energy. Each permitted state in the perpendicular direction therefore gives rise to a discrete energy band, whose lower band edge, characterized by zero movement in the parallel directions, is determined by the energy necessary to give a wavelength of $2d/n$, where n is an integer and d is the well thickness plus an effective ('tunnelling') penetration

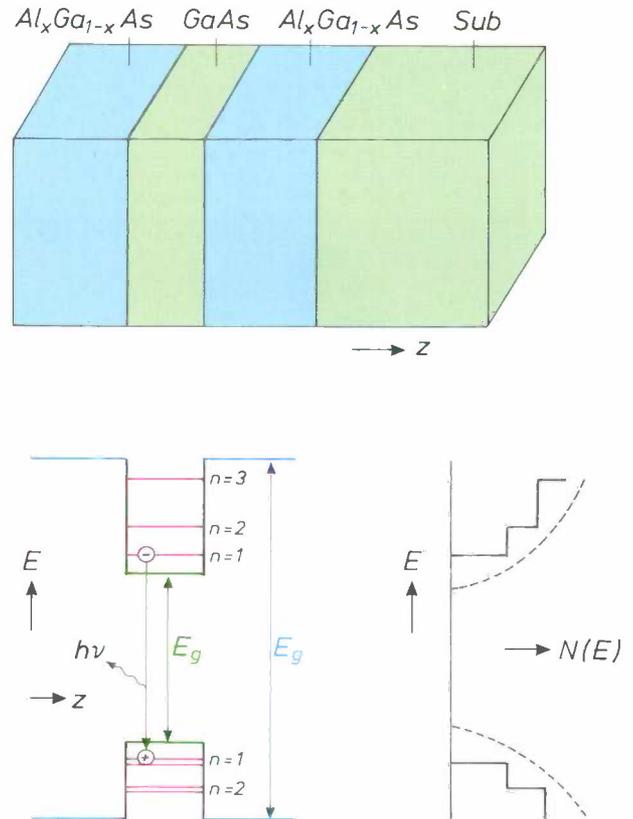


Fig. 6. Simplified energy-level diagram and density of states for a quantum well formed on a GaAs substrate *Sub* by a very thin GaAs layer between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers with a larger band gap E_g . The confinement in the narrow potential well leads to discrete energy levels for electrons and holes that move perpendicular to the interfaces. These levels are characterized by the quantum number n . For each quantum state two valence levels exist, corresponding to 'heavy' holes and 'light' holes. The energy bands corresponding to movement parallel to the interfaces are not shown, their edges represent zero movement. The variations in the density of states $N(E)$ with E is not continuous as in the bulk material (*dashed curves*), but occurs in steps (*continuous curves*). This means that the luminescence resulting from the recombination of electrons and heavy holes with $n = 1$ can have a high efficiency. The energy difference between the states associated with the luminescence is larger than in the bulk material, so that the luminescence occurs at shorter wavelengths.

width of the wave into the barriers. The lowest-energy state will be the state with $n = 1$. There is an analogous effect for the holes, with the complication that two types of holes exist in GaAs, with light and heavy effective mass. The resulting energy-level diagram and density of states are shown in fig. 6.

Luminescence

If we create electron-hole pairs in a quantum well by shining light on it, the electrons and holes will lose energy through a number of rapid relaxation pro-

[12] R. C. Miller, D. A. Kleinman and A. C. Gossard, Energy-gap discontinuities and effective masses for GaAs- $\text{Al}_x\text{Ga}_{1-x}\text{As}$ quantum wells, *Phys. Rev. B* **29**, 7085-7087, 1984.

[13] The possible quantum states in heterostructures with ultra-thin semiconductor layers have been compiled by R. Dingle, *Festkörperprobleme* **15**, 21-48, 1975.

cesses, and they will thermalize in the $n = 1$ states, with the particle distribution approaching a Fermi-Dirac distribution around a quasi-Fermi level. After this there is a much slower (e.g. 1000 times) recombination process. At each recombination, a photon will be emitted at the wavelength corresponding to the energy difference between an electron and a hole, both within a few times kT from the band edges of their respective $n = 1$ levels. With a thinner quantum well these band edges are more widely separated and the emitted photons will therefore have a shorter wavelength. A theoretical relation between the quantum-well thickness and the emission wavelength is given in fig. 7. This relation now permits us to assess the thickness of the GaAs quantum well directly by

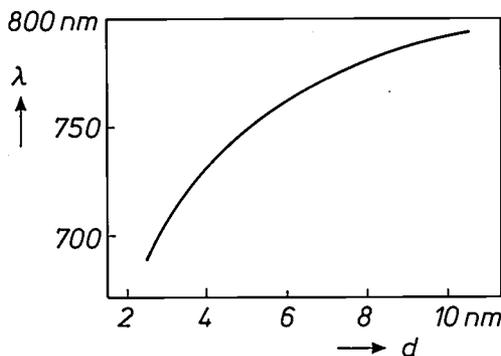


Fig. 7. Calculated relation between well thickness d and emission wavelength λ corresponding to the recombination of electrons and heavy holes in the $n = 1$ state, for a rectangular GaAs quantum well between $\text{Al}_{0.54}\text{Ga}_{0.46}\text{As}$ barriers. For a decrease in thickness from 12 to 2.5 nm the calculated reduction in wavelength is from 800 to 680 nm.

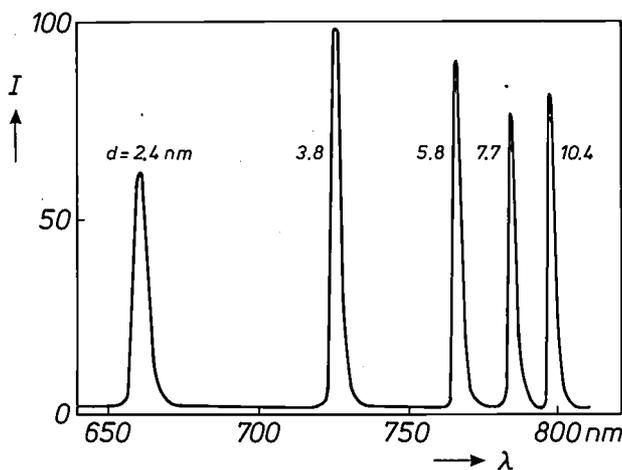


Fig. 8. Photoluminescence spectrum of a multiple-quantum-well sample on excitation at 4 K with radiation at 514.5 nm from an argon laser. The emission intensity I is plotted in arbitrary units as a function of the emission wavelength λ . The quantum wells have various thicknesses d and are separated by 50 nm of $\text{Al}_{0.5}\text{Ga}_{0.5}\text{As}$. The thinner wells were grown last, to prevent re-absorption of emitted radiation. Each well gives an emission peak due to the recombination of electrons and heavy holes with $n = 1$. The peak positions do not differ greatly from the theoretical values indicated in fig. 7. Broader peaks are observed with the thinner wells.

measuring the emission wavelength for small values of kT , e.g. at 4 K. A small correction is necessary because electrons and holes will first form excitons (mobile electron-hole 'pairs') before recombining completely. This effect reduces the photon energy by about 5-10 meV, depending on the thickness of the quantum well.

Fig. 8 shows the emitted spectrum of a sample containing five quantum wells of different sizes, which was excited by light from an argon laser at 514.5 nm. For each quantum well we only see the emission from the excitons associated with the $n = 1$ states. The polished surface of the substrates we used was at an angle of 6 degrees to the (001) plane of the crystal. The surface is 'step-like' during growth, with monolayer steps of 'height' 0.283 nm and about 2.8 nm apart. Addition of one GaAs molecule at each step site increases the mean grown thickness by 0.028 nm. In the plane of the quantum well, the excitons extend over about 30 nm, i.e. an order of magnitude larger than the distance between the steps, and the excitons will therefore not 'see' the steps but only the mean thickness of the GaAs layer, which varies in units of not more than 0.028 nm.

The width of the peaks is determined by the amount of (unintentionally incorporated) impurities in the epitaxial layer (5×10^{14} to 10^{15} cm^{-3}) and also by the amount by which the thickness of the quantum well varies over the light spot whose photoluminescence is measured. The tail of each peak on the low-energy side can be attributed to band-bending and to localized states due to impurities and thickness variations with a magnitude of a few tens of nm in the plane of the quantum well [14].

The width of the peaks is not related to the actual shape of the quantum wells, insofar as this is the same everywhere in the light spot, and does not therefore give information about the transition width of the composition profile at the interfaces. This transition will be fairly gradual because of the time necessary to change the gas composition over the wafer, and possibly because of adsorption or desorption at the walls of the reactor and tubing. Another possible source of gradual transition at the interface could be interdiffusion during growth, which may be enhanced when the interface is still close to the growing surface.

It has been shown [15] that the contribution of interdiffusion without surface-proximity enhancement is less than 0.1 nm for normal growth conditions. Since the time required for changing the gas concentration without adsorption or desorption is of the order of 0.1 s, its contribution will also be less than 0.1 nm. The two remaining mechanisms, adsorption/desorption and surface-proximity enhanced diffusion, are

most likely to provoke exponential tails at each compositional transition.

The effect of gradual compositional changes on the energy-level diagram is illustrated in *fig. 9*. Here the conduction-band edge in the well is lowered so little that the minimum is higher than the $n = 1$ electron level in pure GaAs. It is clear that the $n = 1$ electron level will shift to a higher energy and the $n = 1$ heavy-hole level to a lower energy. Thus the energy difference between the $n = 1$ states will be larger and the photoluminescence peak will be shifted to shorter wavelengths. For a given transition width, this shift is relatively larger for narrower quantum wells. This made it possible to assess the transition width by comparison of the relative positions of the photoluminescence peaks of a multiple quantum-well sample with the theoretical values for perfectly rectangular wells. The required exact knowledge of the ratios of the physical magnitude of the quantum wells was obtained by ensuring that the growth rate remained constant, and that there was no perturbation in the growth rate due to flow transients occurring when the gas flows are switched. The flow through the sources was stabilized by insertion of capillary tubes at the exit of the bubbling vessels.

The usefulness of these capillary tubes may be illustrated by the fact that, in our reactor, 50 cm³ of saturated trimethyl-gallium vapour in the bubbler is sufficient to grow 500 nm of GaAs. An accuracy of 0.1 nm in the thickness of a quantum well means that we have to be able to displace a smallest unit of just 0.02% of this volume (0.01 cm³) into the reactor chamber at the required instant. The result was checked by growing samples with superlattices consisting of alternating layers of GaAs and Al_xGa_{1-x}As. The ripple wavelengths in the depth profiles obtained by secondary-ion mass spectrometry (SIMS)^[16] agree well with the periods in the superlattices.

From the observed relative shifts of the photoluminescence peaks in the spectra of multiple quantum wells that were subsequently grown, we concluded, assuming an exponential decay profile of the Al content at the interfaces, that the transition width is no larger than about 0.5 nm for all interfaces^[17].

By forced depletion of the gas stream in the horizontal reactor chamber, a replica of a structure grown

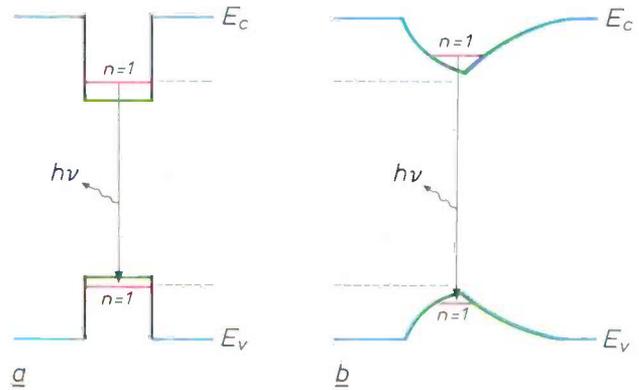


Fig. 9. The effect of gradual changes in composition on the energy-level diagram of a GaAs quantum well. With abrupt changes ideal interfaces are obtained, resulting in a rectangular well (*a*). If the composition changes slowly near the interfaces, the well becomes less deep (*b*). The edge of the conduction band E_c and the edge of the valence band E_v are modulated in such a way that the minimum of E_c becomes higher than the lowest conduction level of a rectangular well, whereas the maximum of E_v becomes lower than the highest valence level. The energy difference between the $n = 1$ states will therefore increase, so that the luminescence occurs at shorter wavelengths.

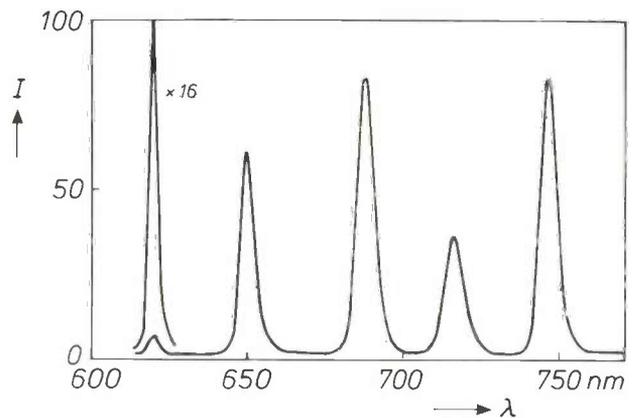


Fig. 10. Photoluminescence spectrum of a sample with five GaAs quantum wells, excited at 4 K with an argon laser emitting at 514.5 nm. The thinnest well, with an estimated thickness of only about 0.6 nm, has a weak emission peak at a very short wavelength (620 nm).

upstream can be grown downstream with layers less than a third of the thickness. In this way it was possible to show that even a 0.6-nm GaAs layer embedded in Al_{0.35}Ga_{0.65}As acts as a quantum well and gives a well-defined photoluminescence peak^[14], in the absence of significant interdiffusion of the layers during growth. The spectrum of a multilayer structure containing such a quantum well is shown in *fig. 10*. The emission wavelength of this well is very short, 620 nm.

Application in lasers

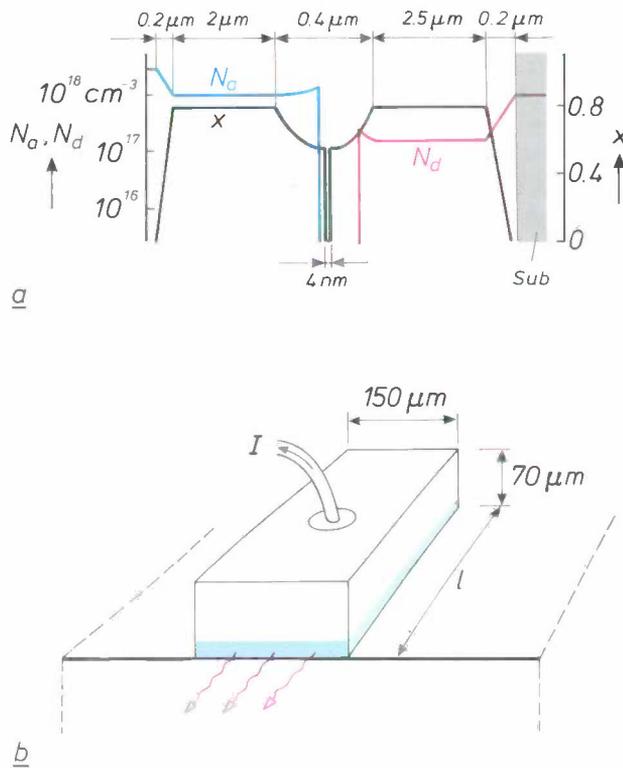
The use of one or more quantum wells as the active layer permits the manufacture of separate confinement semiconductor lasers with shorter emission wavelengths, lower threshold currents, higher output power and a longer lifetime than conventional semi-

[14] P. M. Frijlink, J. P. André and J. L. Gentner, Very narrow interface multilayer III-V heterostructures by organometallic vapor phase epitaxy, *J. Cryst. Growth* **70**, 435-443, 1984.

[15] R. M. Fleming, D. B. McWhan, A. C. Gossard, W. Wiegmann and R. A. Logan, X-ray diffraction study of interdiffusion and growth in (GaAs)_n(AlAs)_m multilayers, *J. Appl. Phys.* **51**, 357-363, 1980.

[16] See for example H. H. Brongersma, F. Meijer and H. W. Werner, Surface analysis, methods of studying the outer atomic layers of solids, *Philips Tech. Rev.* **34**, 357-369, 1974.

[17] P. M. Frijlink and J. Maluenda, MOVPE growth of Ga_{1-x}Al_xAs-GaAs quantum well heterostructures, *Jap. J. Appl. Phys.* **21**, L574-L576, 1982.



conductor lasers^[18]. Fig. 11a shows the compositional and doping profiles of the epitaxial layers for such a laser, which consists of a single central GaAs quantum well that collects and confines electrons and holes coming from n- and p- $\text{Al}_x\text{Ga}_{1-x}\text{As}$ respectively. The well is surrounded by an almost parabolic aluminium concentration profile (with x increasing from 0.5 to 0.8) of 0.2 μm on both sides, which confines the light by the gradient in the refractive index. Cladding layers of $\text{Al}_{0.8}\text{Ga}_{0.2}\text{As}$ with a low refractive index prevent coupling between the light in the central guide and the GaAs in the substrate and the top layer. To obtain good electrical contacts, a heavily n-doped GaAs substrate and a thin heavily p-doped GaAs top layer are used, and the transitions from $\text{Al}_{0.8}\text{Ga}_{0.2}\text{As}$ to GaAs at the extremities occur gradually within 0.2 μm .

After the deposition, the substrate is reduced in thickness by grinding and polishing the back, and both sides are metallized for electrical contacts.

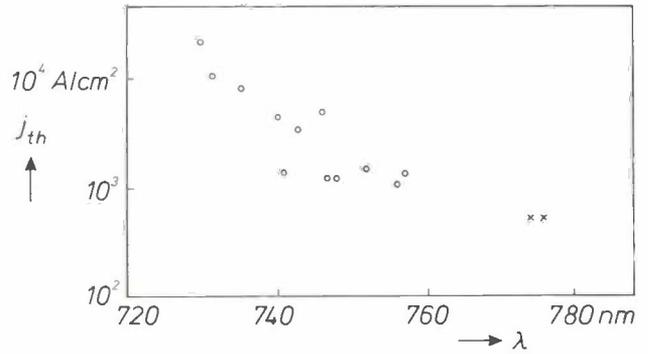


Fig. 12. Threshold current density j_{th} plotted against laser emission wavelength λ at room temperature for four-well lasers obtained from the same wafer (circles) and for two single-quantum-well lasers (crosses). The pulse duration was 2 μs with a repetition rate of 500 Hz. The shorter the emission wavelength, the higher the threshold for laser operation. No laser operation was detected below 730 nm.

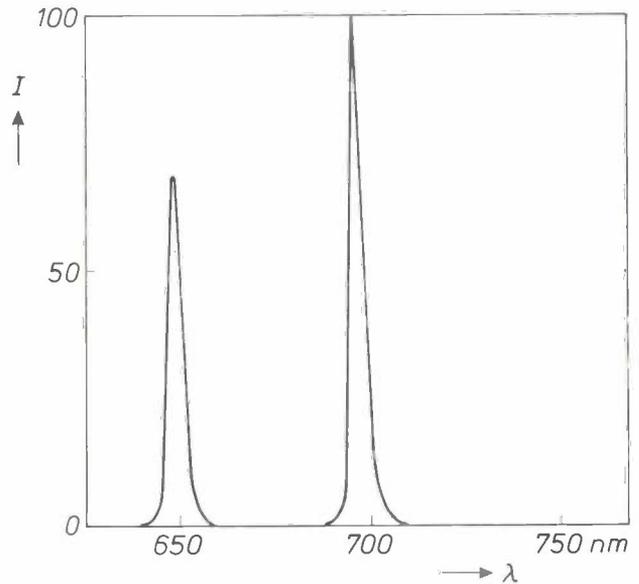


Fig. 13. Photoluminescence spectrum at 4 K on excitation with 514.5-nm argon-laser radiation of a sample containing two $\text{Al}_{0.2}\text{Ga}_{0.8}\text{As}$ quantum wells of thickness 3 and 8 nm and $\text{Al}_{0.5}\text{Ga}_{0.5}\text{As}$ barriers. The combination of a narrow well and Al incorporation gives a very short emission wavelength.

Cleaving and sawing provides individual laser chips, typically 500 by 150 μm , with the cleaved planes at the extremities acting as mirrors for obtaining optical resonance. The laser crystals are soldered, p-side down, to an indiated copper block (fig. 11b), which acts as a heat sink.

The threshold current for laser operation increases as the quantum well becomes thinner. In view of the thickness dependence of the emission wavelength (fig. 7) this implies that lasers with a shorter emission wavelength have a higher threshold current. Fig. 12 shows the threshold current density as a function of the emission wavelength for quantum-well lasers with different well thicknesses^[14]. No laser action was obtained for wavelengths shorter than 730 nm, prob-

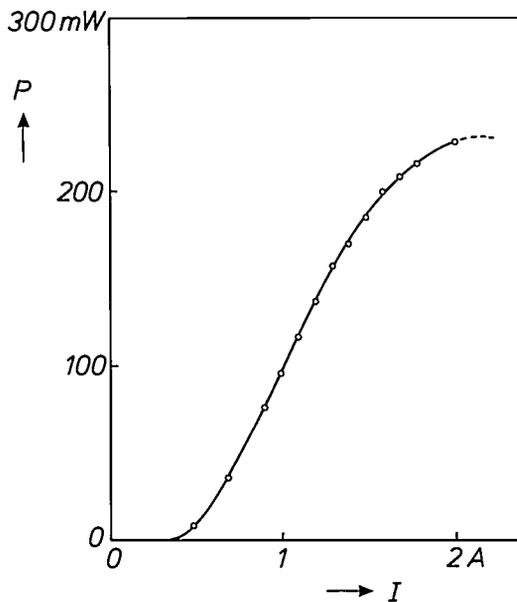


Fig. 14. Power P emitted per facet as a function of current I for a quantum-well laser of the type shown in fig. 11b with $l = 500 \mu\text{m}$, operating at room temperature and emitting at 775 nm. The threshold current is 0.50 A, corresponding to a current density of 660 A/cm². A maximum power of 230 mW is obtained.

ably because carrier capture becomes too inefficient for thinner wells.

The answer to this problem may be to use a large number of quantum wells, or $\text{Al}_x\text{Ga}_{1-x}\text{As}$ instead of GaAs for the quantum well, which can then be thicker for the same emission wavelength. For a given thickness, the incorporation of Al in the active layer shortens the emission wavelength. Fig. 13 shows the photoluminescence of 3-nm and 8-nm quantum wells with 20% of aluminium in the active layer. The emission wavelengths are considerably shorter than for pure GaAs wells of comparable thickness.

A laser structure that was more extensively studied, with a single quantum well and an emission at 775 nm, had a threshold current density of 530 A/cm² in the pulsed mode (2- μs pulses with a repetition rate of 500 Hz) and 660 A/cm² in CW operation. In fig. 14 the CW light output at room temperature is plotted against the input current. The maximum available light output power per facet was about 230 mW. The cavity length was 500 μm . Still lower threshold current-density values, down to 300 A/cm² for a cavity length of 440 μm , were obtained using 15-nm quantum wells, emitting at 860 nm.

Modulation-doped heterostructures

Modulation-doped $\text{Al}_x\text{Ga}_{1-x}\text{As}/\text{GaAs}$ heterostructures have become a subject of intensive study in many laboratories. In these structures, GaAs is undoped whereas $\text{Al}_x\text{Ga}_{1-x}\text{As}$ is n-doped with a typical donor

concentration of between 10^{17} and $5 \times 10^{18} \text{ cm}^{-3}$, except for a thin undoped layer next to the heterojunction, called a spacer layer. Owing to the difference in electron affinity and band gap between GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$, there is a step in the conduction and valence bands at the heterojunction; see fig. 15. Induced electron transfer creates an electron-accumulation region on the GaAs side of the heterojunction and an electron-depletion region on the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ side. The potential well in the accumulation region gives quantization of the electron states in the direction perpendicular to the heterojunction. Since movement parallel to the interface remains free, a two-dimensional electron gas is formed.

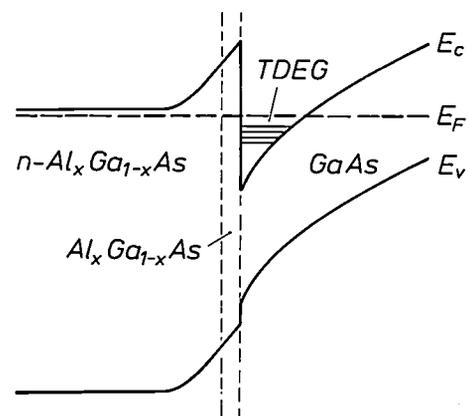


Fig. 15. Schematic energy-band structure of a modulation-doped heterostructure. E_v top of valence band; E_F Fermi level; E_c bottom of conduction band. The structure consists of an n-doped $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layer and an undoped GaAs layer, separated by a thin undoped $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layer or spacer. In the potential well on the GaAs side near the heterojunction, conduction electrons accumulate, forming a two-dimensional electron gas (TDEG). These electrons are spatially separated from the parent donor atoms in the n-doped layer.

The spatial separation between the accumulated electrons and the donor impurities greatly decreases the rate of impurity scattering. This results in a high electron mobility parallel to the heterojunction, which is increased at low temperatures as there is less optical phonon scattering. The combination of high electron density with high electron mobility in the accumulation region makes these structures very attractive for application in high-frequency devices.

An example of the modulation-doped heterostructures we have investigated^[19] is shown in fig. 16.

[18] W. T. Tsang, Extremely low threshold (AlGa)As modified multiquantum well heterostructure lasers grown by molecular-beam epitaxy, *Appl. Phys. Lett.* **39**, 786-788, 1981; D. R. Scifres, C. Lindström, R. D. Burnham, W. Streifer and T. L. Paoli, Phase-locked (GaAl)As laser diode emitting 2.6 W CW from a single mirror, *Electron. Lett.* **19**, 169-171, 1983.

[19] Investigations on modulation-doped heterostructures have also been described by J. Maluenda and P. M. Frijlink in *Jap. J. Appl. Phys.* **22**, L127-L129, 1983, and in *J. Vac. Sci. & Technol. B* **1**, 334-337, 1983.

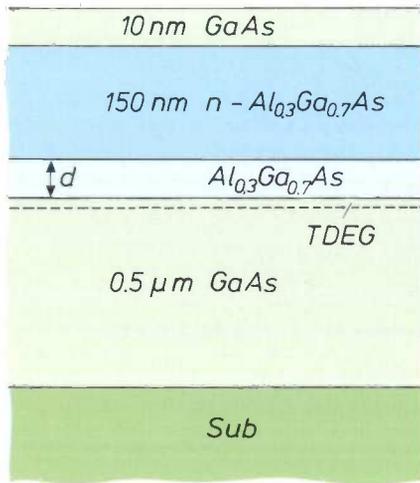


Fig. 16. Cross-section of a typical modulation-doped heterostructure on a Cr-doped semi-insulating GaAs substrate. Below the spacer of thickness d a two-dimensional electron gas (TDEG) is formed. The $n\text{-Al}_{0.3}\text{Ga}_{0.7}\text{As}$ layer contains $7.5 \times 10^{17} \text{ cm}^{-3}$ of Si atoms.

On a semi-insulating Cr-doped GaAs substrate there is a $0.5\text{-}\mu\text{m}$ non-doped GaAs layer and a $0.15\text{-}\mu\text{m}$ $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ layer doped with $7.5 \times 10^{17} \text{ at. Si/cm}^3$, with a thin non-doped $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ spacer layer of variable thickness. A 10-nm cap layer of GaAs facilitates the formation of ohmic contacts. Hall measurements were made on these structures: orthogonal external electric and magnetic fields induce a voltage across the sample at right angles to both fields, and this voltage was determined. The reciprocal of this voltage is proportional to the free carrier concentration and the ratio of this to the electrical resistivity gives the electron mobility.

Results of the Hall measurements on these structures are shown in *fig. 17*. The sheet electron density and the electron mobility determined at 77 K are given as a function of the spacer thickness. The density decreases with larger spacer thickness, owing to a decreasing accumulation effect. The electron mobility has a maximum of $8 \times 10^4 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ at a thickness of 12 nm , where the electron density is $8.3 \times 10^{11} \text{ cm}^{-2}$.

The electron mobility in this structure at 77 K is determined by the scattering at residual impurities in the direct vicinity of the electrons and at intentionally incorporated donors with Coulomb fields extending into the accumulation region. A screening effect also has to be considered. Electrons tend to stay longer near positive scattering centres than elsewhere. These centres are therefore surrounded by a net negative charge, which partly screens their fields and therefore reduces the scattering. If the spacer is very thin, most of the scattering occurs at intentionally incorporated donors. With a larger thickness this scattering diminishes owing to the increasing separation between

electrons and donors. This leads to a higher mobility, although the screening effect becomes weaker because of the lower electron density. With a thick spacer, scattering occurs mainly at residual impurities. The only result of an increase in thickness is then a lower electron density, so that there is less screening and hence a lower mobility.

The dependence of the electron mobility on the temperature is illustrated in *fig. 18* for one of our samples containing a 15-nm undoped spacer layer and a 45-nm $n\text{-doped}$ layer of $\text{Al}_{0.35}\text{Ga}_{0.65}\text{As}$, with a 30-nm GaAs cap layer^{[14][20]}. The sheet electron density was found to be $5 \times 10^{11} \text{ cm}^{-2}$. At room temperature a mobility of $7.4 \times 10^3 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ was measured. At lower temperatures, the decrease of optical

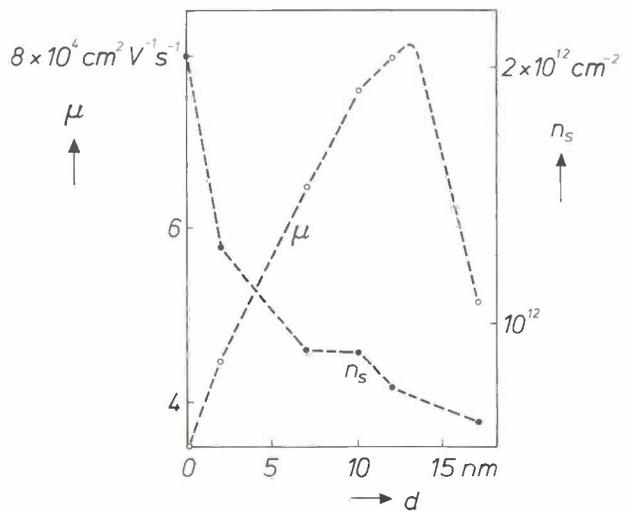


Fig. 17. Effect of spacer thickness d on the electron mobility μ at 77 K and the sheet electron density n_s , determined by Hall measurements on the modulation-doped heterostructures of *fig. 16*^[19]. With increasing spacer thickness the sheet electron density decreases monotonically, whereas there is a maximum in the electron mobility.

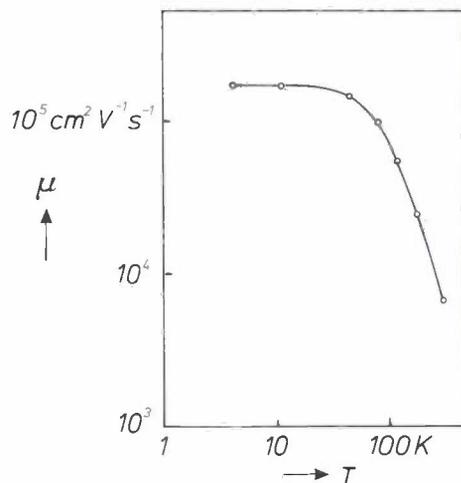


Fig. 18. Decrease in the electron mobility μ with temperature T for a modulation-doped heterostructure based on $\text{Al}_{0.35}\text{Ga}_{0.65}\text{As}$ and GaAs with a spacer thickness of 15 nm .

phonon scattering leads to markedly higher mobilities: $1.21 \times 10^5 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ at 77 K and $1.73 \times 10^5 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ at 4 K. The highest mobility obtained for one of our best samples was $2.7 \times 10^5 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ at 4 K.

Application in transistors

The idea of a modulation-doped heterostructure has been used in the fabrication of a high-electron-mobility transistor (HEMT), a special type of field-effect transistor. Here the electrons are transferred between the source and drain electrodes via a very thin channel near the substrate surface. The current conductance between source and drain is modulated by the transverse electric field between the gate and the substrate. A characteristic feature of a HEMT device is that the conductance takes place in the undoped GaAs layer in the accumulation region near the heterojunction.

The technology used in the fabrication of HEMT devices is based on planar processes. Here again MO-VPE has been found to be an appropriate method for depositing the various layers. Device isolation is obtained by boron implantation. The ohmic drain and source contacts are formed by evaporated alloys of gold and germanium. Several device geometries have been considered, differing in the length and width of the gate and in the source-gate and drain-gate spacings.

The gain that can be obtained is proportional to the transconductance, i.e. the change in the drain-source current induced by a given change in the gate voltage. The transconductance can be derived experimentally by measuring the drain-source current as a function of drain-source voltage at different values of the gate voltage. In *fig. 19* the curves measured at 300 and 77 K are given for a non-optimized device in which the Al gate had a length of $100 \mu\text{m}$ and a width of $1000 \mu\text{m}$ and the source-gate and drain-gate spacings were $1 \mu\text{m}$ [19]. This particular device was in fact the first HEMT made with MO-VPE (instead of MBE), thus demonstrating the feasibility of this technique. The transconductance, determined by the distance between the curves for different gate voltages, increases markedly on going from 300 K to 77 K. This clearly demonstrates the higher electron mobility at lower temperatures.

The best values obtained so far for the transconductance are 250 S/m at 300 K and 450 S/m at 77 K. These values were measured for a HEMT device with a gate length of $1.1 \mu\text{m}$, a gate width of $240 \mu\text{m}$ and a source-drain spacing of $4 \mu\text{m}$. The test versions of the HEMT devices were found to have a good high-frequency performance, mainly because of the modulated doping and the small gate length. Microwave

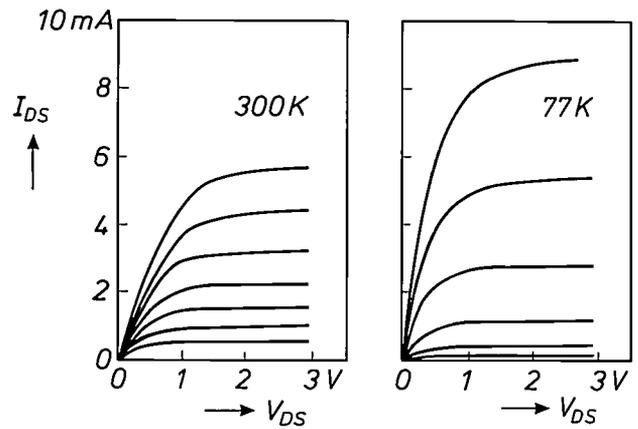


Fig. 19. Characteristic curves of drain-source current I_{DS} plotted against drain-source voltage V_{DS} for a HEMT device at 300 and 77 K, with the gate voltage varying in steps of 200 mV. The transconductance, derived from the spacings between the curves, increases threefold from 300 to 77 K.

measurements indicate a noise figure of 2.1 dB with an associated gain of 5.1 dB at 12 GHz and a cut-off frequency of 30 GHz.

Quantized Hall effect

An interesting effect associated with the presence of a two-dimensional electron gas in modulation-doped heterostructures is the 'quantized Hall effect'. This can be observed when the Hall voltage is measured at very low temperatures ($\leq 4 \text{ K}$) as a function of the

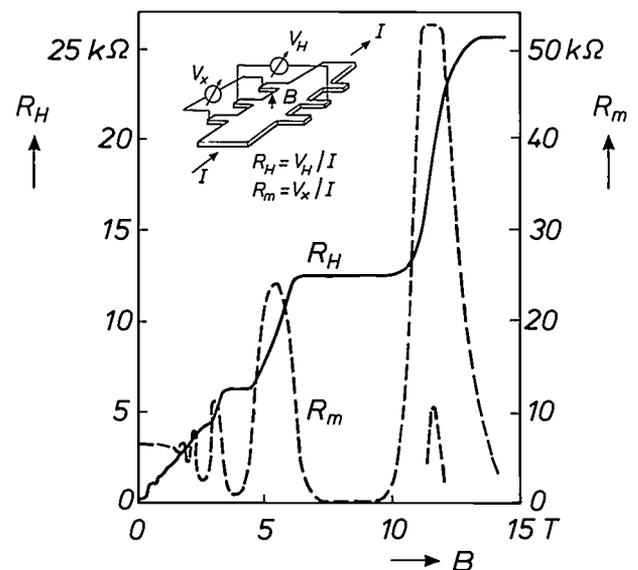


Fig. 20. Quantized Hall effect for a modulation-doped heterostructure based on $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ and GaAs, with a two-dimensional electron gas. The Hall resistance R_H (V_H/I) and magnetoresistance R_m (V_x/I), measured at 4 K, are given as a function of the magnetic flux density B (in the z -direction). At certain values of the flux density there is a plateau in the curve for R_H and there is a minimum for R_m .

[20] J. P. André, A. Brière, M. Rocchi and M. Riet, Growth of (Al,Ga)As/GaAs heterostructures for HEMT devices, *J. Cryst. Growth* 68, 445-449, 1984.

strong magnetic field. The Hall resistance R_H — the ratio between the measured Hall voltage and the electric current — does not increase linearly with the magnetic field but has a number of plateaus^[21]. The resistance at these plateaus appears to be independent of the electron mobility, sample geometry and impurity concentration and is given to an accuracy of 1 in 10^8 by:

$$R_H = h/e^2n, \quad (10)$$

where h is Planck's constant, e the electronic charge and n an integer^[22]. A plateau in the Hall resistance is associated with a minimum in the magnetoresistance.

The quantized Hall effect has been observed for several modulation-doped heterostructures grown in our Laboratories^[23]. An example of the results is given in *fig. 20*, where the Hall resistance and magnetoresistance measured at 4 K are plotted against the magnetic field. The sample had a 7-nm spacer of undoped $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ between 500-nm undoped GaAs and 150-nm Si-doped $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$. The value derived for the sheet carrier density in the accumulation region was $4 \times 10^{11} \text{ cm}^{-2}$ and the electron mobility measured at 77 K was $6 \times 10^4 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$. The plateaus in the Hall resistance and the close correspondence with the variation in the magnetoresistance are clearly visible.

[21] The quantized Hall effect was first observed in silicon-MOS field-effect transistors, see K. von Klitzing, G. Dorda and M. Pepper, *Phys. Rev. Lett.* **45**, 494-497, 1980.

Shortly afterwards this effect was also observed in $\text{Al}_x\text{Ga}_{1-x}\text{As}/\text{GaAs}$ heterostructures, see D. C. Tsui and A. C. Gossard, *Appl. Phys. Lett.* **38**, 550-552, 1981.

[22] With strong pulsed magnetic fields plateaus also occur at fractions of n , see D. C. Tsui, H. L. Stormer and A. C. Gossard, *Phys. Rev. Lett.* **48**, 1559-1562, 1982, and J. Wolter, this issue, pp. 111-117.

[23] The measurements of the quantized Hall effect were made by colleagues at Philips Research Laboratories, Eindhoven, and also at the Catholic University of Nijmegen. See for example R. E. Horstman, E. J. van den Broek, J. Wolter, R. W. van der Heijden, G. L. J. A. Rikken, H. Sigg, P. M. Frijlink, J. Maluenda and J. Hallais, *Solid State Commun.* **50**, 753-756, 1984;

R. Woltjer, J. Mooren, J. Wolter and J. P. André, *Solid State Commun.* **53**, 331-333, 1985.

Summary. Vapour-phase epitaxy with metal-organic reactants (MO-VPE) is a good method for growing multilayer structures of III-V semiconductor materials such as GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$. Precise control of the growth parameters can give interfaces that are abrupt on an atomic scale. The quality of the interfaces and the depth profiling have been assessed by spectroscopic ellipsometry and confirmed by photoluminescence experiments on quantum wells consisting of a thin GaAs layer between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers with a larger band gap. The emission wavelength of these wells can be varied in a controlled manner from 800 to 620 nm by decreasing the well width or incorporating aluminium in the well, or both. Laser operation has been obtained for quantum wells with emission wavelengths down to 730 nm. In modulation-doped heterostructures with GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ a two-dimensional electron gas is formed on the GaAs side near the heterojunction. The spatial separation of carriers and donor impurities gives high electron mobilities, particularly at low temperatures. This favours the application of such structures for obtaining transistors with a good high-frequency performance. At low temperatures the presence of a two-dimensional electron gas is responsible for the 'quantized Hall effect'.

Metal-organic vapour-phase epitaxy with a novel reactor and characterization of multilayer structures

M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft

Introduction

The previous article^[1] explains how multilayer structures with $\text{Al}_x\text{Ga}_{1-x}\text{As}$ can be grown on a GaAs substrate by metal-organic vapour-phase epitaxy (MO-VPE). The great importance of abrupt transitions in such structures prompted us to optimize the growth method to produce such transitions. A novel reactor was therefore developed that could be used to bring about very fast changes in the gas composition. Transmission-electron-microscope (TEM) investigations of the resulting structures have shown that the transition between layers of different composition can take place within approximately one monolayer. Investigations of the optical quality of the structures have shown that quantum wells consisting of a very thin GaAs layer between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers can have a high luminescence efficiency. We have studied the temperature dependence of the radiative recombination coefficient for various quantum wells; this quantity is important for optoelectronic applications of quantum-well lasers. The measured temperature dependence and the variation with well thickness showed reasonably good agreement with calculations based on a simple model.

In this article we first describe the new reactor and the growth procedure. We then discuss the characterization by TEM, with special attention to the method of obtaining a good contrast between layers of different composition. TEM micrographs of structures that have been grown are then shown. Finally the temperature dependence of the radiative recombination coefficient in quantum wells is compared with the value in the bulk material and the consequences for laser applications are discussed.

Growth of multilayer structures with abrupt transitions

MO-VPE can be used to produce a multilayer structure on a substrate by appropriately changing the composition of the vapour phase above the substrate surface during epitaxial growth^[1]. To obtain struc-

tures with abrupt transitions and very thin layers, good control of the gas composition is essential. This means that special precautions have to be taken with the design of the reactor and the processing conditions. A rapid change in the supply of the components does not necessarily produce the same rapid change near the

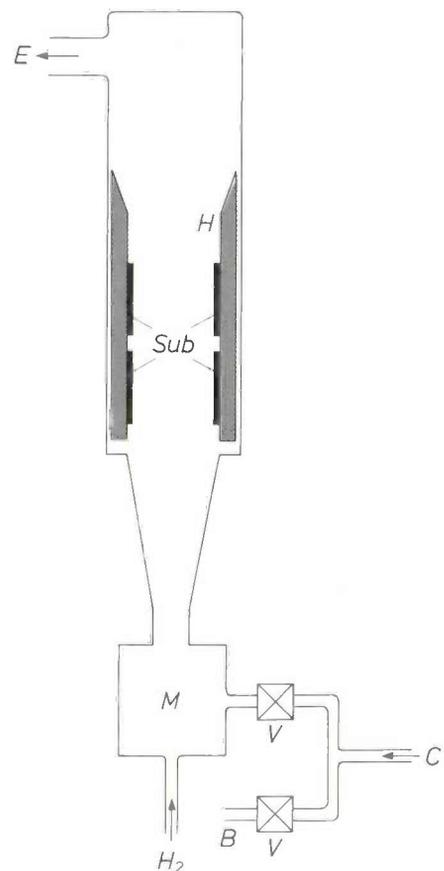


Fig. 1. Diagram of a novel reactor for MO-VPE of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers. The carrier gas (H_2) is supplied to the mixing chamber M at the bottom of the vertical reactor. A component C is admitted by a pneumatic valve V to the mixing chamber or switched by another valve to a bypass B . The substrates Sub are attached to the inner wall of a hollow susceptor H , which is heated by r.f. induction to a temperature of about 700°C . Thin films are deposited on the substrates from the gas mixture (carrier gas plus components). E exhaust.

Drs M. R. Leys, formerly with Philips Research Laboratories, Eindhoven, is now with the University of Lund, Sweden; Dr M. P. A. Vieggers and Dr G. W. 't Hooft are with Philips Research Laboratories, Eindhoven.

[1] P. M. Frijlink, J. P. André and M. Erman, Metal-organic vapour-phase epitaxy of multilayer structures with III-V semiconductors, this issue, pp. 118-132.

substrate surface. While the components are being transported through the reactor various recirculation effects can occur, with the result that the original fast change in the gas composition can be 'spread out' to give a gradual change in the composition in the actual structure. One important factor here is gas recirculation in the reactor; another is diffusion, especially from 'stagnant' volumes in the valves and in a boundary layer near the substrate surface. Our reactor has been designed to reduce this 'smoothing out' to a minimum [2].

The new reactor

Fig. 1 is a diagram of the new reactor used for producing multilayer structures with sharp transitions. *Fig. 2* shows part of the experimental arrangement used

chamber through a feed line or is switched to a bypass line. These lines have the same flow resistance, so that the transit times are the same and there are no pressure fluctuations when the valves are opened or closed. The inlet valves and the lines are purged constantly with hydrogen so that there are no stagnant volumes in the gas-feed system.

The geometry of the reactor is optimized for fast throughput of the gas mixture (carrier gas plus components). This means that the distances between inlet valves, mixing chamber and substrates are kept to a minimum. Recirculation of the gas mixture is largely eliminated by streamlining the input region, and by the vertical arrangement of the reactor. The gas mixture passes through a hollow susceptor with a rectangular internal cross-section; the substrates are attached

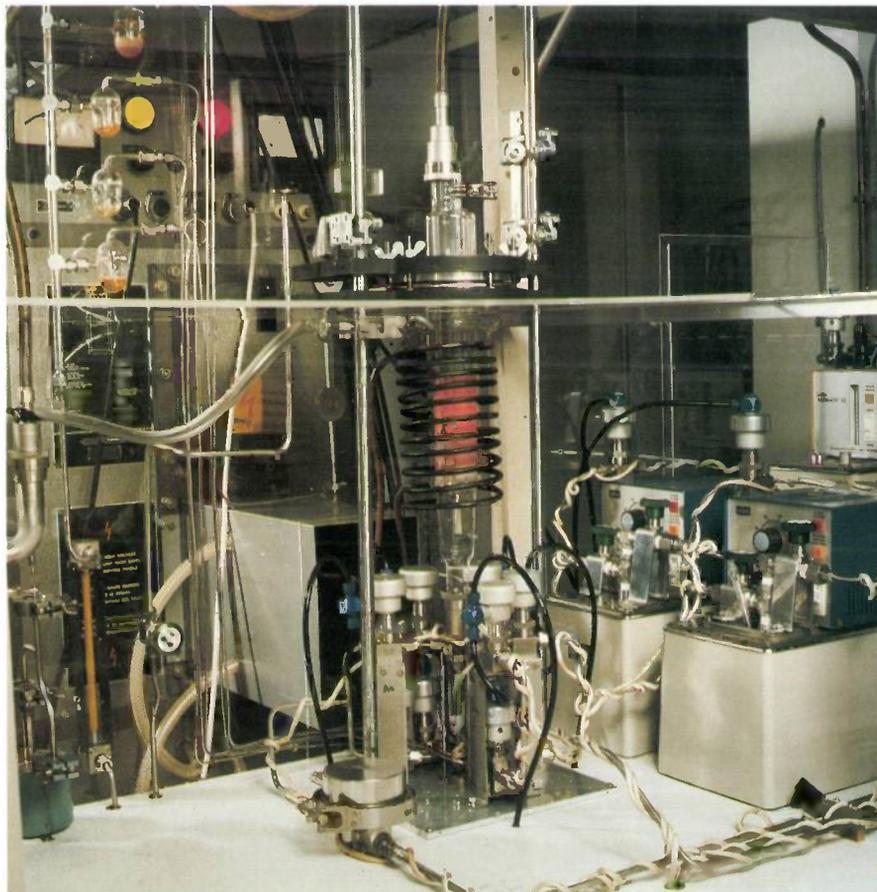


Fig. 2. Part of the experimental MO-VPE arrangement for growing layers of $\text{Al}_x\text{Ga}_{1-x}\text{As}$. The reactor with gas inlet at the bottom is shown at the centre, and the valves for admitting the various components are arranged around it radially. The substrate is heated to the required temperature in the reactor by r.f. induction in the susceptor. The thermostatic baths with bubbling vessels containing the metal-organic compounds are on the right.

for producing the structures described in this article. The carrier gas (hydrogen) and the reactive components enter the vertical reactor from below after passing through a mixing chamber. The different components are supplied to the mixing chamber by separate lines (five in total). Each component enters the mixing

to the inner wall of the susceptor. With this geometry the free convection — the tendency of hot (less dense) gases to rise, as in a chimney — helps rather than hinders the forced convection. This minimizes the effect of gas recirculation and reduces the thickness of the diffusion-boundary layer near the substrate surface.

The graphite susceptor contains two areas (each $35 \times 70 \text{ mm}^2$) for positioning the substrates. It is heated by r.f. induction. The hydrogen acting as carrier gas is purified by a palladium diffusor. The system for the growth of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ also contains sources of the reactive components trimethyl aluminium ($\text{Al}_2(\text{CH}_3)_6$), trimethyl gallium ($\text{Ga}(\text{CH}_3)_3$) and arsine (AsH_3), as well as the reactive components for n-type or p-type doping. A microcomputer controls the supply of the various components during growth and keeps the susceptor at the correct temperature.

Growth procedure

The multilayer structures are grown on the (001) surface of a GaAs substrate. The susceptor is kept at 700°C during the growth. The total gas flow through the reactor is 7.5 l/min and the total pressure in the reactor is 10^5 Pa (1 atm). The partial pressure of $\text{Ga}(\text{CH}_3)_3$ is 3 Pa, the partial pressure of AsH_3 is 75 Pa. Under these conditions the growth rate is 0.35 nm/s, corresponding to slightly more than one monolayer per second. The partial pressure of $\text{Al}_2(\text{CH}_3)_6$ can be varied from 0 to 12 Pa to grow $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers with different compositions. The partial pressure of $\text{Al}_2(\text{CH}_3)_6$ for the growth of AlAs is 3.6 Pa.

To demonstrate how flexible the growth process is and how well the layer thickness can be controlled, we made a complicated test structure. The growth pattern of the structure is shown in fig. 3. At the start we grew alternate layers of GaAs and $\text{Al}_{0.55}\text{Ga}_{0.45}\text{As}$ by continuous introduction of $\text{Ga}(\text{CH}_3)_3$ and AsH_3 while alternately switching the flow of $\text{Al}_2(\text{CH}_3)_6$ from the mixing chamber to the bypass line. The first GaAs layer (called the buffer layer) and the first $\text{Al}_{0.55}\text{Ga}_{0.45}\text{As}$ layer are relatively thick; the others are extremely thin (down to about 0.7 nm grown in two seconds). Next some layers of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ with ascending x were grown. Before growing each layer the component supply was briefly interrupted (for about 5 s) until the partial $\text{Al}_2(\text{CH}_3)_6$ pressure had reached the appropriate value. In the final part of the growth $\text{Ga}(\text{CH}_3)_3$ and $\text{Al}_2(\text{CH}_3)_6$ were supplied alternately, to produce a quasi-superlattice, consisting of alternate 3.5-nm layers of GaAs and AlAs. Finally, a layer of AlAs and a top layer of $\text{Al}_{0.55}\text{Ga}_{0.45}\text{As}$ were grown.

Characterization by transmission electron microscopy

It is difficult to obtain direct and unambiguous information about layer thickness and transition width in these multilayer structures. The previous article explains how luminescence spectra can be used for estimating the thicknesses of the layers and transition regions [1]. The model used, however, only applies for

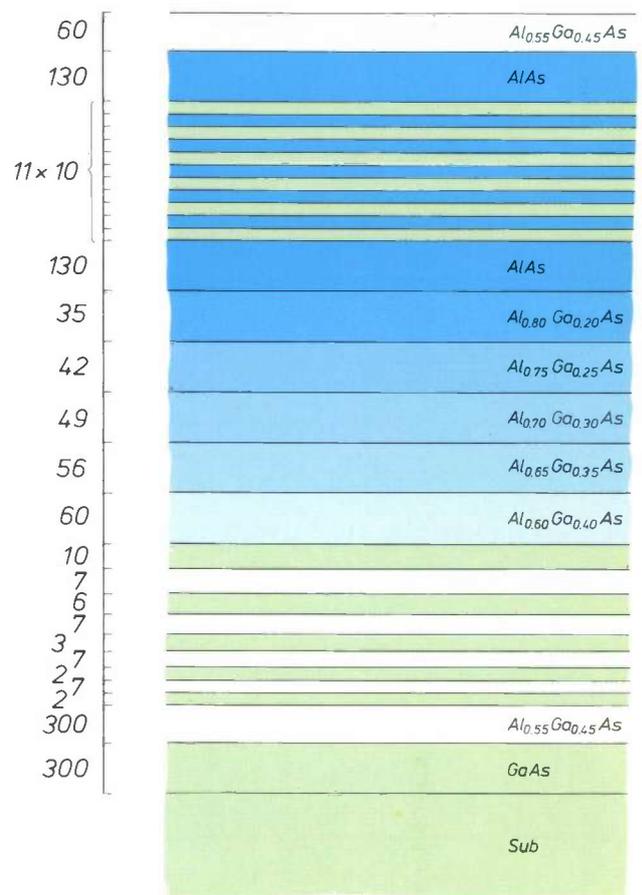


Fig. 3. Diagram of a test structure grown by MO-VPE, consisting of a large number of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers with different values of x . The growth time (in seconds) for each layer is plotted along the vertical axis. The layers and growth times are not shown to scale.

a direct-gap semiconductor ($\text{Al}_x\text{Ga}_{1-x}\text{As}$ with $x \leq 0.35$) with a layer thickness between 2 and 20 nm. Analytical techniques such as Auger electron spectroscopy and secondary-ion mass spectroscopy [3] are limited by the fact that the depth resolution is no better than about 2 nm. Also, since these methods only give the composition as a function of depth, details at interfaces cannot be observed.

The method we used for characterization in our investigations was transmission electron microscopy (TEM). In principle a resolution of 0.2 nm can be obtained with TEM, so that structures can be studied on an atomic scale. The method can also reveal details with atomic dimensions at interfaces.

The basic requirement for TEM is a specimen thin enough to be transparent to electrons. This means that a multilayer structure has to be studied in cross-

[2] M. R. Leys, C. van Opdorp, M. P. A. Vieggers and H. J. Talen-van der Mheen, Growth of multiple thin layer structures in the GaAs-AlAs system using a novel VPE reactor, *J. Crystal Growth* **68**, 431-436, 1984.

[3] See H. H. Brongersma, F. Meijer and H. W. Werner, Surface analysis, methods of studying the outer atomic layers of solids, *Philips Tech. Rev.* **34**, 357-369, 1974.

section, which requires a special technique for specimen preparation [4]. The TEM micrographs shown here were made with a Philips EM 420 ST microscope operated at 120 kV. It is possible to distinguish between layers of different composition in the structure because they diffract the electrons differently [6]. We shall now show how the composition of the layers affects the contrast in a TEM micrograph [6].

Contrast

In making a TEM micrograph of a single crystal a coherent electron beam is diffracted at the crystal planes in directions given by Bragg's law:

$$2d \sin \theta = \lambda, \quad (1)$$

where d is the distance between the planes, θ the angle between the incident beam and the planes (about 0.5°) and λ is the wavelength of the electrons. The diffracted beams are focused to form a diffraction pattern at the diaphragm in the back focal plane of the objective lens; see fig. 4. When the specimen is tilted so that only one set of planes exactly satisfies the Bragg relation, then apart from the transmitted beam only one other beam effectively reaches the back focal plane. An image can then be made through the aperture of a diaphragm at the back focal plane, either including the transmitted beam (bright field) or with the diffracted beam only (dark field).

Local deviations from the Bragg relation and local differences in the composition of the specimen produce a contrast (amplitude contrast) in the TEM micrograph. The multilayer structures considered here are always oriented in such a way that the Bragg condition is satisfied, so that the contrasts in the TEM photograph are a direct representation of the local differences of composition.

The amplitude of an electron beam that is diffracted from the (hkl) planes is determined by the structure factor F_{hkl} , defined as [6]:

$$F_{hkl} = \sum_i^n f_i \exp 2\pi j(hx_i + ky_i + lz_i), \quad (2)$$

where n is the number of atoms in the unit cell of the crystal structure, j the imaginary unit ($\sqrt{-1}$), f_i the scattering factor of atom i , and x_i , y_i and z_i are the coordinates of this atom, expressed in fractions of the dimensions of the unit cell. For the Al and Ga atoms in $\text{Al}_x\text{Ga}_{1-x}\text{As}$ the atomic scattering factors have to be weighted with the relative concentrations:

$$f_{\text{Al/Ga}} = x f_{\text{Al}} + (1 - x) f_{\text{Ga}}. \quad (3)$$

Since the scattering factor is proportional to the atomic number, f_{Al} is much smaller than f_{As} whereas f_{Ga} and f_{As} are not very different. The coordinates

(x_i, y_i, z_i) of the atoms in the unit cell can be derived from the crystal structure of $\text{Al}_x\text{Ga}_{1-x}\text{As}$. This corresponds to the zinc-blende structure of ZnS: the Al/Ga atoms form a cubic close-packed structure in which As atoms occupy half the number of tetrahedral interstices. Fig. 5 shows the unit cell of $\text{Al}_x\text{Ga}_{1-x}\text{As}$. The coordinates (x_i, y_i, z_i) of the atoms are:

$$\begin{aligned} (0,0,0), (0, \frac{1}{2}, \frac{1}{2}), (\frac{1}{2}, 0, \frac{1}{2}), (\frac{1}{2}, \frac{1}{2}, 0) & \text{ for Al/Ga} \\ (\frac{1}{4}, \frac{1}{4}, \frac{1}{4}), (\frac{1}{4}, \frac{3}{4}, \frac{3}{4}), (\frac{3}{4}, \frac{1}{4}, \frac{3}{4}), (\frac{3}{4}, \frac{3}{4}, \frac{1}{4}) & \text{ for As.} \end{aligned} \quad (4)$$

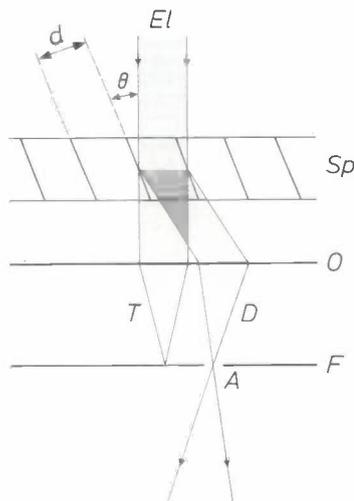


Fig. 4. Ray diagram for dark-field imaging with transmission electron microscopy (TEM). The single-crystal specimen Sp is oriented in such a way that the angle θ for the incident electron beam EI only satisfies the Bragg relation given in eq. (1) for one set of planes (with spacing d). The diffracted beam D and the transmitted beam T are focused by an objective lens O on to the back focal plane F , where there is a diaphragm. The aperture A in the diaphragm only passes the diffracted beam. Contrasts in the micrograph are due to local deviations from the Bragg relation and local differences in composition.

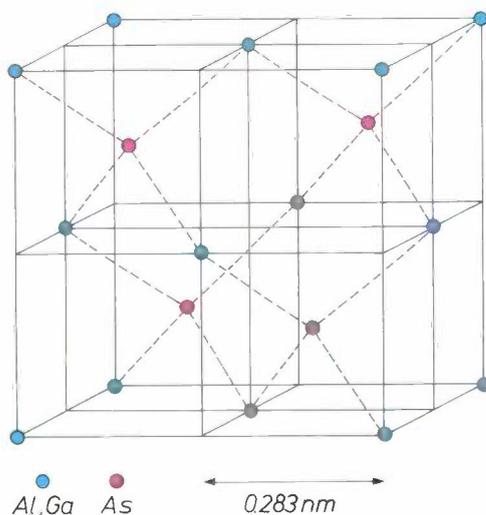


Fig. 5. Unit cell of $\text{Al}_x\text{Ga}_{1-x}\text{As}$. The Al/Ga atoms form a face-centred cubic lattice: they are situated at the corners of the unit cell and at the centre-points of the six faces. The As atoms fill half the number of the available tetrahedral interstices. The coordinates of the atoms in the unit cell are given in equation (4). The lattice constant is 0.566 nm.

On substituting f_i and (x_i, y_i, z_i) in eq. (2) for different values of h, k and l it is found that the structure factor depends most strongly on x for diffraction from (002) planes^[6]. In this case the structure factor is given by:

$$F_{002} = \sum_i^n f_i \exp(4\pi jz_i). \quad (5)$$

Substituting z_i from eq. (4), $f_{\text{Al/Ga}}$ from eq. (3) and f_{As} , we obtain:

$$F_{002} = 4xf_{\text{Al}} + 4(1-x)f_{\text{Ga}} - 4f_{\text{As}}. \quad (6)$$

With $f_{\text{As}} \approx f_{\text{Ga}}$ this becomes:

$$F_{002} \approx 4x(f_{\text{Al}} - f_{\text{Ga}}). \quad (7)$$

The composition dependence of the intensity I_{002} of the diffracted beam is then given approximately by:

$$I_{002} \propto F_{002}^2 \propto x^2. \quad (8)$$

This means that layers with a small x (for example GaAs, $x = 0$) will be dark compared with layers with a large x (for example AlAs, $x = 1$), if we make a dark-field image with I_{002} alone.

Another contribution to the contrast is connected with the number of electrons taking part in the imaging. In the presence of heavier atoms in the specimen, many electrons are scattered over such large angles

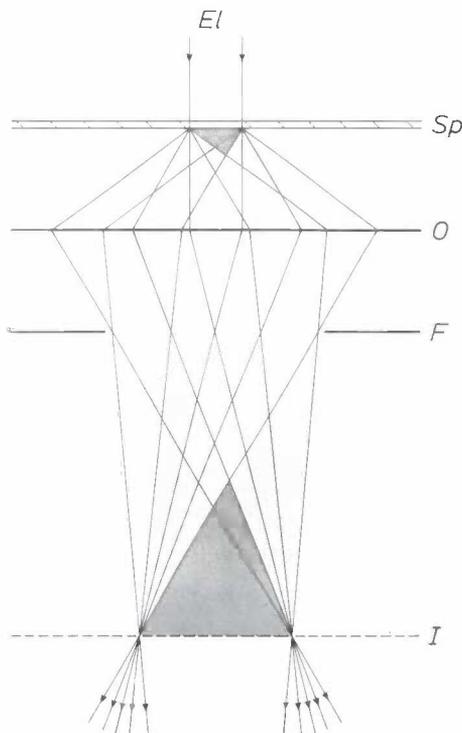


Fig. 6. Ray diagram for a high-resolution TEM micrograph. The electron beam EI is incident on the specimen Sp in the direction of (say) the [110] axis. The beams produced by diffraction from the crystal planes are focused by an objective lens O on to the back focal plane F . The diffraction pattern formed by the combination of a large number of diffracted beams results in an image of the crystal lattice at the image plane I .

that they fall outside the lens aperture. This effect is not so pronounced in specimens consisting of lighter atoms (they have less scattering power). Although this effect does not greatly enhance the contrast, it operates in the same sense as the amplitude contrast: areas with GaAs appear darker and those with AlAs appear brighter.

The resolution of the dark-field image is determined by the diameter of the aperture of the diaphragm in the back focal plane of the objective lens. This aperture cannot be made arbitrarily large: the maximum diameter is equal to the distance between the transmitted and the diffracted beams. This means that the resolution cannot be better than the spacing of the (002) planes; in GaAs this is 0.283 nm. In practice it is about 0.5 nm.

A higher resolution can be obtained by quite a different method, in which the electron beam is incident parallel to one of the crystal axes with a low index, e.g. the [110] axis. In this case the beam is diffracted from many crystal planes, all nearly parallel to the incident beam. An image of the crystal lattice is formed by combination of a large number of these diffracted beams, as shown in fig. 6. Calculation of the contrast is very complicated in this case because there are so many beams. Nevertheless, because of the difference in scattering between Al and Ga, the GaAs regions appear darker than the AlAs regions.

Results

Several TEM micrographs were made of the structure that was grown in the pattern shown in fig. 3. Fig. 7 shows a dark-field image made with I_{002} . The different layers are easily distinguished from each other. At the top of relatively thick AlAs layers a wavy interface can be seen; the 'waves' have an amplitude of about 2 nm and a wavelength of about 100 nm. This effect has not yet been completely understood. It is probable that the interface structure is due to surface instability during growth, caused by adsorption of impurity atoms, notably oxygen and carbon. These may accumulate at 'steps' already present on the growing crystal surface. Such impurities slow down

[4] M. P. A. Vieggers, A. F. de Jong and M. R. Leys, Characterization of structural features in thin layers of GaAs, Al(x)Ga(1-x)As and AlAs by means of structure factor imaging and high resolution electron microscopy, Spectrochim. Acta 40B, 835-845, 1985.

[6] P. M. Petroff, Transmission electron microscopy of interfaces in III-V compound semiconductors, J. Vac. Sci. & Technol. 14, 973-978, 1977.

[6] Good books on diffraction theory for TEM include: P. B. Hirsch, A. Howie, R. B. Nicholson, D. W. Pashley and M. J. Whelan, Electron microscopy of thin crystals, Butterworth, London 1965; G. Thomas and M. J. Goringe, Transmission electron microscopy of materials, Wiley, New York 1979.

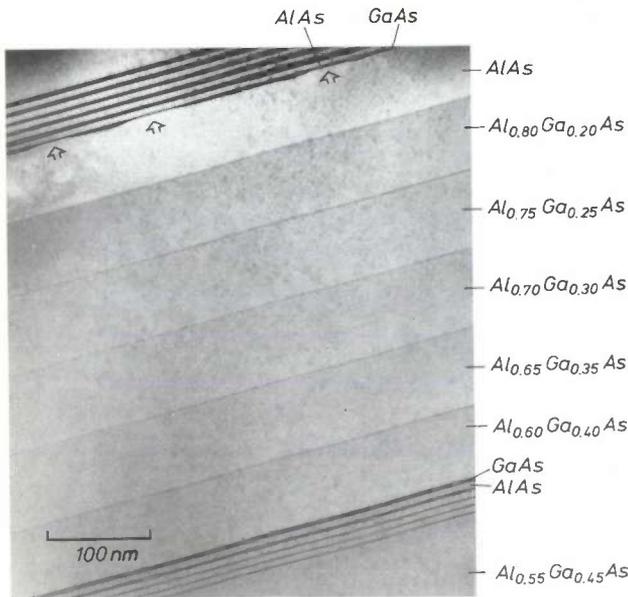


Fig. 7. TEM dark-field image, with (002) diffraction, of a multilayer structure of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ grown as in the diagram of fig. 3. Layers with low x (Al-poor) appear dark compared with layers with high x (Al-rich), as indicated by eq. (8). The arrows indicate undulations at the AlAs surface.

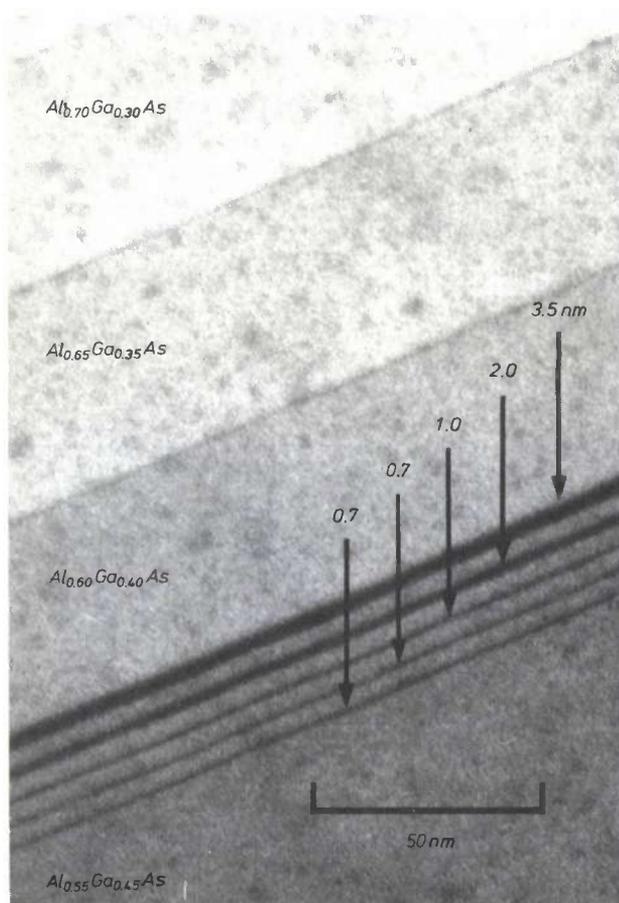


Fig. 8. TEM dark-field image, with I_{002} , of the bottom part of the multilayer structure of fig. 3. The GaAs layers (dark) have a thickness ranging from 0.7 to 3.5 nm. The $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers become brighter as x increases; an increase of 0.05 is clearly visible.

the advance of a step over the surface. During the growth of 'thick' layers (≥ 60 nm), with prolonged adsorption of impurities, this may have the effect of disturbing the two-dimensional growth required for a planar surface. Pure GaAs is much less sensitive to oxygen and carbon than AlAs, and will therefore maintain a planar interface much more easily. It is remarkable that the growth of a thin layer of 3.5 nm of GaAs is nevertheless sufficient to produce a smooth interface on a surface with undulations 2 nm high. It appears that the flatness of a surface is largely kinetically determined, e.g. by the surface mobility and the adsorption of impurities. The flatness is therefore not determined by the abruptness of the changes in composition, which depends on the technology of the MO-VPE reactor.

The lower part of the structure is shown in more detail in fig. 8. The stepwise changes of 5% in the aluminium content are clearly visible in the TEM micrograph as differences in grey tone. It can also be seen that the changes in composition take place within 0.5 nm, which is approximately equal to the attainable resolution and corresponds to a transition region of one or two monolayers. In considering the minimum thickness of the GaAs layers, we have to remember that only an integral number of monolayers can be deposited. This means that the thickness of the GaAs layers must always be a multiple of half the lattice constant (0.283 nm). The thinnest layer, grown in 2 s at an average rate of 0.35 nm/s will therefore have a thickness varying in steps between 0.57 and 0.85 nm.

By combining several beams (fig. 6) a more detailed image of the structure is obtained. Fig. 9 shows a TEM micrograph made with the electron beam incident along a [110] direction for an image of the quasi-superlattice of GaAs and AlAs (fig. 3). The structure in the micrograph bears a direct relation to the atomic order in the GaAs and AlAs layers. Each round 'spot' visible in the micrograph corresponds to a combination of a Ga (or Al) atom and an As atom; see fig. 10. The growth steps at the interfaces are clearly visible, especially at the transition from AlAs to GaAs. It is clear that the transition from GaAs to AlAs, and vice versa, takes place within approximately one monolayer (0.283 nm). The growth steps observed at the in-

[7] W. K. Burton, N. Cabrera and F. C. Frank, The growth of crystals and the equilibrium structure of their surfaces, *Phil. Trans. R. Soc. London A* **243**, 299-358, 1950/51.

[8] J. H. Neave, B. A. Joyce, P. J. Dobson and N. Norton, Dynamics of film growth of GaAs by MBE from RHEED observation, *Appl. Phys. A* **31**, 1-8, 1983.

See also the article by B. A. Joyce and C. T. Foxon, Molecular beam epitaxy of multilayer structures with GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$, this issue, pp. 143-153.

[9] G. W. 't Hooft, M. R. Leys and H. J. Talen-van der Mheen, Temperature dependence of the radiative recombination coefficient in GaAs-(Al,Ga)As quantum wells, *Superlattices & Microstructures* **1**, 307-310, 1985.

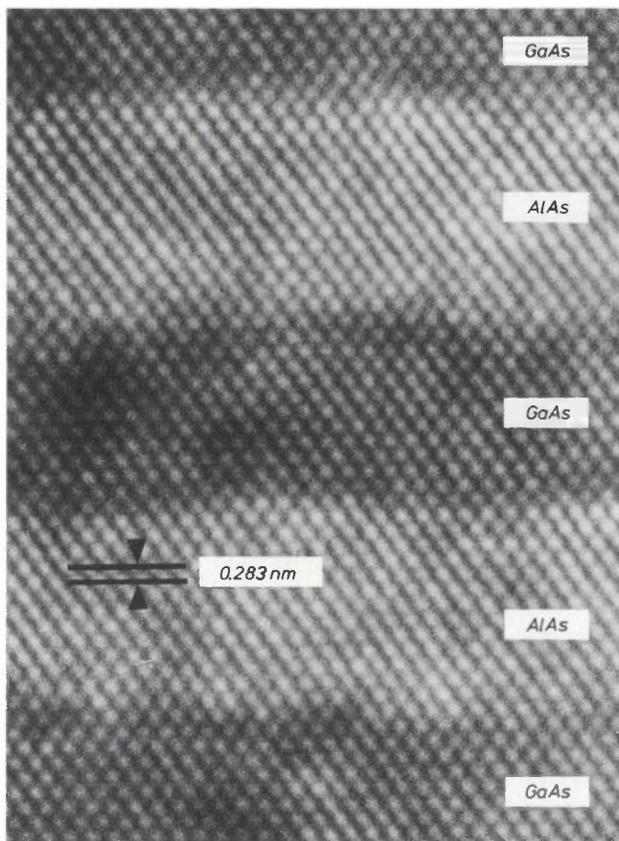


Fig. 9. TEM high-resolution micrograph, with the electron beam in a [110] direction, of part of the quasi-superlattice of GaAs and AlAs, grown as shown in fig. 3. The transition from GaAs to AlAs (and vice versa) takes place within about one monolayer. Growth steps are visible at the interfaces.

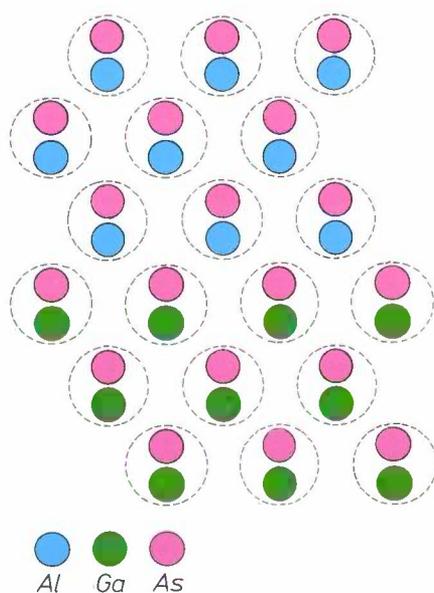


Fig. 10. Projection in the (110) plane of a structure with a GaAs layer and an AlAs layer. Dashed circles are drawn around the Ga (or Al) and As atoms that lie close together in this projection. These correspond to the round spots visible in the TEM image in fig. 9.

terfaces are due to the presence of surface islands during the (two-dimensional) growth. The occurrence of surface islands during epitaxial growth is known from crystal-growth theory^[7] and has also been demonstrated in the growth of GaAs with MBE^[8].

Temperature dependence of the radiative recombination coefficient in quantum wells

The reactor described in this article can be used for the fabrication of quantum-well lasers of the type described in the previous article^[1]. Under optimum growth conditions it is possible to obtain quantum wells that have a high value of the luminescence efficiency, defined as the number of emitted photons divided by the number of injected charge carriers. A high efficiency means that the probability of a spontaneous radiative electron-hole recombination, given by the coefficient B , is large compared with the probability of non-radiative transitions. The dependence of these probabilities on temperature is of practical importance for optoelectronic devices. In general, the coefficient B decreases with increasing temperature whereas the probability of non-radiative transitions increases. This leads to a fall in the luminescence efficiency and increases the threshold current for laser operation at higher temperatures.

Theoretically, B is less strongly dependent on the temperature in a quantum well than in the bulk material; this has an advantageous effect on the luminescence efficiency and on the threshold current for laser operation when quantum wells are incorporated as the active layer. We have studied the temperature dependence of B experimentally by measuring the luminescence decay times at different temperatures^[9]. The measurements were performed on bulk GaAs and on quantum wells, which have such a high efficiency at these temperatures that, as a first approximation, the effect of non-radiative transitions can be neglected. Before looking at the results, we shall first show, in simple theoretical terms, how the temperature dependence of B in a quantum well differs from that in the bulk material^[9].

Theory

The radiative recombination coefficient B is defined as:

$$B = L/np, \quad (9)$$

where L is the photon generation rate (number of photons per unit time and per unit volume) and n and p are the concentrations of electrons and holes in the conduction and valence bands, respectively. To establish the temperature dependence we assume that holes

of only one type are present in the valence band. The values of n , p and L are determined by the densities of states ρ_e and ρ_h and by their occupation probabilities f_e and f_h in the conduction and valence bands ^[10]:

$$n = \int_{E_c}^{\infty} \rho_e f_e dE_e, \quad (10)$$

$$p = \int_{-\infty}^{E_v} \rho_h (1 - f_h) dE_h, \quad (11)$$

$$L \propto \int_{E_g}^{\infty} \rho_{\text{red}} f_e (1 - f_h) d(E_e - E_h), \quad (12)$$

where E_e is the energy of an electron in the conduction band with minimum energy E_c and E_h is the energy of a hole in the valence band with maximum energy E_v ; see fig. 11. In eq. (12) E_g is the band gap ($E_c - E_v$) and ρ_{red} is the reduced density of states for electron-hole pairs:

$$\rho_{\text{red}} = \frac{1}{2}(\rho_e^{-1} + \rho_h^{-1})^{-1}. \quad (13)$$

The occupation probabilities f_e and f_h depend on temperature in accordance with Fermi-Dirac statis-

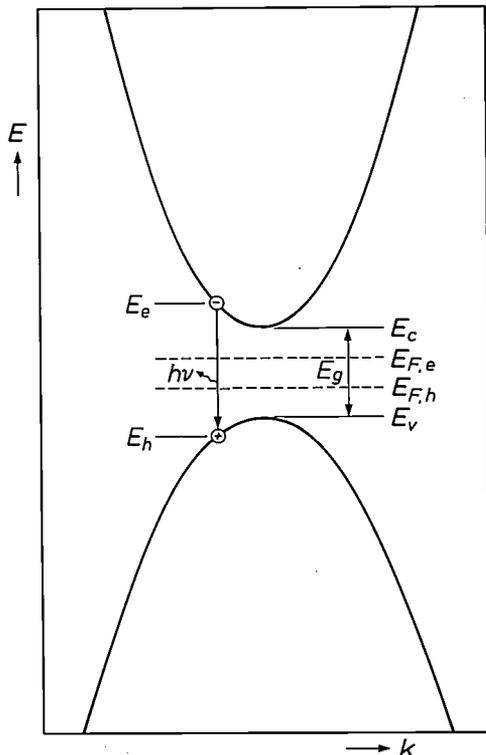


Fig. 11. Parabolic variation (schematic) of the energy E of the electrons in GaAs, as a function of the wave number k (momentum). GaAs is a direct-gap semiconductor: the minimum E_c of the conduction band and the maximum E_v of the valence band lie at about the same value of k . The difference between E_c and E_v is the band gap E_g . $E_{F,e}$ and $E_{F,h}$ are the quasi-Fermi energies for electrons and holes, respectively. A photon is generated when an electron (with energy E_e) in the conduction band recombines with a hole (energy E_h) in the valence band. This recombination must comply with the selection rule that the electron and the hole should have the same value of k (conservation of momentum).

tics. At low concentrations of conduction electrons and holes, approximate values of f_e and f_h are given by Boltzmann's formula:

$$f_e = \exp\{(E_{F,e} - E_e)/kT\}, \quad (14)$$

$$f_h = 1 - \exp\{(E_h - E_{F,h})/kT\}, \quad (15)$$

where $E_{F,e}$ and $E_{F,h}$ are the 'quasi-Fermi energies' for electrons and holes.

The difference in the temperature dependence of the radiative recombination coefficient B for a quantum well and for bulk material can be traced back to the differences in ρ_e and ρ_h . In bulk GaAs the conduction and valence bands are a parabolic function of the wave number k (fig. 11). For a conduction electron this means that the energy E_e is given by:

$$E_e = E_c + \frac{\hbar^2 k^2}{8\pi^2 m_e}, \quad (16)$$

where \hbar is Planck's constant and m_e is the effective electron mass. The number of energy levels at E_e in an interval dE_e is equal to the number of k -states in the corresponding interval dk . This number is proportional to the contents of a spherical shell in the k -space (hence proportional to $k^2 dk$), so that:

$$\rho_e \propto k^2 dk/dE_e. \quad (17)$$

Since it follows from eq. (16) that k is proportional to $(E_e - E_c)^{1/2}$ it is clear that:

$$\rho_e \propto (E_e - E_c)^{3/2}. \quad (18)$$

In the same way we find:

$$\rho_h \propto (E_v - E_h)^{3/2}, \quad (19)$$

and taking account of the k -selection rule (conservation of momentum) in the recombination:

$$\rho_{\text{red}} \propto (E_e - E_h - E_g)^{1/2}. \quad (20)$$

Substituting eqs (14) and (18) in eq. (10) gives:

$$n \propto \int_{E_c}^{\infty} (E_e - E_c)^{3/2} \{\exp(E_{F,e} - E_e)/kT\} dE_e. \quad (21)$$

If we put $\varepsilon = (E_e - E_c)/kT$, we obtain:

$$n \propto (kT)^{3/2} \exp\{(E_{F,e} - E_c)/kT\} \int_0^{\infty} \varepsilon^{3/2} \exp(-\varepsilon) d\varepsilon. \quad (22)$$

The integral here is independent of temperature, so that the temperature dependence of n is given by:

$$n \propto T^{3/2} \exp\{(E_{F,e} - E_c)/kT\}. \quad (23)$$

In the same way we find:

$$p \propto T^{3/2} \exp\{(E_v - E_{F,h})/kT\}, \quad (24)$$

$$L \propto T^{3/2} \exp\{(E_v - E_c + E_{F,e} - E_{F,h})/kT\}. \quad (25)$$

When eqs (23), (24) and (25) are substituted in eq. (9)

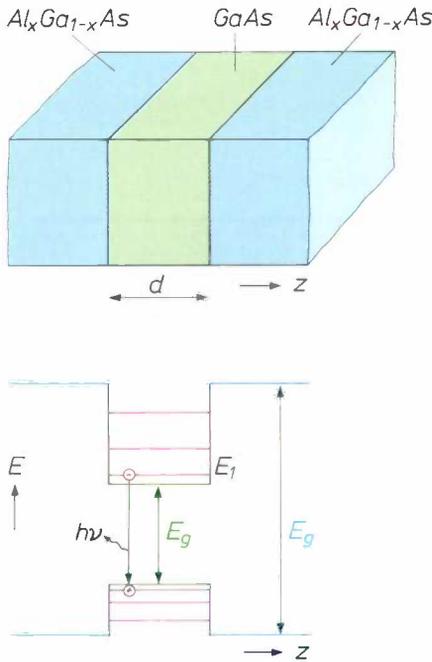


Fig. 12. Diagram of a quantum well consisting of a GaAs layer (thickness d) between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers with a larger band gap E_g . The layers are grown in the z -direction on a GaAs substrate. In a very thin GaAs layer the movements of electrons and holes in the z -direction correspond to discrete levels with energy differences that are much greater than the thermal energy kT . The electrons and holes produce photons on recombination. Almost all the electrons have the energy of the lowest conduction level (E_1) and almost all the holes have the energy of the highest valence level, in addition to the thermal energy for free movement parallel to the interfaces.

the exponential factors cancel, so that the temperature dependence of the radiative recombination coefficient B_b in bulk GaAs is given by:

$$B_b \propto T^{-\frac{3}{2}}. \tag{26}$$

In a quantum well, consisting of a thin layer of GaAs between two layers of $\text{Al}_x\text{Ga}_{1-x}\text{As}$, with a larger band gap, discrete energy levels exist for the electron movements perpendicular to the interfaces^[1]; see fig. 12. At these levels the increase in the density of states ρ_e with energy then takes place in steps. For a very thin well the energy difference between these levels is much greater than kT , so that almost all the electrons occupy the lowest level at energy E_1 . The density of states may then be assumed to be constant over the energy range where the occupancy f_e has a reasonable value, so that eq. (10) becomes:

$$n \propto \int_{E_1}^{\infty} \exp\{(E_{F,e} - E_e)/kT\} dE_e. \tag{27}$$

After substituting $\varepsilon = (E_e - E_1)/kT$, we then find for the temperature dependence:

$$n \propto T \exp\{(E_{F,e} - E_1)/kT\}. \tag{28}$$

The pre-exponential factor here is now T^1 instead of $T^{\frac{3}{2}}$, as in eq. (23). The same applies for p and L with respect to equations (24) and (25). Since the exponential factors are eliminated again on substituting in eq. (9), it follows that the temperature dependence of the radiative recombination coefficient B_q in a quantum well is given by:

$$B_q \propto T^{-1}. \tag{29}$$

As the well thickness increases the energy levels come closer together, so that the densities of states may no longer be assumed to be constant and the temperature dependence of B_q will change from T^{-1} to $T^{-\frac{3}{2}}$.

Experimental results

The temperature dependence of B can be determined experimentally by measuring the decay time τ of the luminescence. This is given by:

$$\tau^{-1} = Bn_0 + \tau_{nr}^{-1}, \tag{30}$$

where n_0 is the concentration of the majority carriers and τ_{nr}^{-1} is the probability of non-radiative transitions. If the luminescence efficiency is high, the second term can be neglected, so that τ^{-1} at constant n_0 has the same temperature dependence as B .

The measurements of τ were made for quantum wells of GaAs between $\text{Al}_{0.33}\text{Ga}_{0.67}\text{As}$, grown by MO-VPE with the reactor described here, producing

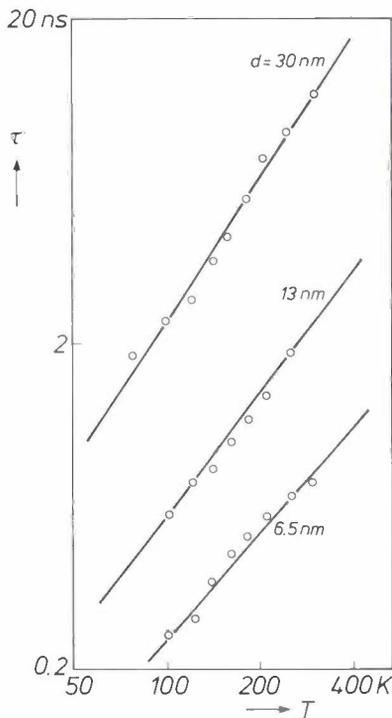


Fig. 13. Log-log plots of the luminescence decay time τ against the absolute temperature T , for three quantum wells of GaAs (thickness d) between two layers of $\text{Al}_{0.33}\text{Ga}_{0.67}\text{As}$. The exponent of the temperature dependence can be derived from the slopes of the lines.

^[10] See for example C. Kittel, Introduction to solid state physics, 4th edition, Wiley, New York 1971.

transition regions of one or two monolayers. *Fig. 13* shows a log-log plot of the measured values of τ against temperature for three wells of different thicknesses. In all three cases a reasonably straight line is found, so that $\tau \propto T^a$. The value of a corresponding to the slope is 1.07 for a thickness d of 6.5 nm, 1.22 for $d = 13$ nm and 1.41 for $d = 30$ nm. There is therefore good qualitative agreement with the predicted temperature dependence as a function of the layer thickness. A quantitative interpretation of the results is difficult because the variation of the energy levels with layer thickness is not accurately known.

Because of the smaller temperature dependence of τ in a quantum well the threshold current I_{th} for laser operation is also less dependent on temperature. In the ideal case ($\tau_{nr}^{-1} = 0$, no non-radiative processes) the value of I_{th} is determined solely by τ , so that the temperature dependence is given by:

$$I_{th} \propto T^a. \quad (31)$$

Usually an empirical relation is used to represent the temperature dependence of the threshold current:

$$I_{th} \propto \exp(T/T_0), \quad (32)$$

where the value of T_0 characterizes the temperature

behaviour. From equations (31) and (32) it can be shown that:

$$I_{th} (dI_{th}/dT)^{-1} = T/a = T_0. \quad (33)$$

For a conventional semiconductor laser at room temperature a is about 1.5, so that the maximum value of T_0 is about 200 K. In an ideal quantum well with high luminescence efficiency the value of a may approach 1, so that T_0 is about 300 K. The smaller increase in threshold current with temperature will lengthen the life of quantum-well lasers.

In the work described here Ing. H. J. Talen-van der Mheen took part in the growth experiments and C. W. T. Bulle-Lieuwma and Drs A. F. de Jong contributed to the TEM investigations.

Summary. A novel type of reactor has been developed for making MO-VPE multilayer structures, in which very sharp transitions can be obtained between layers of different composition. The perfection of the structures produced is very well demonstrated by means of transmission electron microscopy. It is demonstrated for instance that the transition regions are no thicker than approximately one monolayer. The optical quality of the quantum wells is determined by measuring the decay time of the luminescence. It appears that the decay time in quantum wells is less dependent on temperature than in bulk material, so that the threshold current for laser operation is also less temperature-dependent.

Molecular beam epitaxy of multilayer structures with GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$

B. A. Joyce and C. T. Foxon

Introduction

As indicated in the opening article of this issue ^[1], thin epitaxial films of some III-V semiconductors have properties which are important in high-frequency and optoelectronic applications. Such films can be prepared by conventional growth techniques such as liquid-phase epitaxy (LPE) or vapour-phase epitaxy (VPE). For particular applications it may be advantageous to use metal-organic vapour-phase epitaxy (MO-VPE), discussed in the previous articles ^{[2][3]}, or molecular beam epitaxy (MBE) ^[4].

MBE is a refined form of ultra-high vacuum evaporation. In this technique thermally generated collision-free molecular (or atomic) beams of the constituent elements, formed in Knudsen sources, are deposited on a heated substrate where they react to form an epitaxially related crystalline film. Growth temperatures are usually somewhat lower than those used in the more conventional growth techniques, while growth rates are in the range 0.01 to 1 nm/s. Because the molecular beam fluxes can be started or stopped rapidly using a simple mechanical shutter, atomically abrupt interfaces can be obtained. Using *in situ* analysis it is possible to measure directly the intensity of the various molecular beams and as a result films of high crystalline quality with precise control of the thickness, composition and doping level can be obtained. Changes in composition and doping level on an atomic scale facilitate the formation of well-defined multilayer structures with special properties ^[5].

The growth of III-V semiconductor films by MBE has been studied extensively at Philips Research Laboratories in Redhill, England. This study included the preparation of various types of III-V compounds and layered structures, investigation of the mechanisms controlling the growth and dopant incorporation, characterization of surfaces and interfaces,

measurements of film or layered-structure properties and the improvement of the MBE equipment. The film properties have been studied *in situ* using non-destructive surface analysis. The III-V semiconductors grown by MBE have included mainly the binary compound GaAs and the related ternary alloy $\text{Al}_x\text{Ga}_{1-x}\text{As}$ but some work has also been carried out both on the fundamental properties controlling growth and on measurements of film properties for other alloys such as InGaAs, InGaP, GaAsP, InAsP and InGaAsP. In all cases single-crystal films were grown epitaxially on GaAs or for InGaAs on InP.

The progress of MBE technology has led to films and multilayer structures having adequate quality for device applications. Background donor and acceptor levels can be as low as $2 \times 10^{14} \text{ cm}^{-3}$. In a hetero-structure of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ a two-dimensional electron gas is formed in the GaAs material near the interface, showing extremely high electron mobilities at low temperatures: $> 3 \times 10^6 \text{ cm}^2\text{V}^{-1}\text{s}^{-1}$ at 4 K. With structures of GaAs between $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers multiple quantum-well lasers have been grown in which the GaAs wells are as thin as 1.3 nm, corresponding to about five monolayers of material. At

[1] J. Wolter, Research on layered semiconductor structures, this issue, pp. 111-117.

[2] P. M. Frijlink, J. P. André and M. Erman, Metal-organic vapour phase epitaxy of multilayer structures with III-V semiconductors, this issue, pp. 118-132.

[3] M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft, Metal-organic vapour phase epitaxy with a novel reactor and characterization of multilayer structures, this issue, pp. 133-142.

[4] See for example: L. L. Chang, L. Esaki, W. E. Howard, R. Ludeke and G. Schul, Structures grown by molecular beam epitaxy, *J. Vac. Sci. & Technol.* **10**, 655-662, 1973;

A. Y. Cho and J. R. Arthur, Molecular beam epitaxy, *Progr. Solid State Chem.* **10**, 157-191, 1975.

[5] A. C. Gossard, P. M. Petroff, W. Wiegmann, R. Dingle and A. Savage, Epitaxial structures with alternate-atomic-layer composition modulation, *Appl. Phys. Lett.* **29**, 323-325, 1976. See also: G. H. Döhler, Solid-State superlattices, *Sci. Am.* **249** (No. 5), 118-126, 1983;

D. W. Shaw, Advanced multilayer epitaxial structures, *J. Cryst. Growth* **65**, 444-453, 1983.

this well width a visible laser emission is obtained at 704 nm^[6]. To our knowledge this is the shortest wavelength reported for a room-temperature injection laser with only GaAs in the wells.

In this article the discussion on MBE will be restricted to certain aspects of the growth of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$. Among the III-V compounds and alloys these materials are the most widely investigated and the most promising for practical devices. Before dealing with some results of our investigations a general description of the MBE process will be given. The results to be discussed concern the surface chemistry of growth, the control of the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ composition and the growth dynamics. Special attention will be given to the interface between GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ films. Finally some properties and possible applications will be discussed.

The MBE process

It will be useful to define at the outset the important features of the MBE process in terms of the apparatus used, process parameters and source materials^[7].

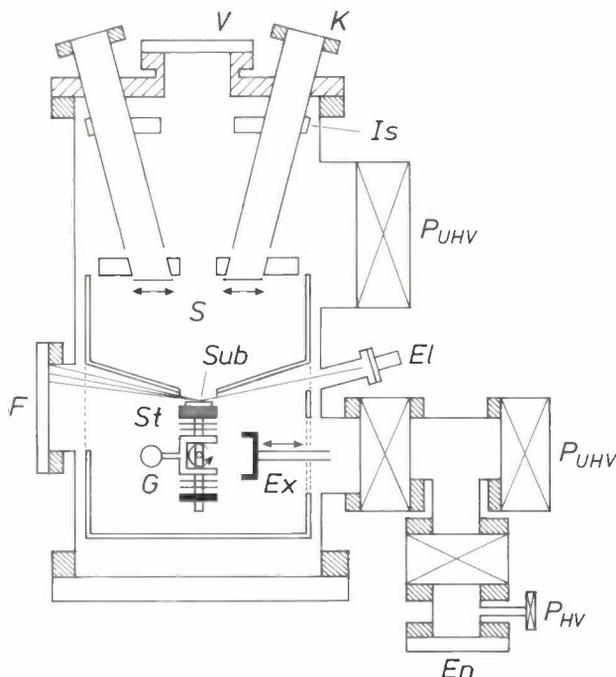


Fig. 1. Schematic diagram of an MBE system showing the essential features required for the growth of III-V compound semiconductor films. Atomic or molecular beams of the film constituents are produced in the heated Knudsen effusion cells *K*. By opening the shutters *S* the beams are directed to the heated substrate *Sub* on the stage *St*, which can be rotated. The beam fluxes are monitored by the ion gauge *G*. Substrates are introduced via a sample entry lock *En* and a sample exchange mechanism *Ex*. A RHEED system consisting of an electron gun *El* and a fluorescent screen *F* is used to study the substrate surface before growth and the film surface during growth. Other facilities for *in situ* analysis such as a quadrupole mass spectrometer and an Auger electron spectrometer are not shown. *Is* water-cooled thermal isolator. *P_{HV}* pump for obtaining intermediate vacuum. *P_{UHV}* pump for obtaining ultra-high vacuum. *V* view port.

Although MBE is simply a refined form of vacuum evaporation in which directed neutral thermal atomic and molecular beams impinge on a heated substrate under ultra-high vacuum conditions, the apparatus required to achieve this with the necessary degree of control has become rather complex. The essential elements of a system suitable for growth of III-V compound films are illustrated schematically in *fig. 1*. It is based on a two- or three-chamber stainless-steel ultra-high-vacuum system ($< 10^{-9}$ Pa), usually ion-pumped or cryopumped, and incorporates large areas of liquid-nitrogen-cooled panels. The provision of a vacuum interlock, to enable the substrate to be introduced into the growth chamber without breaking the vacuum, is essential if films having high-quality electrical and optical properties are to be prepared.

The atomic or molecular beams are formed in heated Knudsen effusion cells. An example of such a cell is shown in *fig. 2*. The material to be evaporated is contained in a crucible, which is heated by radiation from a separate winding. In the cell we attempt to ensure that the liquid or solid phase and the vapour are in equilibrium (Knudsen evaporation). The mean free path of the vapour (atoms or molecules) is much

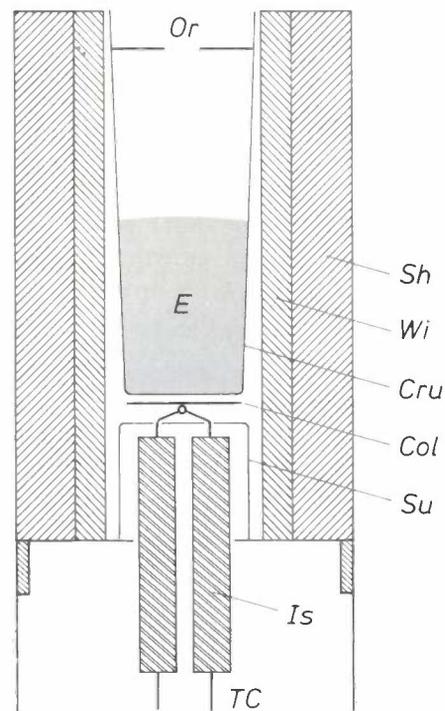


Fig. 2. Schematic diagram of a typical Knudsen effusion cell, acting as source of the Group III element (Ga or Al) for the MBE growth of III-V-films. The evaporation source *E* is contained in a crucible *Cru* of boron nitride. The heater winding *Wi* is surrounded by heat-shields *Sh* of tantalum. The beam of evaporated particles leaves the cell through the orifice *Or*. The temperature is read by a thermocouple *TC* passing through the insulators *Is* and connected to a radiation collector *Col* of tantalum. *Su* thermocouple support of boron nitride.

larger than the cell orifice, ensuring molecular flow. The inner crucible is made either of reactor-grade graphite or, preferably, pyrolytic boron nitride, so that the possible contamination of the molecular beams is minimized. Beam intensities are controlled by the cell temperatures and for molecular flow are given by [7]:

$$J_i = \frac{ap_i \cos \theta}{\pi d^2 (2\pi m_i kT)^{\frac{1}{2}}},$$

where J_i is the flux per unit area at a distance d from

surface structure [9]. Additionally there may be an Auger electron spectrometer to determine surface composition and purity [10]. Some or all of these facilities may be housed in a separate analysis/preparation chamber mounted between the sample-insertion interlock and the growth chamber. A general view of the equipment used for our MBE experiments is shown in fig. 3.

A RHEED arrangement consists simply of a 5-20 keV electron gun and a fluorescent screen. The electron beam is incident at a very shallow angle ($1^\circ - 3^\circ$) to

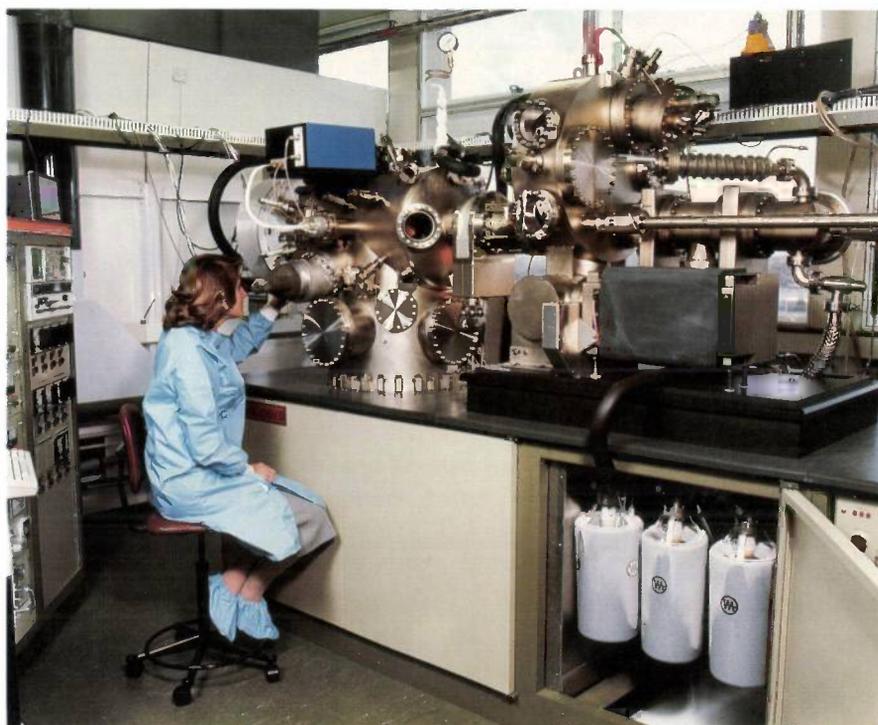


Fig. 3. General view of the MBE equipment (Varian GEN II) at the Philips Research Laboratories in Redhill. The growth occurs in the vacuum chamber on the left by effusing molecular beams from the tubular Knudsen cells. The growth rates for GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ can be deduced from RHEED intensity variations measured in the growth chamber. The vacuum system on the right is used for the introduction and preparation of substrates.

the source, which has an orifice of area a and contains atoms (molecules) of mass m_i having an equilibrium vapour pressure p_i at temperature T (in K). θ is the angle between the beam and the substrate surface normal. Individual cells are thermally isolated by a succession of heat shields and the flux is regulated by a shutter operating in front of each aperture. A single microcomputer can be used to control all of the cell temperatures (within ± 0.25 K) and all of the shutters to achieve any growth sequence.

Fluxes are usually monitored by an ion gauge which can be rotated in and out of the beams. Other analytical facilities may include a quadrupole mass spectrometer for residual gas analysis and leak-checking [8], and a reflection high-energy electron diffraction (RHEED) arrangement for assessment of

- [6] K. Woodbridge, P. Blood, E. D. Fletcher and P. J. Hulyer, Short wavelength (visible) GaAs quantum well lasers grown by molecular beam epitaxy, *Appl. Phys. Lett.* **45**, 16-18, 1984.
- [7] Various MBE aspects have been extensively reviewed by K. Ploog, *Molecular beam epitaxy of III-V compounds*, in: H. C. Freyhardt (ed.), *Crystals — growth, properties and applications*, Vol. 3, Springer, Berlin 1980, pp. 73-162. A review more directly related to the work described in this article has been given by C. T. Foxon, *Molecular beam epitaxy*, *Acta Electron.* **21**, 139-150, 1978.
- [8] In a quadrupole mass spectrometer the ions formed in the ionizer are mass-separated using r.f. and d.c. fields generated by four electrodes.
- [9] The importance of *in situ* RHEED studies in the MBE development has been demonstrated by (for example): A. Y. Cho, *J. Appl. Phys.* **41**, 2780-2786, 1970, and **42**, 2074-2081, 1971, and by J. H. Neave and B. A. Joyce, *J. Cryst. Growth* **44**, 387-397, 1978.
- [10] A survey of various surface-analysis methods, including electron diffraction and Auger electron spectroscopy has been given by H. H. Brongersma, F. Meijer and H. W. Werner, *Surface analysis, methods of studying the outer atomic layers of solids*, *Philips Tech. Rev.* **34**, 357-369, 1974.

the substrate surface. The diffracted electron beams form a characteristic pattern on the fluorescent screen, which provides information on the morphology and the symmetry of the surface. Because of the small incidence angle the information obtained in RHEED is derived from the first few atomic layers. Temporal variation in the intensity of diffraction features during MBE growth can be related to growth dynamics.

Since the atomic or molecular beams are neutral and collimated they are necessarily divergent, and usually non-axial with respect to the substrate. To obtain uniform growth rate or composition (in the case of alloy films) over a large substrate (diameter 5 to 8 cm), the substrate stage is rotated at speeds between 0.03 and 2.0 Hz. It was shown that with this technique a growth-rate variation of $< 1\%$ could be obtained across a 5-cm substrate, with comparable control of alloy composition^[11]. For alloy growth, however, the time taken for a complete revolution must be shorter than the time to grow a monolayer, to avoid a periodic modulation of alloy composition in the direction of growth.

To provide the basis of high-quality films, a substrate surface free of crystallographic and other defects and clean on an atomic scale (≤ 0.01 monolayer of impurities) must be prepared. This usually involves free etching by an oxidative process which removes any carbon and leaves the surface covered with a protective volatile oxide, which is subsequently removed in the vacuum system by heating in a beam of arsenic. The final quality of the substrate surface can be checked by RHEED and Auger spectrometry, for example.

Growth is initiated by bringing the substrate to the appropriate temperature (typically in the range 500-700 °C) in a beam of arsenic and then opening the shutter of the Group III element source. Growth rates are in the range 0.03-3 nm/s (0.1-10.0 $\mu\text{m}/\text{h}$), corresponding to beam fluxes from 5×10^{13} to 5×10^{15} at. $\text{cm}^{-2}\text{s}^{-1}$. The arsenic flux is typically 3-5 times greater than that of the Group III element, but it is the latter which determines the growth rate.

Source materials are usually elemental, and all of the Group III elements produce monoatomic beams. The arsenic source is rather more complex, however, both in terms of the species produced and the methods used to produce them. When evaporation takes place directly from the element, the flux consists entirely of tetratomic molecules (As_4), but if GaAs is used as the source, the arsenic flux is dimeric^[12]. Alternatively, dimers can be produced from the element by using a two-zone Knudsen cell in which a tetramer flux is formed conventionally and passed through an optically

baffled high temperature stage, which can be designed to produce a complete conversion to a dimeric flux^[13].

Dopant beams are also produced from Knudsen effusion cells. The elements used for doping in MBE are principally Be for the preparation of p-type films and Si or Sn for the preparation of n-type films. All these elements evaporate as monomers. Free-carrier concentrations can be controlled over the range between 10^{14} and 10^{19} carriers per cm^3 . The upper limit is set by the limited solid solubility in the lattice, the lower by the presence of background impurities.

Surface chemistry of growth

Surface kinetic data on growth and doping can be obtained by using modulated molecular beam techniques^{[12][14]}. The flux of an incident species can be modulated by opening and closing the shutter of the Knudsen cell periodically or by placing a beam chopper inside the vacuum system. The perturbation of the incident flux gives rise to a time-varying concentration of chemisorbed atoms and molecules on the substrate surface. This leads to a time-dependent desorption rate, which can be measured mass-spectrometrically. Alternatively, by modulating the desorption flux with a beam chopper, atomic and molecular species which are leaving the surface directly can be identified and distinguished from background vapour species.

The kinetic parameters to be extracted from modulated molecular-beam experiments include the surface residence time, desorption energy, sticking coefficient and order of reaction in which the adsorbed and desorbed species are involved. From these data detailed interaction mechanisms have been deduced, especially for the reaction of gallium and arsenic on (100)

[11] A. Y. Cho and K. Y. Cheng, Growth of extremely uniform layers by rotating substrate holder with molecular beam epitaxy for applications to electro-optic and microwave devices, *Appl. Phys. Lett.* **38**, 360-362, 1981;

K. Y. Cheng, A. Y. Cho and W. R. Wagner, Molecular-beam epitaxial growth of uniform $\text{Ga}_{0.47}\text{In}_{0.53}\text{As}$ with a rotating sample holder, *Appl. Phys. Lett.* **39**, 607-609, 1981.

[12] C. T. Foxon, J. A. Harvey and B. A. Joyce, The evaporation of GaAs under equilibrium and non-equilibrium conditions using a modulated beam technique, *J. Phys. & Chem. Solids* **34**, 1693-1701, 1973.

[13] J. H. Neave, P. Blood and B. A. Joyce, A correlation between electron traps and growth processes in *n*-GaAs prepared by molecular beam epitaxy, *Appl. Phys. Lett.* **36**, 311-312, 1980.

[14] C. T. Foxon, M. R. Boudry and B. A. Joyce, Evaluation of surface kinetic data by the transform analysis of modulated molecular beam measurements, *Surf. Sci.* **44**, 69-92, 1974.

[15] C. T. Foxon and B. A. Joyce, Interaction kinetics of As_2 and Ga on {100} GaAs surfaces, *Surf. Sci.* **64**, 293-304, 1977.

[16] C. T. Foxon and B. A. Joyce, Interaction kinetics of As_4 and Ga on {100} GaAs surfaces using a modulated molecular beam technique, *Surf. Sci.* **50**, 434-450, 1975.

[17] C. T. Foxon and B. A. Joyce, Surface processes controlling the growth of $\text{Ga}_x\text{In}_{1-x}\text{As}$ and $\text{Ga}_x\text{In}_{1-x}\text{P}$ alloy films by MBE, *J. Cryst. Growth* **44**, 75-83, 1978.

substrates of GaAs. It has been demonstrated that the use of an As_4 flux (from an elemental arsenic source) and that of an As_2 flux (from a GaAs source) lead to substantially different surface reactions [15] [16].

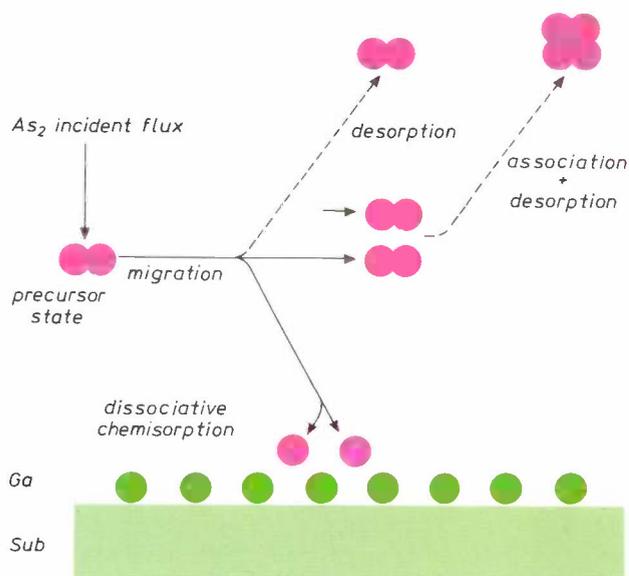


Fig. 4. Model of the chemistry of growth of a GaAs film from molecular beams of Ga and As_2 . The As sticking coefficient approaches unity when the substrate is already covered with a monolayer of Ga atoms. An As_2 molecule is first adsorbed into a weakly bound precursor state, in which it can migrate on the surface. Dissociative chemisorption can occur on Ga adatoms with a sticking coefficient ≤ 1 . Excess As_2 molecules can either desorb from the precursor state or associate at low temperatures, leading to desorption of an As_4 molecule.

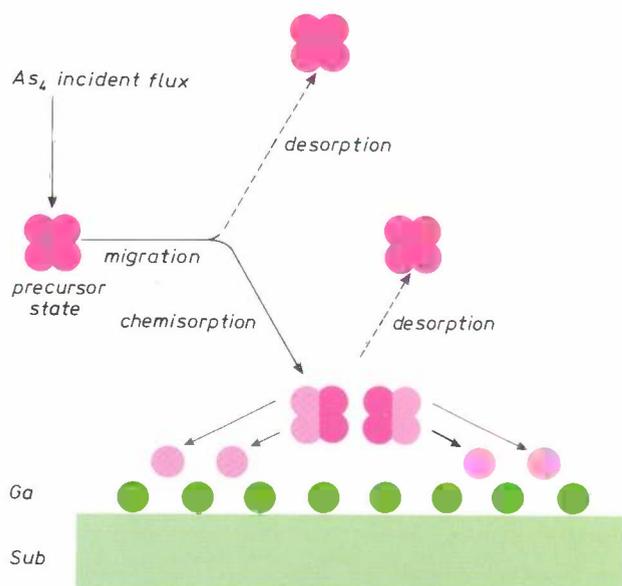


Fig. 5. Model of the chemistry of growth of a GaAs film from molecular beams of Ga and As_4 . Here an As_4 molecule, adsorbed into a weakly bound precursor state, migrates on the substrate. The molecules can desorb or be chemisorbed on adjacent Ga adatoms. Pairwise interaction of two chemisorbed As_4 molecules on adjacent Ga atoms leads to the incorporation of four As atoms in the lattice and to desorption of another four As atoms as an As_4 molecule. As a result the sticking coefficient of As_4 is ≤ 0.5 .

The interactions of the As_2 flux and gallium on a GaAs substrate are summarized in the growth model shown in *fig. 4* [15]. According to this model the As_2 molecules are first adsorbed into a mobile, weakly bound precursor state. The surface residence time of molecules in this state is less than 10^{-5} s. The basic process for As incorporation during thin-film growth is a simple first-order dissociative chemisorption on surface Ga atoms. When the Ga flux is lower than half the As_2 flux, one As atom sticks for each Ga atom supplied. The maximum value of the As_2 sticking coefficient is unity, which is obtained when the surface is covered with a complete monolayer of Ga atoms (as in *fig. 4*). At relatively low substrate temperatures (< 600 K) there is an additional association reaction to form As_4 molecules. These desorb very slowly by a first-order reaction. At substrate temperatures above 600 K some dissociation of GaAs may occur.

By contrast, the reaction mechanism with an incident As_4 flux is significantly more complex [16]. The model for this is shown in *fig. 5*. Here the As_4 molecules are first adsorbed into a mobile precursor state. The migration of the adsorbed As_4 molecules has an activation energy of about 0.25 eV. The surface residence time is temperature dependent, from which a desorption energy of about 0.4 eV may be determined. A crucial finding is that the sticking coefficient of arsenic never exceeds 0.5, even when the Ga flux is much higher than the As_4 flux or when the surface is completely covered with a monolayer of Ga atoms (as in *fig. 5*). This is explained by assuming a pairwise dissociation of As_4 molecules chemisorbed on adjacent Ga atoms. This second-order reaction is the key feature of the thin-film growth of GaAs with an As_4 flux. From any two As_4 molecules four As atoms are incorporated in the GaAs lattice, while the other four desorb as an As_4 molecule.

Composition control for $\text{Al}_x\text{Ga}_{1-x}\text{As}$ films

For the MBE growth of the ternary III-III-V alloy $\text{Al}_x\text{Ga}_{1-x}\text{As}$, beams of Al, Ga and excess As_2 or As_4 are directed simultaneously at the substrate surface. In this way films of good crystallographic perfection can be obtained. A critical factor, however, in alloy-film growth by MBE is the production of films having a homogeneous composition.

Growth reactions for III-III-V alloys, in so far as the Group V element is concerned, are similar to those observed in the growth of binary compounds. The only problem is to establish those kinetic factors which can influence the ratio of the Group III elements [17]. An important point to note, however, is that the thermal stability is determined by the less stable of the two

binary compounds of which the alloy may be considered to be composed. In the case of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ films AlAs is more thermally stable than GaAs. Thus by using growth parameters comparable to those for GaAs, films of good compositional uniformity and control can be prepared.

At comparatively low substrate temperatures ($< 575^\circ\text{C}$) the sticking coefficients of Al and Ga are both unity, so the alloy composition is determined simply by the ratio of the two atomic fluxes. However, to obtain adequate electrical and/or optical properties it is frequently necessary to use significantly higher substrate temperatures. In that case preferential desorption of the more volatile Ga occurs, which means that the sticking coefficient of Ga is less than unity. The value for Al, on the other hand, does not change significantly over the entire temperature range used ($\approx 650\text{--}750^\circ\text{C}$). The effective loss of Ga can be readily calculated from its known vapour pressure over GaAs.

Growth dynamics

It is possible to explore some features of the growth dynamics of MBE by monitoring temporal variations in the intensity of various features in the RHEED pattern. It has been found that damped oscillations in the intensity of both the specular and diffracted beam occur immediately after initiation of growth^[18]. A typical example for the specular beam is shown in *fig. 6*. The period of oscillation corresponds exactly to the growth of a single monolayer, i.e. a complete

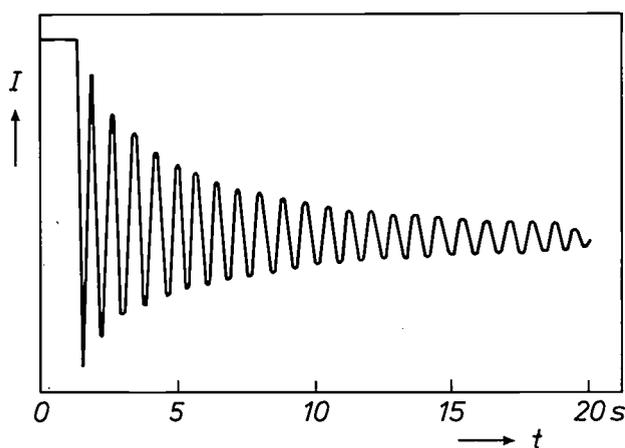


Fig. 6. Oscillations of the intensity of the specular beam in the RHEED pattern from a (2×4) reconstructed surface in the $[110]$ azimuth, during the growth of a thin GaAs film on a (001) substrate of GaAs^[17]. The notation (2×4) indicates an enlargement of the surface unit cell by a factor of two and four in the $[110]$ and $[110]$ directions respectively. The growth is initiated by an incident flux of Ga atoms at the heated substrate which is maintained in a flux of either As_2 or As_4 . The intensity I is given in arbitrary units as a function of the growing time t . The oscillation period corresponds to the growth of a monolayer of GaAs.

layer of Ga and As atoms, which in the $[001]$ direction is equal to the thickness of half the lattice constant. The period is independent of the azimuth of the incident beam and of the particular diffraction feature being measured^[19]. The amplitude, however, is strongly dependent on both of these parameters. For evaluation of growth dynamics most of the information is contained in the specular beam, so we shall limit the discussion to this feature.

Similar oscillatory effects have been observed during epitaxial growth of thin metal or silicon films, when studied *in situ*, for example by Auger spectroscopy or low-energy electron diffraction (LEED)^[20]. The observation of these effects is usually associated with a layer-by-layer growth process (i.e. two-dimensional nucleation). Detailed analysis of the oscillations provides important information on growth dynamics.

If we equate changes in intensity of the specular beam in the RHEED pattern with changes in surface roughness, a smooth equilibrium surface corresponds to high reflectivity. On commencement of growth, clusters are formed at random positions on the crystal surface, leading to a decrease in the reflectivity. This decrease can be predicted for purely optical reasons, since the de Broglie wavelength of the electrons is about 0.012 nm while the bi-layer step height is about 0.28 nm; i.e. the wavelength is at least an order of magnitude less than the size of the scatterer, so diffuse scattering will result. Nucleation is not restricted to a single layer, but can reoccur before the preceding monolayer is complete. In the early stages, however, one monolayer is likely to be almost complete before the next monolayer starts, so the reflectivity will increase as the surface again becomes smooth on the atomic scale, but with subsequent roughening as the next monolayer develops. This repetitive process will cause the oscillations in reflectivity to be gradually damped as the surface becomes statistically distributed over several incomplete monolayers.

In *fig. 7* we show a real-space representation of the formation of two monolayers which illustrates how the oscillations in the intensity of the specular beam

^[18] J. H. Neave, B. A. Joyce, P. J. Dobson and N. Norton, Dynamics of film growth of GaAs by MBE from RHEED observations, *Appl. Phys. A* **31**, 1-8, 1983.

^[19] J. J. Harris, B. A. Joyce and P. J. Dobson, Oscillations in the surface structure of Sn-doped GaAs during growth by MBE, *Surf. Sci.* **103**, L90-L96, 1981.

^[20] V. Bostanov, R. Roussinova and E. Budevski, Multinuclear growth of dislocation-free planes in electrocrystallization, *J. Electrochem. Soc.* **119**, 1346-1347, 1972; Y. Namba, R. W. Vook and S. S. Chao, Thickness periodicity in the Auger line shape from epitaxial (111) Cu films, *Surf. Sci.* **109**, 320-330, 1981; K. D. Gronwald and M. Henzler, Epitaxy of Si (111) as studied with a new resolving LEED system, *Surf. Sci.* **117**, 180-187, 1982.

occur. There is a maximum in reflectivity for the initial and final smooth surfaces and a minimum (or maximum in diffuse scattering) for the intermediate stage when the growing monolayer is approximately half

complete. We have, however, also established that the amplitude of the intensity oscillations is dependent on the direction of the primary electron beam, being greater for the beam incident along $[110]$ than for the beam incident along $[\bar{1}10]$. This would suggest, if we take our optical model one stage further, that most of the steps which develop on the surface are along the $[\bar{1}10]$ direction, as shown in fig. 7, i.e. they cause maximum diffuse scattering when their longer edges are normal to the incident beam.

From the results obtained we may conclude that growth occurs principally by a two-dimensional monolayer-by-monolayer process, but new monolayers are able to start before preceding ones have been completed. The oscillation period provides a continuous and absolute growth rate monitor with atomic-layer precision.

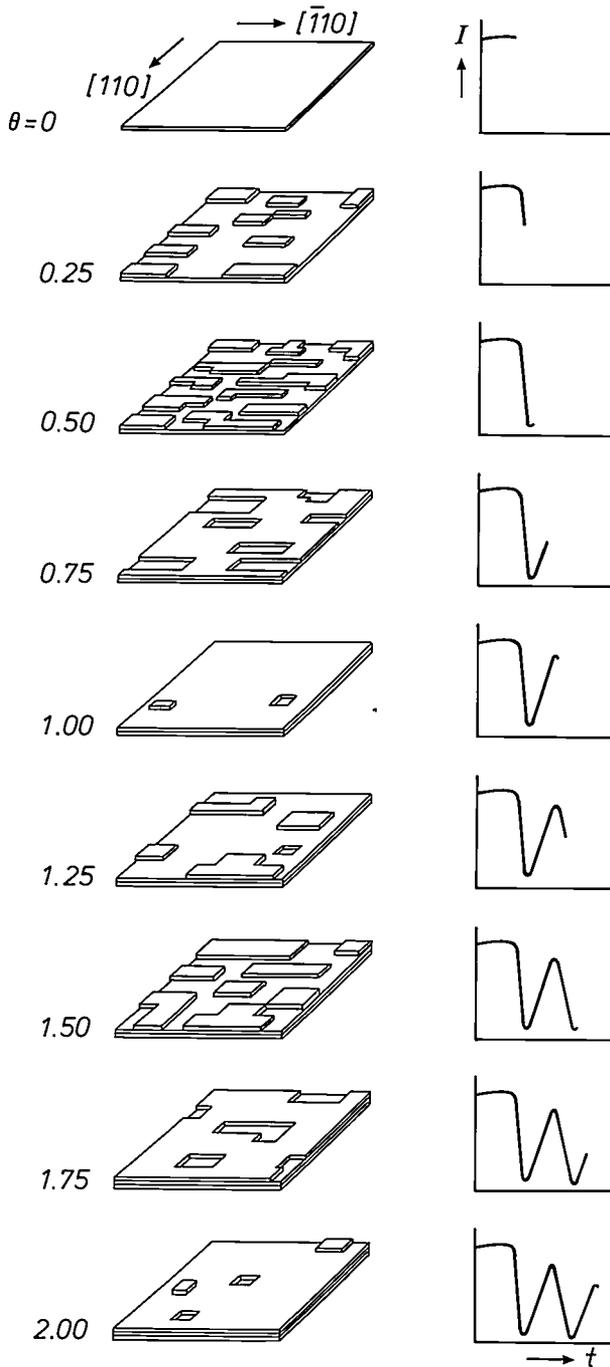


Fig. 7. Real-space representation of the formation of the first two monolayers of GaAs (left) with the corresponding RHEED oscillations (right). The intensity I is given as a function of time t for various numbers θ of monolayers deposited. When a smooth substrate is progressively covered with less than half a monolayer, the reduction in smoothness leads to a continuous decrease of reflectivity and hence of intensity. At higher coverages, up to one monolayer, the surface again becomes smoother. The intensity maximum at $\theta = 1$ is however lower than that at $\theta = 0$ because the growth of a monolayer is not perfectly two-dimensional. Similar intensity variations occur between $\theta = 1$ and $\theta = 2$. The increasing influence of the deviation from two-dimensional growth leads to a progressive damping of the amplitude of the oscillation, as shown in fig. 6.

GaAs- $\text{Al}_x\text{Ga}_{1-x}\text{As}$ interfaces

The MBE method described here has been used to prepare various multilayer structures based on GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$. These include: single-heterojunction two-dimensional electron-gas structures and multiple-quantumwell (MQW) structures formed from thin semiconductor films (usually GaAs) confined by layers of higher band-gap material (usually AlGaAs); when the structure is periodic and the confining barriers are thin enough for the wells to be electronically coupled the structure is known as a superlattice. For the properties of such structures the quality of the interfaces is very important, as has been indicated in the previous articles of this issue.

It is clear from the RHEED study of growth dynamics^[18] that films of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ grow predominantly by a process of two-dimensional nucleation of new monolayers. The damped intensity oscillations are however indicative of some non-ideality, i.e. growth is not perfectly two-dimensional. It is not yet possible to extract quantitative information on the structure and composition of GaAs- $\text{Al}_x\text{Ga}_{1-x}\text{As}$ interfaces from the RHEED observations, but several other techniques have been utilized. It is important to emphasize, however, that experimentally realizable interface perfection is so high that conventional methods of profiling do not have the required resolution to test it. These methods involve ion sputtering to section the material followed by some means of composition determination, such as Auger electron spectrometry and secondary-ion mass spectrometry (SIMS)^[10].

The best claimed (and probably optimistic) resolution limit of 1 nm was based on a 90%-10% peak-height range of the Auger signal of Al, from which an

interfacial width of 1.3 nm was deduced for a GaAs- $\text{Al}_{0.5}\text{Ga}_{0.5}\text{As}$ interface [21]. Even when a superlattice of GaAs and $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ was 'expanded' by a factor of 2400 using a very shallow angle bevelling technique, the minimum well width of GaAs in which the Auger signal of Al went to zero ($< 1\%$ Al) was 5 nm, probably due to ion beam mixing and other sputtering effects [22].

To put these results in context, it has been demonstrated by other characterization techniques that superlattices of alternate monolayers of GaAs and AlAs can be grown successfully with MBE [6] [7]. The available methods of obtaining structural and compositional information with the necessary level of spatial resolution are limited to transmission electron microscopy (TEM) of cross-sections, and to optical techniques using photoluminescence and excitation spectroscopy.

Characterization with TEM

In TEM cross-sections, when the superlattice period is greater than a few atomic layers it becomes very difficult to form a dark-field image from the superlattice spots since they are then too close to the main lattice reflections. However, because Ga and Al have significantly different scattering factors, the superlattice can still be imaged [3], as shown in *fig. 8a*. Here the dark bands correspond to GaAs and the light bands to $\text{Al}_x\text{Ga}_{1-x}\text{As}$, but it is clear that the resolution is not adequate to define the interface on an atomic scale.

It is also possible to produce lattice-plane images of superlattices at extremely high resolution, but it is only possible to obtain a reasonable contrast distinction between the two materials for GaAs-AlAs structures. For GaAs- $\text{Al}_x\text{Ga}_{1-x}\text{As}$ structures the interfaces are not clearly defined, but it is nevertheless evident that there is lattice-point alignment with no crystallographic disorder across the interface. An example of this is shown in *fig. 8b*, a cross-section image of a quantum well consisting of GaAs between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers. The compositional variation cannot be precisely determined, however.

Luminescence investigations

Probably the best available method for investigating interface disorder effects in superlattices or multiple quantum-well structures relies on their optical characterization [23]. The measurements are of photoluminescence and excitation spectra, typically at about 4 K, and the features observed are due to free or bound excitons, i.e. weakly bound mobile electron-hole pairs. Typical photoluminescence and excitation spectra for a sample which was grown in our Varian GEN II system at 650 °C are shown in *fig. 9*. The

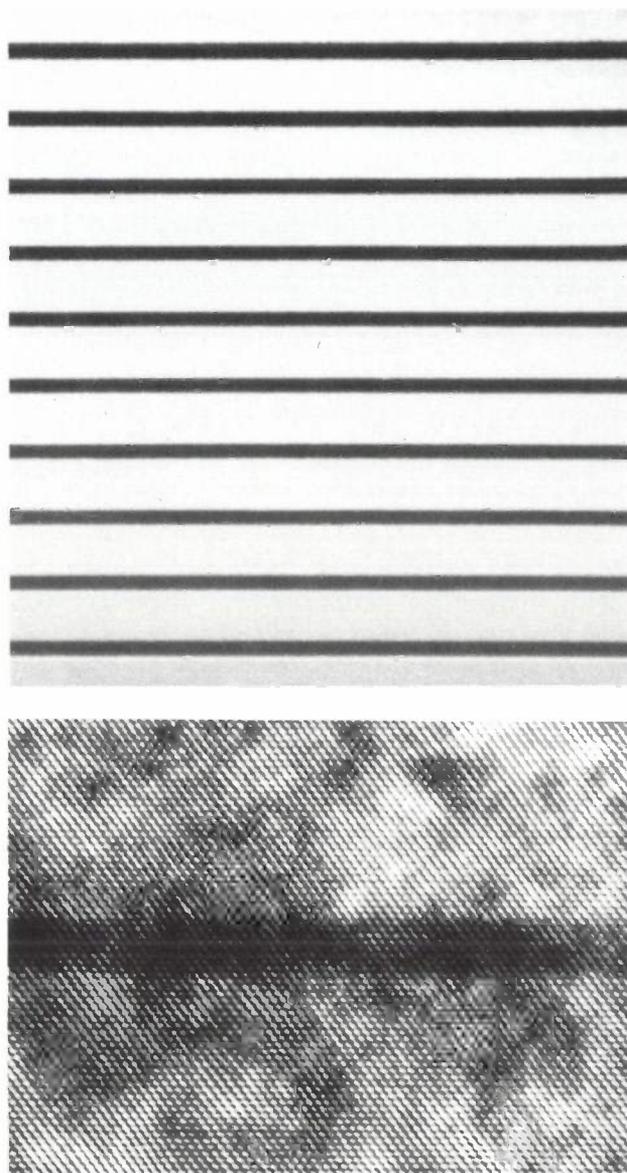


Fig. 8. Above: (110) cross-sectional TEM micrograph of a multiple quantum-well structure. The GaAs wells (dark) are 5.5 nm thick, the $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers (light) 17.5 nm. The micrograph was taken under dark-field conditions using the (002) diffraction. *Below:* high resolution (110) cross-sectional TEM micrograph (dark-field image) of a 2.7-nm quantum well of GaAs (dark) between two $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers (light). The samples were grown by K. Woodbridge and the micrographs were taken by J. P. Gowers and D. J. Smith (University of Cambridge).

sample consists of five GaAs quantum wells 5.5 nm thick between barriers of $\text{Al}_x\text{Ga}_{1-x}\text{As}$.

Exciton linewidths become broader as the well width decreases, but this could result from two possible effects. The first is a layer-to-layer thickness variation, leading to different energy levels for the conduction electrons in the different layers. The confinement of the conduction electrons in GaAs wells between the potential barriers of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ gives rise to the occurrence of discrete levels, as in the theoretical example of a particle confined in a box. The approximate value of the confinement energy E_n , as-

suming infinitely deep wells, is given by:

$$E_n = \frac{n^2 \hbar^2}{8m_e d^2},$$

where n is an integer, m_e is the effective mass of conduction electrons, \hbar is Planck's constant and d the well thickness. For a thickness variation Δd which is independent of d , the variation of the confinement energy (ΔE_n) is given by:

$$\Delta E_n = c d^{-3},$$

where c is a constant.

The second alternative is for a thickness variation within each layer, but with a constant average layer thickness. This is the case corresponding to a growth process which is not perfectly two-dimensional. It will lead to a modification of the exciton energy provided the lateral scale of the interface roughness is greater than the exciton diameter, which is about 30 nm in bulk GaAs. When the lateral scale of the interface roughness becomes smaller than the diameter it will 'feel' an average potential, so that the energy levels will be sharp and narrow luminescence peaks will be observed. A study of the linewidths alone can therefore give information only if the lateral scale of the roughness is greater than the exciton diameter.

It is not possible to distinguish between these alternatives in transmission and luminescence measurements, since all the layers are then probed simultaneously, but a distinction can be made by excitation spectroscopy. In this technique the luminescence intensity at a fixed wavelength is measured as the exciting wavelength is varied. If there is a layer-to-layer thickness variation, luminescence peaks will be observed at energies corresponding to each thickness, since recombination at a given energy can only originate from wells of a specific thickness. If there is no average thickness variation, however, identical peaks will be observed in the excitation spectra for the different luminescence energies at which the spectra are obtained. Apparently this is the case in the sample of fig. 9. From the peak width of 3.5 meV a thickness variation of less than one monolayer can be deduced.

Excellent agreement between theory and experiment was obtained for well widths between 8 nm and 15 nm for the model in which there was no average thickness variation, but with an interfacial roughening on the scale of one (0.28 nm thick) monolayer^[23]. In the experiment the wells of GaAs and the barriers of $\text{Al}_{0.24}\text{Ga}_{0.76}\text{As}$ were deposited at a substrate temperature of 690 °C. For very similar structures of GaAs and $\text{Al}_{0.21}\text{Ga}_{0.79}\text{As}$ barriers, deposited at 670 °C, the halfwidths of free exciton luminescence were studied as a function of well thickness^[24]. It was deduced

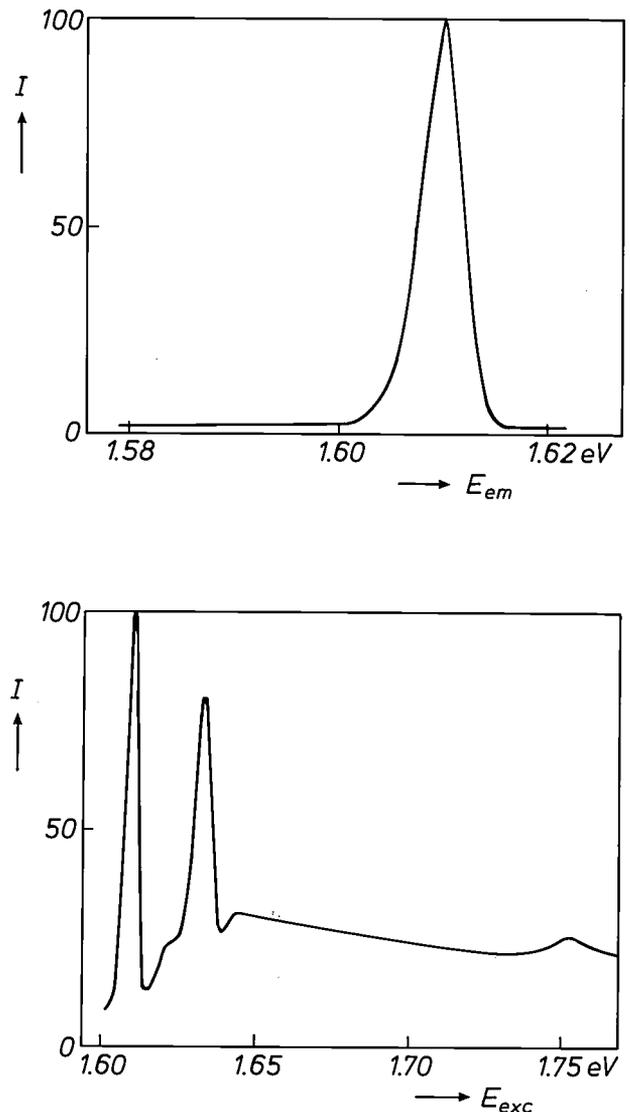


Fig. 9. a) Low temperature photoluminescence spectrum (relative intensity I as a function of the emission energy E_{em}) from a five-period quantum-well structure with 5.5-nm GaAs wells and 17.5-nm $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers. b) Corresponding excitation spectrum, obtained by measuring the relative photoluminescence intensity I as a function of the excitation energy E_{exc} at a fixed emission energy of 1.601 eV. The well-resolved peaks correspond to the formation of conduction electrons in the quantum state $n = 1$ and heavy and light holes in the $n = 1$ state. Note that the emission peak and the first excitation peak occur at the same energy. The peak width of 3.5 meV corresponds to a thickness variation which is less than one monolayer. The spectra were measured by P. J. Dobson and K. J. Moore.

[21] C. M. Garner, C. Y. Su, Y. D. Shen, C. S. Lee, G. L. Pearson, W. E. Spicer, D. D. Edwall, D. Miller and J. S. Harris, Jr., Interface studies of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ -GaAs heterojunctions, *J. Appl. Phys.* **50**, 3383-3389, 1979.

[22] L. P. Erickson and B. F. Phillips, Examination of MBE GaAs/ $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ superlattices by Auger electron spectroscopy, *J. Vac. Sci. & Technol. B* **1**, 158-161, 1983.

[23] C. Weisbuch, R. Dingle, A. C. Gossard and W. Wiegmann, Optical characterization of interface disorder in GaAs- $\text{Ga}_{1-x}\text{Al}_x\text{As}$ multi-quantum well structures, *Solid State Commun.* **38**, 709-712, 1981.

[24] H. Jung, A. Fischer and K. Ploog, Photoluminescence of $\text{Al}_x\text{Ga}_{1-x}\text{As}$ /GaAs quantum well heterostructures grown by molecular beam epitaxy; II. Intrinsic free-exciton nature of quantum well luminescence, *Appl. Phys. A* **33**, 97-105, 1984.

that more than 65% of the interface within the photo-excited area (100 μm diameter) was atomically sharp, and any roughness in the remaining area only extended over one monolayer.

In conclusion, the analytical methods with sufficient resolution have undoubtedly demonstrated the high structural and compositional perfection of the interfaces between some MBE-grown layers. For GaAs-Al_xGa_{1-x}As heterojunctions we can be confident that it is possible to prepare interfaces by MBE so that compositional changes occur over no more than one monolayer. They are free of extended defects.

Some properties and possible applications

Important properties of MBE-grown layers are the donor and acceptor concentrations and the free carrier concentration and mobility. These properties can be obtained by electrical characterization methods such as measurements of the Hall coefficient [26] and the electrical resistivity. The donor and acceptor concentrations can be derived by careful analysis of the temperature dependence of the free carrier concentration and mobility [26].

GaAs layers grown in our laboratory by MBE in a Varian GEN II equipment were found to have a low background acceptor level. This was determined by studying the electron mobility in 10-12 μm thick films lightly n-doped with silicon. The mobilities measured are typically $7.8 \times 10^3 \text{ cm}^2 \text{V}^{-1} \text{s}^{-1}$ at 300 K and $1 \times 10^5 \text{ cm}^2 \text{V}^{-1} \text{s}^{-1}$ at 77 K. The free electron concentrations allowing for depletion are $3 \times 10^{14} \text{ cm}^{-3}$ and the sum of the donor and acceptor concentrations

can be estimated to be about $7 \times 10^{14} \text{ cm}^{-3}$ [26]. Hence a value of $2 \times 10^{14} \text{ cm}^{-3}$ is derived for the background acceptor concentration. Unintentionally doped films of similar thickness are fully depleted and this implies that the background donor concentration is also about $2 \times 10^{14} \text{ cm}^{-3}$. Two-dimensional electron-gas structures prepared in the same MBE equipment were found to have mobilities as high as $3 \times 10^6 \text{ cm}^2 \text{V}^{-1} \text{s}^{-1}$ at 4 K, which also indicates how pure such samples are.

Multiple-quantum-well GaAs-AlGaAs injection lasers have been grown by MBE in our laboratory-constructed equipment. The GaAs well width was varied from 5.5 to 1.3 nm to study its effect on the wavelength of the laser emission [6]. The smallest well width (1.3 nm) corresponds to only about five monolayers of GaAs. The wells are separated by 8 nm wide Al_xGa_{1-x}As barriers.

In *fig. 10* the laser emission spectrum is given for quantum-well structures with various well widths and for a conventional laser with GaAs in the active region, which emits at about 880 nm. A decrease of the well width from 5.5 to 1.3 nm results in a progressive reduction in operating wavelength from 837 to 704 nm [6]. The emission at 704 nm is in the visible part of the spectrum. To our knowledge it is the shortest emission wavelength achieved with a room-temperature injection laser with only GaAs (and no Al) in the wells. Quantum-well lasers with such a short-wavelength emission are of interest in view of their possible use in optical information read-out systems.

The reduction of well width, particularly below about 3 nm, leads however to a strong increase of the threshold current for laser emission. This forms the only limitation for the realization of practical short-wavelength quantum-well lasers: as demonstrated, extremely thin wells with sufficiently abrupt interfaces can readily be obtained by MBE. Recently the influence of the number of wells on the threshold current was investigated, using a structure with a fixed waveguide [27]. It was shown that a small number of wells (preferably a single well) is favourable for obtaining lower threshold currents. For broad-area devices with a well of only 2.5 nm in a wide optical waveguide a threshold current density of about 1 kA/cm² was measured. This implies that it should

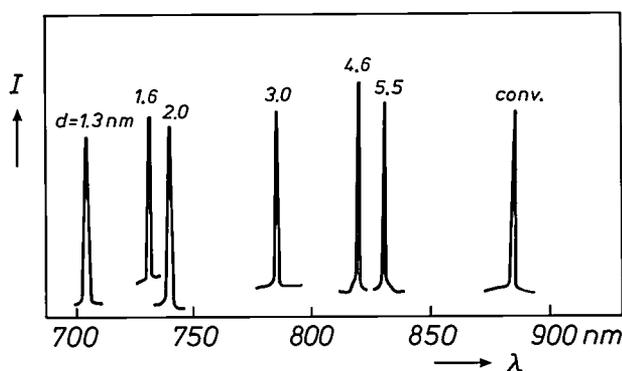


Fig. 10. Laser emission (relative intensity I plotted against emission wavelength λ) for a series of multiple quantum-well structures of GaAs wells and Al_xGa_{1-x}As barriers and for a conventional semiconductor laser. The structures show emission at progressively shorter wavelengths as the well thickness d is reduced. The structure with the thinnest well (1.3 nm) operates in the visible part of the spectrum, the emitting wavelength being 704 nm. The multiple quantum well lasers were made by P. Blood, E. D. Fletcher and P. J. Hulyer from layers grown by K. Woodbridge.

[26] The Hall coefficient represents the voltage induced by external electric and magnetic fields orthogonal to each other, which appears across a sample in a direction at right angles to both fields. The reciprocal of this coefficient is proportional to the free carrier concentration, and the ratio of the Hall coefficient to the electrical resistivity gives the free carrier mobility.

[26] G. E. Stillman and C. M. Wolfe, *Electrical characterization of epitaxial layers*, Thin Solid Films 31, 69-88, 1976.

[27] P. Blood, E. D. Fletcher, K. Woodbridge and P. J. Hulyer, Short wavelength (visible) quantum well lasers grown by molecular beam epitaxy, *Physica* 129B, 465-468, 1985.

be possible to obtain narrow stripes (5 μm wide) of production-type lasers with threshold currents below 50 mA.

In our MBE research and development at PRL invaluable assistance was given by Mr J. H. Neave. Many other members of the MBE group at PRL also contributed to the work described in this article, and there were many other contributions from members of other groups and other Philips colleagues.

Summary. The growth of thin films of GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ on GaAs substrates by molecular beam epitaxy (MBE) has been investigated extensively within the Philips Research Laboratories at Redhill. High-quality films and multilayer structures with controlled thickness and composition have been deposited. Detailed information has been obtained on the surface chemistry and the dynamics of growth. It is possible to produce abrupt interfaces between GaAs and $\text{Al}_x\text{Ga}_{1-x}\text{As}$ layers, where the compositional changes occur over no more than one monolayer. GaAs layers have been grown with low background impurity concentrations ($2 \times 10^{14} \text{ cm}^{-3}$) and this has enabled us to achieve high electron mobilities at low temperatures ($3 \times 10^6 \text{ cm}^2 \text{ V}^{-1} \text{ s}^{-1}$ at 4 K) in two-dimensional electron-gas structures. Multiple quantum-well structures prepared by growing very thin GaAs wells between $\text{Al}_x\text{Ga}_{1-x}\text{As}$ barriers showed laser emission at wavelengths down to 704 nm.

Silicon molecular beam epitaxy on GaP and GaAs

P. C. Zalm, C. W. T. Bulle-Lieuwma and P. M. J. Marée

Introduction

Most products emerging from today's semiconductor industry are manufactured from silicon. Only for applications such as light-emitting diodes, microwave devices^[1] and lasers for optical communication^[2], Compact Disc players^[3] and digital optical recording^[4] are III-V semiconductors employed. New and technologically interesting possibilities can be envisaged if single-crystal wafers containing both types of semiconductor can be prepared. Such monolithic wafers can be used for the fabrication of devices combining the high-frequency and optoelectronic features of III-V compounds with modern silicon VLSI technology. It might even be possible to integrate a laser diode, a detector (for the reflected laser light) and the circuits for drive and signal processing all on a single chip. But there are many difficulties to be mastered in achieving this goal.

Because of the similarity in crystal structure, silicon can be deposited epitaxially on a substrate such as the well-known III-V semiconductor GaAs. A problem is the difference in interatomic distances: the lattice constant of Si is about 4% smaller than that of GaAs. Such a mismatch may lead to incorporation of crystal defects in the layer or film, and these can adversely affect the electrical properties. Deposition of III-V material on Si poses an additional problem, connected with the presence of atomic steps and 'kinks' in the surface and with the different bonding of the two components to Si; see *fig. 1*. A slight bonding preference for one component (e.g. As) may induce 'anti-phase boundaries', dependent on the growth mechanism. This may also affect the electrical properties. It was therefore decided to investigate the growth of Si on III-V substrates. The intention was to obtain a better understanding of the effect of the lattice mismatch on the crystal quality and to assess the potential and limitations of these systems.

A large number of experiments were performed with GaP as the substrate material. This only has a

slight lattice mismatch (0.36%) with Si. We expected that a study of the combination of Si and GaP would shed some light on the onset of crystal defects due to lattice mismatch. Substrates that give a much larger mismatch are also of interest. A good representative of this category, which includes all the III-V materials of technological significance, is GaAs. The growth of silicon on GaAs was therefore extensively investigated as well.

A suitable growth technique is molecular beam epitaxy (MBE), essentially a controlled deposition by evaporation in ultra-high vacuum^[5]. It has been shown that MBE of Si films on an Si substrate below 800 °C gives material of excellent quality for device manufacture^[6]. Low deposition temperatures reduce diffusion across the interface, so that very abrupt changes in composition can be obtained. In addition, the decomposition of III-V compounds at higher temperatures limits the temperature range for epitaxial growth.

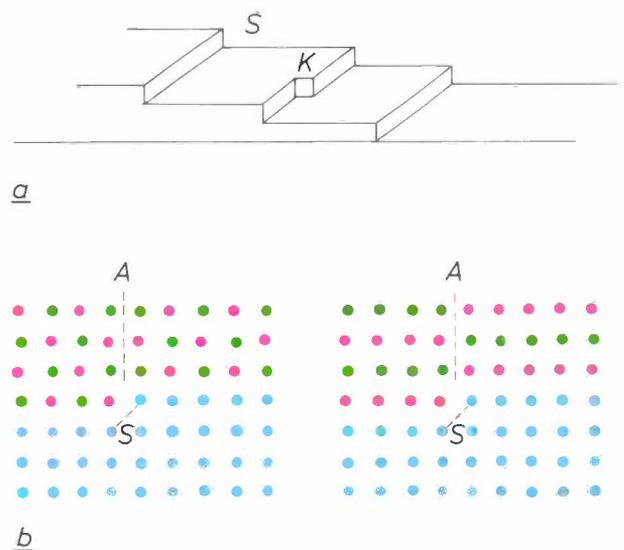


Fig. 1. *a*) Schematic representation of a surface of a single-crystal substrate showing steps *S* and a kink *K*. *b*) Two examples of atomic arrangements for a binary compound (e.g. GaAs) grown epitaxially on a substrate containing only one type of atom (e.g. Si). The difference in bonding leads to anti-phase boundaries *A* at steps *S* in the substrate surface.

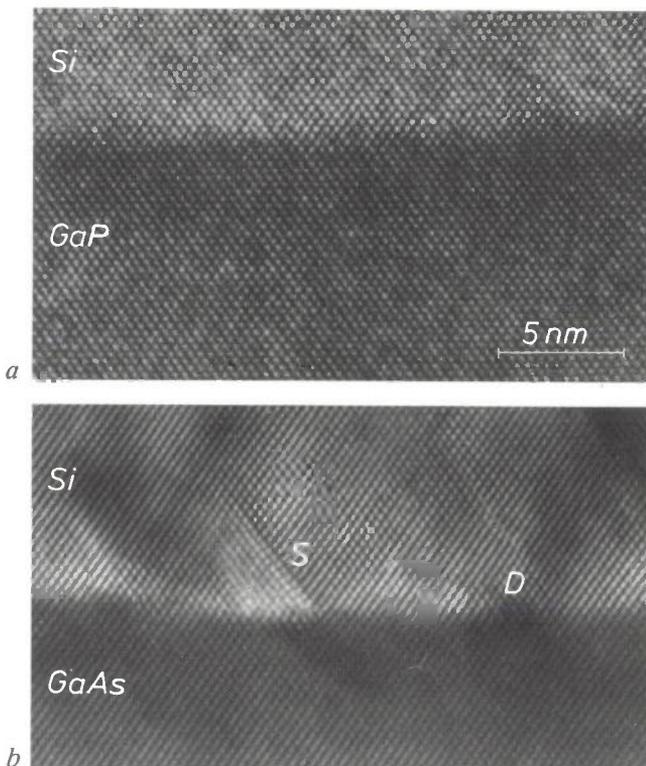


Fig. 2. High-resolution lattice images, made along the [110] direction by transmission electron microscopy (TEM), of MBE-grown Si on (001) GaP (a) and (001) GaAs (b). The interface region on GaP does not have any crystal imperfections. The interface region on GaAs has 60° dislocations *D* close to the interface and stacking faults *S* on (111) planes. The gradual interface transitions of about five monolayers can be attributed to an intermixing of elements or to the roughness of the substrates after the cleaning procedure.

The growth experiments described in this article were performed at the FOM Institute for Atomic and Molecular Physics (AMOLF) in Amsterdam, where a silicon MBE facility with associated analysis equipment is available for research purposes. Supplementary information on the properties of the films was obtained with various surface-analysis techniques at Philips Research Laboratories in Eindhoven.

Our studies showed that Si could be deposited epitaxially on (001) GaP, with no lattice defects for film thicknesses up to 75 nm. A cross-sectional lattice image made with a high-resolution transmission electron microscope is shown in fig. 2a. In films thicker than 75 nm misfit dislocations are introduced. The critical thickness for the onset of dislocation formation is much larger than predicted by current theory. In Si films deposited on (001) GaAs many crystal defects such as misfit dislocations and stacking faults are observed; see fig. 2b. The way in which the lattice mismatch induces the observed crystal imperfections is now fairly well understood. This may be important for the fabrication of thick perfectly epitaxial heterostructures. Before discussing the results in more detail, we shall first give a description of the MBE equipment and the procedure.

Experimental arrangement and procedure

Because of the fairly high temperature ($> 1400^\circ\text{C}$) required for significant evaporation, and the high reactivity of molten Si, effusion cells of the Knudsen type [6] cannot be used as a source. A good alternative is the electron-beam evaporator depicted in fig. 3. Here a high-power beam of electrons, with an energy of about 10 keV and a current of about 0.1 A, is focused by a magnetic field on to a Si 'slug', resulting in localized heating, melting and evaporation. The beam is scanned over an area of about 0.5 cm^2 .

A problem with such a source is the limited evaporation rate. Above a certain power input 'spluttering' occurs: droplets or clusters containing many Si atoms are ejected and deposited on to the substrate. Whereas atoms that arrive individually can diffuse over the surface until they arrive at a suitable lattice position, these droplets form immobile amorphous islands, because they do not disintegrate sufficiently well on impact.

To prevent this effect the deposition rate must be limited to less than about 2 nm/s, corresponding to a silicon flux of less than $10^{16}\text{ at.cm}^{-2}\text{s}^{-1}$. This in turn necessitates growth in ultra-high vacuum to prevent

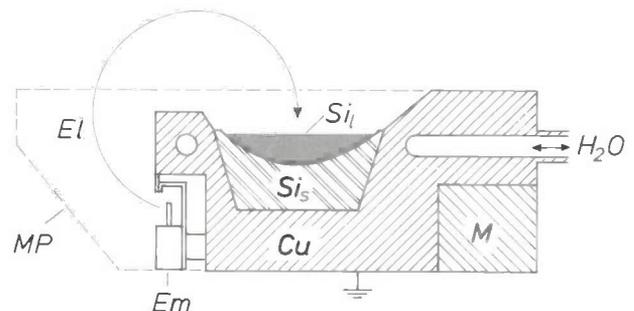


Fig. 3. Schematic representation of the electron-beam evaporator used as source for silicon MBE. Evaporation is produced by heating with an electron beam *El* emitted by the emitter assembly *Em* and deflected by a magnetic field on to a silicon slug. *Si_s* solid silicon, *Si_l* molten silicon, *M* magnet, *MP* magnet pole pieces. A deflection of 270° is used to avoid tungsten contamination from the filament. The copper holder *Cu* for the silicon slug is water-cooled to prevent melting and outgassing.

- [1] P. Baudet, M. Binet and D. Boccon-Gibod, Low-noise microwave GaAs field-effect transistor, *Philips Tech. Rev.* **39**, 269-276, 1980.
- [2] G. A. Acket, J. J. Daniele, W. Nijman, R. P. Tijburg and P. J. de Waard, Semiconductor lasers for optical communication, *Philips Tech. Rev.* **36**, 190-200, 1976.
- [3] Special issue 'Compact Disc Digital Audio', *Philips Tech. Rev.* **40**, 149-180, 1982.
- [4] K. Bulthuis, M. G. Carasso, J. P. J. Heemskerk, P. J. Kivits, W. J. Kleuters and P. Zalm, Ten billion bits on a disk, *IEEE Spectrum* **16**, No. 8 (August), 26-33, 1979.
- [5] Details of the MBE growth of III-V compounds have been given in the previous article in this issue by B. A. Joyce and C. T. Foxon.
- [6] Various aspects of silicon molecular beam epitaxy have been described in an extensive review article by Y. Ota, *Thin Solid Films* **106**, 3-136, 1983.

excessive contamination. Even at 10^{-8} Pa the rate of incidence of molecules from the residual gas is about $3 \times 10^{10} \text{ cm}^{-2} \text{ s}^{-1}$. If all these molecules were to 'stick' and remain on the surface the relative contamination level would be 3×10^{-6} , corresponding to an undesired doping level of $1.5 \times 10^{17} \text{ cm}^{-3}$. This would not only be disastrous for the electrical properties, but would also prevent perfect crystal growth. Although the sticking probability is generally much less than unity (typically 10^{-3}), great care has to be taken with the design and operation of an MBE system, so that these problems can be avoided.

For our experiments we used an earlier version of the present silicon-MBE system at the FOM institute AMOLF, designed by T. de Jong *et al.* [7] and shown schematically in *fig. 4*. It has the basic features found in many research-oriented systems, and consists of three linked vacuum chambers. Substrates mounted on holders are introduced twelve at a time into the load-lock annexe to the storage chamber; the twelve holders are placed on a carousel. After loading, the storage chamber is sealed, pumped down and baked for about ten hours at 150°C . The final pressure is a few times 10^{-7} Pa. Then the valve to the main chamber is opened and a holder with substrate is picked up and transferred with the aid of a magnetically coupled transfer rod. During the transfer the main chamber remains at low pressure (10^{-8} Pa). The desired substrate temperature is obtained by direct-current heating and checked with an infrared pyrometer to an accuracy of $\pm 25^\circ\text{C}$.

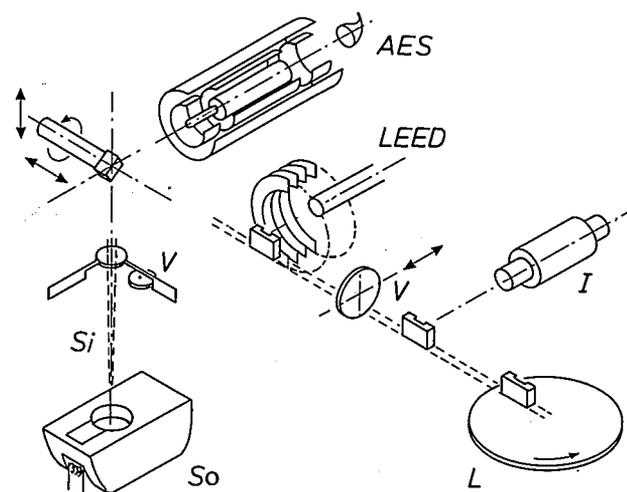


Fig. 4. Schematic diagram of an earlier version of the silicon MBE system at AMOLF. The system contains three vacuum chambers which can be connected by the gate valves *V*. A holder with a substrate is supplied via a load lock *L*. An ion gun *I* is used for cleaning the substrate surface. In the main vacuum chamber the surface is investigated by low-energy electron diffraction (LEED) and Auger electron spectroscopy (AES). The growth is started by introducing silicon vapour from the chamber with the Si source *So*. The diagram does not show the equipment for substrate heating and temperature measurement, the microbalance for growth-rate monitoring, or the quadrupole mass spectrometer for residual-gas analysis.

The source chamber is mounted below the main chamber. When the evaporator is turned on, residual gas adsorbed on the slug surface and the filament is released. This is pumped away and the valve to the main chamber is only opened when the source is stabilized at the desired evaporation or deposition rate, measured by a quartz-crystal oscillator microbalance. The pressure in the main chamber rises by less than an order of magnitude during the growth; the additional gas is mainly relatively harmless H_2 from the bulk of the Si slug.

Substrate preparation and analysis

Successful MBE growth of Si on GaP and GaAs requires careful attention to the preparation and cleaning of the substrate and the analysis of its surface. It is known that epitaxial Si films can be deposited on (001) Si from 200°C , whereas the deposition on (111) Si requires much higher temperatures ($> 600^\circ\text{C}$) [6]. It was therefore decided to use (001) substrates only. These are cut from liquid-encapsulated Czochralski-grown crystals, polished and degreased. The GaP substrates of thickness 1 mm are too small to be mounted directly on the holder, and are therefore attached by small amounts of indium to a $7 \text{ mm} \times 25 \text{ mm}$ silicon carrier. The indium provides a uniform thermal and electrical contact for the substrate heating, and does not affect the purity of the films. The GaAs substrates are 0.4 mm thick and can be cut immediately to the dimensions required for direct mounting in the holder.

Although the substrates are treated with the utmost care, their surface will always be covered by a thin film of oxide. Nor can slight contamination by carbon, from CO_2 in the atmosphere or from organic cleaning fluids, be avoided in practice. These contaminants are removed by sputter-etching with an Ar^+ beam at an energy of 600 to 800 eV and a total dose of about 10^{16} cm^{-2} . As ion bombardment makes the substrate surface amorphous, post-annealing is necessary to restore the crystallinity by solid-phase epitaxial regrowth [8]. This is done at 550°C for 30 min (GaP) and at 600°C for 90 min (GaAs).

A problem with this cleaning procedure is that the two components may be sputtered at different rates. In addition, at the temperatures necessary for reasonable regrowth rates ($> 1 \text{ nm/min}$) there will be some evaporation. Since the two components evaporate at different rates, there will be a depletion of one of the components (usually the group V element) at the surface. Consequently the surface, although crystalline, will have a non-stoichiometric composition after cleaning. In principle this deficiency can be corrected by supplying the 'missing' component through adsorption from the gas phase, or by first growing a

stoichiometric crystalline film of the III-V material on the substrate, as described in the previous article [5].

The effectiveness of the cleaning procedure was checked by Auger Electron Spectroscopy (AES) [9]: an energetic (a few keV) electron beam removes a strongly bound electron from a target atom, and the 'hole' created is 'filled' by a more weakly bound electron. The energy gained is used to eject another weakly bound (Auger) electron with an element-specific energy between 10 and 2000 eV. Determining the number of Auger electrons therefore provides information about the composition of the surface. This technique is surface-sensitive since the escape depth (without energy loss) of the Auger electrons is less than about 2 nm. Residual impurities with a surface concentration greater than 1 % can be detected. Our AES measurements revealed that the cleaned surfaces of GaP and GaAs contain less than 0.01 of a monolayer of carbon, nitrogen or oxygen.

Another surface-sensitive technique, low-energy electron diffraction (LEED) [9], was used to check the recrystallization of the surface. In LEED, a low-energy (< 200 eV) electron beam is reflected by the sample. Diffraction occurs, and if the surface is crystalline there will be constructive interference of the reflected beams in certain directions only. A phosphor screen placed in the path of the reflected and subsequently post-accelerated electrons will then emit light locally. Spot patterns thus obtained are characteristic of the ordering of the surface atoms. The distance between the various spots is inversely proportional to the inter-atomic distances on the surface in a specific direction.

Methods of investigating the films.

Once the substrates had been cleaned sufficiently, Si films were deposited with thicknesses ranging from a few tenths of a nanometer to a few hundred nanometres. The films were investigated by *in situ* analysis with the integrated AES and LEED equipment (fig. 4) and further characterization was performed outside the MBE system [10] [11].

The structural quality of the layers was investigated by transmission electron microscopy (TEM). The observed defects were correlated to the strain induced by the lattice mismatch with the substrate. The occurrence of strain and defects was studied further by measuring the back-scattering of light ions in various directions, which is sensitive to the exact positions of the atoms in the lattice. Additional information was obtained by Raman scattering, where the frequency shift of the scattered light is a measure of the stress caused by atomic displacements from bulk lattice sites. The variation of composition with depth was studied in some samples by secondary-ion mass spec-

troscopy (SIMS). In this technique the ions that are successively sputtered from the surface by an ion bombardment are analysed [9].

Not all of the Si films were grown on GaP and GaAs. Some were grown 'homo-epitaxially', i.e. on an Si substrate. Investigations by TEM revealed that these films were completely free of defects, thus verifying the quality of the MBE system and the procedure.

Results of *in situ* analysis

Figs 5 and 6 show LEED patterns from cleaned GaP and GaAs surfaces and from Si films deposited

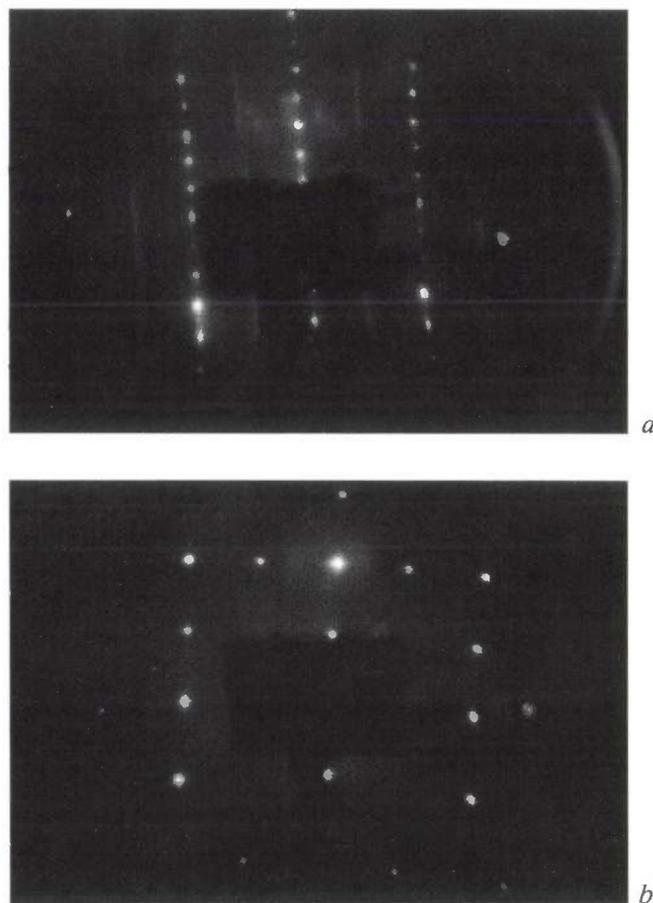


Fig. 5. LEED patterns from a (001) GaP substrate (a) and a silicon film (b). The patterns indicate that when silicon is deposited the surface reconstruction changes from (4×2) to (2×1) .

[7] A description of this MBE system has been given by T. de Jong, W. A. S. Douma, L. Smit, V. V. Korablev and F. W. Saris, *J. Vac. Sci. & Technol. B* 1, 888-898, 1983.

[8] T. de Jong, F. W. Saris, Y. Tamminga and J. Haisma, *Solid phase epitaxy of silicon on gallium phosphide*, *Appl. Phys. Lett.* 44, 445-446, 1984.

[9] See for example H. H. Brongersma, F. Meijer and H. W. Werner, *Surface analysis, methods of studying the outer atomic layers of solids*, *Philips Tech. Rev.* 34, 357-369, 1974.

[10] Si MBE on GaP has been reported previously by T. de Jong, W. A. S. Douma, J. F. van der Veen, F. W. Saris and J. Haisma, *Appl. Phys. Lett.* 42, 1037-1039, 1983.

[11] Some of the investigations on Si MBE on GaAs have been reported by P. C. Zalm, P. M. J. Marée and R. I. J. Olthoff, *Appl. Phys. Lett.* 46, 597-599, 1985.

on them. For (001) GaP the surface reconstruction is characterized as (4×2) , indicating that the periodicities along the $[110]$ and $[1\bar{1}0]$ directions are four times as large and twice as large as in the bulk material. The LEED pattern from (001) GaAs gives a $c(8 \times 2)$ reconstruction indicating an eight times two enlarged unit cell as compared with the bulk material, with a centre of symmetry. These results are typical of crystalline

Normal bulk diffusion of substrate atoms through the Si film is negligible at the temperatures used here, but there are other possible causes for the segregation observed. Detailed inspection of AES spectra and LEED patterns of various films have revealed two competing mechanisms. First there is surface diffusion of Ga and P or As from the sides and base of the substrate, an 'infinite' supply. This can be elim-

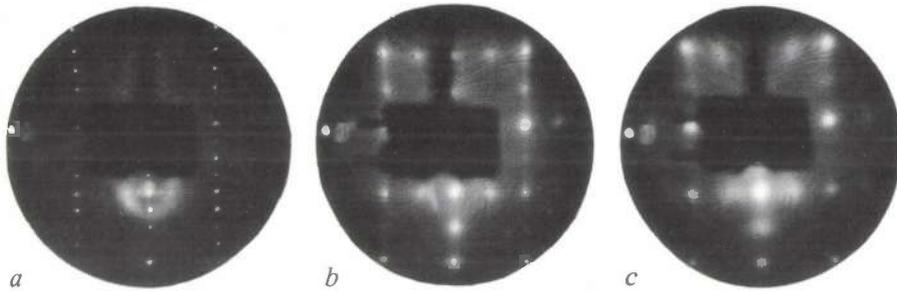


Fig. 6. LEED patterns from a (001) GaAs substrate (a) and Si films of thickness 1.4 nm (b) and 14 nm (c). On deposition, the surface reconstruction changes from $c(8 \times 2)$ to (2×1) ; the diffraction spots become more diffuse, indicating more crystal defects. The distance between the spots indicates that the 1.4-nm film has almost the same lattice constant as GaAs, whereas the lattice constant of the 14-nm film is 4% less and has the same value as for bulk Si.

GaP and GaAs surfaces after ion bombardment and annealing.

Only a few Si monolayers have to be deposited on GaP or GaAs for the LEED patterns to indicate a (2×1) reconstruction, which is typical of (001) Si surfaces. The distance between corresponding diffraction spots, however, remains unchanged. It appears that the interatomic spacing in the Si film is the same as in GaP or GaAs, i.e. the growth is coherent or pseudomorphic on the substrate. This can only be established reliably for Si on GaAs, of course; the lattices of Si and GaP are too similar (a difference of only 0.36%) for any change in LEED spacings to be noticed. For thicker films on GaAs the diffraction spots broaden and become more diffuse. This is a strong indication of increasing crystal imperfection. Moreover, the distance between the spots increases by 4%, corresponding to a decrease in interatomic spacing to the value for bulk Si. Finally, for very thick films the LEED spots (not shown here) become somewhat sharper again, suggesting an improvement in the crystal quality.

To determine the surface composition, Auger spectra were recorded for various Si coverages, see fig. 7. The Auger signals from the substrate atoms decay exponentially with thickness up to 1 nm, suggesting layer-by-layer growth rather than a formation of islands. At a greater thickness (≥ 10 nm), however, the Auger spectra still indicate the presence of Ga and P (or As) with a total surface concentration of $5 \times 10^{14} \text{ cm}^{-2}$, indicating some segregation.

inated by reducing the growth temperature. Secondly, there is a position-exchange reaction with the incident Si atoms at the surface. The orientation^[12] and shape of the LEED pattern for a silicon monolayer on GaAs can only be explained if the Si atoms are assumed to expel Ga atoms from their lattice positions, so that

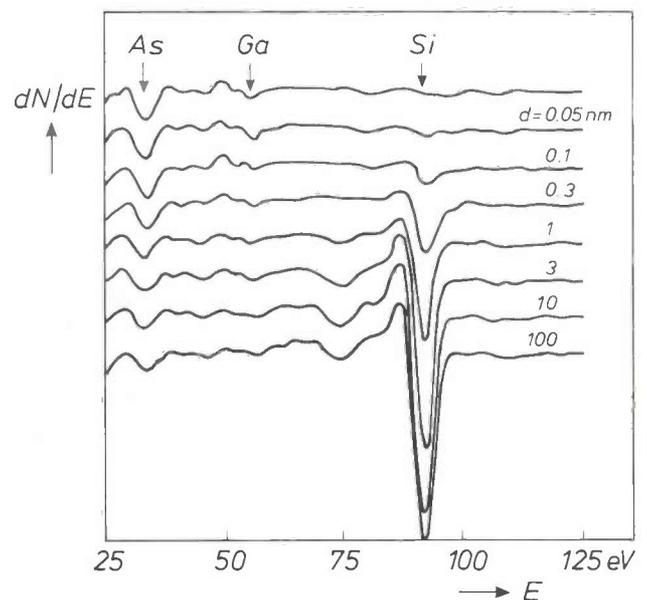


Fig. 7. Auger spectra of a cleaned (001) GaAs substrate and of Si films of thickness d deposited at 600 °C. The derivative dN/dE is shown in arbitrary units as a function of the electron energy E , where N is the number of Auger electrons detected. The arrows denote the peak positions corresponding to the characteristic Auger transitions of As (31 eV), Ga (55 eV) and Si (92 eV). For $0 < d < 1$ nm, the Ga and As peaks decrease exponentially with film thickness and the Si peak increases exponentially. When the thickness exceeds 10 nm the Ga and As peaks do not vanish.

they can bond directly with As atoms. This is rather surprising, since there are unoccupied As sites above the Ga atoms that form the upper layer of the Ga-enriched (001) surface. The exchange between Si and Ga atoms is corroborated by the Auger spectra of a monolayer, which showed that there was some Ga on top of the Si layer. A similar bonding preference has been reported for the deposition of germanium on GaAs [13].

Characterization by transmission electron microscopy

A technique such as LEED can give useful information about crystalline properties but is not very suitable for characterization of the crystal perfection because of its low sensitivity to defects. Transmission electron microscopy (TEM), however, can give more detailed structural information [14] [15]. Two examples have already been given in fig. 2, showing cross-sectional images of Si on GaP and GaAs obtained by high-resolution electron microscopy. Useful information can also be derived by making TEM 'plane-view' images at a much lower magnification. Before discussing some of the results, we shall first give a short description of the procedure for our TEM investigations.

Procedure

The samples are first made transparent to electrons. The preparation techniques required for the two kinds of imaging are very different.

For the preparation of samples for plane viewing, 3-mm discs are drilled from the wafer ultrasonically. The discs are reduced in thickness from the rear by jet-etching with a chlorinated-methanol etchant until the substrate has been completely removed from a small central area of diameter say 0.5 mm. The Si film can be examined here without interference from the substrate. At the edges the interface region can be studied.

For cross-sectional imaging, samples of {110} orientation are prepared [16]. Two strips less than 3 mm wide are cut from the wafer along one of the {110} planes, bonded together with the epitaxial films facing each other and then embedded in resin; see fig. 8. Thin slices are cut, mechanically polished and reduced in thickness by Ar-ion milling on a rotating sample holder at grazing incidence (about 10° to the surface) and a low voltage (≈ 4 kV) to make the surface smooth and minimize the formation of an amorphous top layer.

Images were obtained with a Philips EM 420 ST microscope operating at 120 kV. The contrast observable in TEM images has been described in an earlier

article in this issue [17]. An area of GaP or GaAs appears dark compared with a silicon area because of the difference in atomic scattering. The plane-view images are made in 'bright field' (with the transmitted beam) or in 'dark field' (with one diffracted beam) under dynamic diffraction conditions. Defects can then be observed because of local variations in the diffraction conditions (amplitude contrast). The cross-sectional images are made with the electron beam along the [110] direction, so that there is strong simultaneous excitation of many diffracted beams. Recombination of the transmitted beam and several of the diffracted beams results in a high-resolution image (phase contrast).

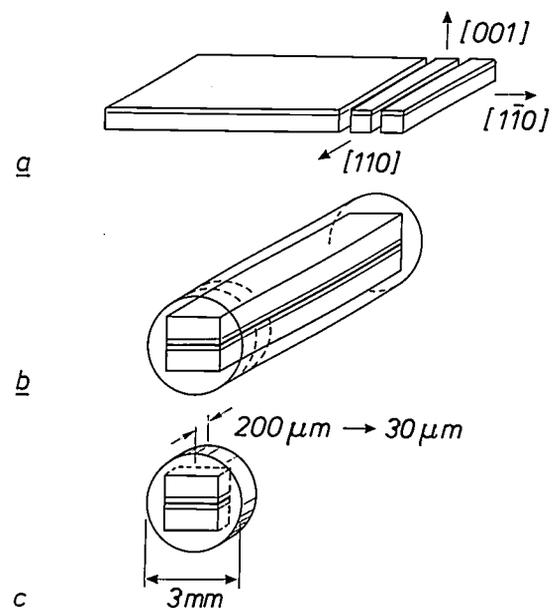


Fig. 8. Schematic illustration of the preparation of {110} samples for cross-sectional TEM imaging. Two bars are cut along one of the {110} planes (a). The bars are bonded together with the epitaxial films facing each other and are embedded in resin (b). They are then sawn into thin slices (c). The slices are further reduced in thickness by Ar ion milling on a rotating sample stage.

[12] A. J. van Bommel and J. E. Crombeen, Experimental determination of the correlation between the LEED pattern and the Ga-As bond vectors in the surface of GaAs (001), *Surf. Sci.* 57, 437-440, 1976.

[13] B. J. Mrstik, LEED and AES studies of the initial growth of Ge epilayers on GaAs (100), *Surf. Sci.* 124, 253-266, 1983.

[14] M. P. A. Vieggers, C. W. T. Bulle-Lieuwma, P. C. Zalm and P. M. J. Marée, Misfit dislocations in epitaxial layers of Si on GaP (001) substrates, *Mater. Res. Soc. Symp. Proc.* 37, 331-336, 1985.

[15] C. W. T. Bulle-Lieuwma, P. C. Zalm and M. P. A. Vieggers, Characterization of MBE grown Si on (001) GaAs by transmission electron microscopy, *Proc. Microscopy of Semiconductor Materials Conf.*, Oxford 1985 (Inst. Phys. Conf. Ser. 76, section 4), pp. 123-128.

[16] C. W. T. Bulle-Lieuwma and P. C. Zalm, Suppressing of surface topography development in ion-milling of semiconductors, to be published in *Surface and Interface Analysis*, Vol. 10, 1987.

[17] M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft, this issue, pp. 133-142.

Si on GaP

Plane viewing of Si on GaP gives almost featureless images for thicknesses up to 75 nm, indicating film growth with no misfit dislocations. It appears that the Si lattice is strained to match GaP, giving a coherent interface. At thicknesses above 75 nm, misfit dislocations are introduced. This is energetically more favourable when the strain energy becomes too high. An example of a TEM image with dislocations is given in *fig. 9*. This shows a rectangular network of misfit dislocations near the interface, and stacking faults intersecting the Si film and also forming a rectangular pattern.

The critical thickness observed (75 nm) for dislocation formation is considerably larger than the value of 14 nm derived from J. H. van de Merwe's theory assuming a thermodynamic equilibrium between an array of misfit dislocations along the interface and the strained layer [18]. A barrier to dislocation formation appears to exist. An understanding of its origin may be important for the production of perfect thick epitaxial heterostructures. We therefore studied the lattice defects of films thicker than 75 nm in more detail.

A complete characterization reveals that the dislocations are of the 60° type: 'Burgers vectors', describing the direction and magnitude of the atomic displacements, are at an angle of 60° to the dislocation lines in [110] and [110] directions in the (001) interface; see *fig. 10a*. This is a common configuration for dislocations nucleating at the surface of a strained layer and gliding towards the interface along {111} slip planes [19]. These intersect the interface along the [110] and [110] directions of the dislocation lines.

The presence of stacking faults (*fig. 9*) can be explained by considering the atomic arrangement in the slip planes. This is shown in *fig. 10b* for a face-centred cubic (fcc) lattice with a biaxial strain field in the (001) plane. The resolved shear stress $\vec{\tau}$ on a (111) plane is in the [112] direction. The unit displacement described by the Burgers vector \vec{b} is achieved by two movements \vec{b}_1 and \vec{b}_2 , producing 90° and 30° Shockley partial dislocations [20]. The force on the 90° dislocation ($\propto \vec{b}_1 \cdot \vec{\tau}$) is twice that on the 30° dislocation ($\propto \vec{b}_2 \cdot \vec{\tau}$). The principle just outlined is the same for all four {111} planes and is applicable to Si [21]. The special geometry in Si on (001) GaP therefore ensures that the 90° dislocation, which experiences the larger force, nucleates first. It moves to the interface and produces a stacking fault there. The defect structure thus initially consists of stacking faults bounded by 90° dislocations at the interface. As the film grows, the 30° dislocations also nucleate because of the stacking-fault energy. This means that the stacking faults disappear, giving perfect 60° dis-

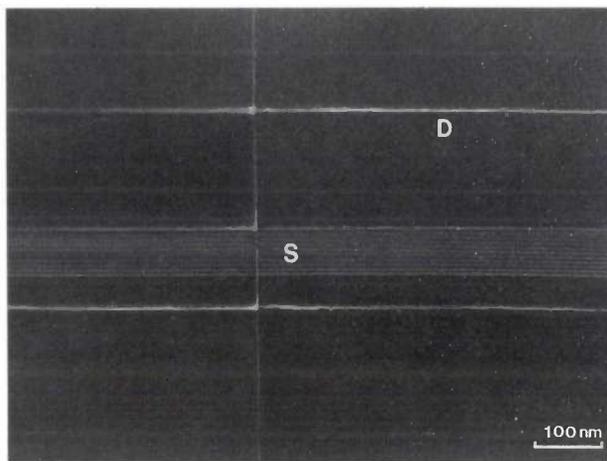


Fig. 9. Transmission electron micrograph of a planar section of an MBE-grown Si film 100 nm thick on (001) GaP, showing dislocations *D* along the interface and stacking faults *S* intersecting the Si film.

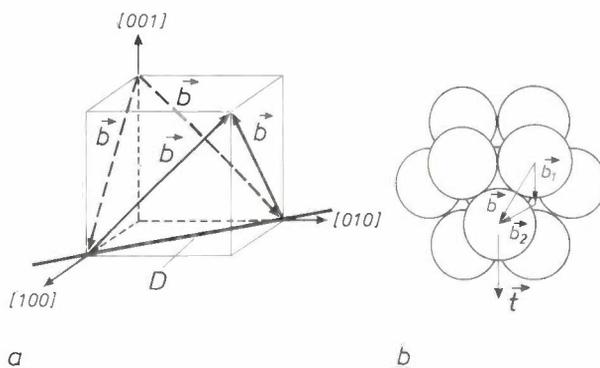


Fig. 10. *a*) Schematic illustration of the four possible Burgers vectors \vec{b} of a 60° dislocation *D*, along one of the two <110> directions on the (001) interface. The vectors, which describe the magnitude and direction of the atomic displacement, are at an angle of 60° to the dislocation line. *b*) Schematic illustration of the atomic arrangement in one of the {111} planes of a face-centred cubic lattice, for a (001) film with a biaxial strain field. The resolved shear stress $\vec{\tau}$ on a (111) plane is in the [112] direction. The unit displacement \vec{b} , which results in a 60° dislocation, is produced by \vec{b}_1 and \vec{b}_2 , giving 90° and 30° Shockley partial dislocations. The vectors are given by: $\vec{b}_1 = \frac{1}{6}a[112]$, $\vec{b}_2 = \frac{1}{6}a[2\bar{1}\bar{1}]$ and $\vec{b} = \frac{1}{2}a[101]$, where *a* is the lattice constant.

locations. In *fig. 9* these form a rectangular network, while the stacking faults can be associated with the situation where the 30° dislocations have not yet nucleated.

An example of incompletely formed (dissociated) 60° dislocations is shown in the high-resolution image

[18] J. H. van der Merwe and C. A. B. Ball in: J. W. Matthews (ed.), *Epitaxial growth*, part B, Academic Press, New York 1975.
 [19] J. W. Matthews in: J. W. Matthews (ed.), *Epitaxial growth*, part B, Academic Press, New York 1975.
 [20] A. H. Cottrell, *Dislocations and plastic flow in crystals*, Oxford Univ. Press, London 1953.
 [21] J. Hornstra, *Dislocations in the diamond lattice*, *J. Phys. & Chem. Solids* 5, 129-141, 1958.
 [22] A. Gomez, D. J. H. Cockayne, P. B. Hirsch and V. Vittek, *Dissociation of near-screw dislocations in germanium and silicon*, *Phil. Mag.* 31, 105-113, 1975.

of *fig. 11*. Close to the interface, the dislocations are found to have an average separation of about 10 nm. This is substantially larger than the equilibrium distance at 880 °C [22]. This may indicate that the local strains in the interface region differ considerably from the mean strain in the bulk of the film.

Finally, we note that the above geometrical effects depend strongly on the crystallographic orientation of the strained film and on the sign of the stress field. For example, in a compressive field it is not the 90° dislocations but the 30° dislocations that will have to nucleate first because of the different atomic arrangement. Since the 30° dislocations are subject to a smaller force from the shear stress, there may be a greater barrier to their nucleation than with a tensile force. Thus, films free from epitaxial defects may be thicker for compressed GaP on (001) Si than for strained Si on (001) GaP. We can assume that the nucleation of the 90° dislocations will immediately follow the nucleation of the 30° dislocations, because they are subject to a larger force from the shear stress. This means that the defect structure would consist only of perfect 60° dislocations.

Si on GaAs

Plane-view images of Si on (001) GaAs reveal Moiré fringes; see *fig. 12*. This indicates a small difference in atomic spacing between film and substrate, with no rotational misorientation on average. The periodicity corresponds to the 4% lattice mismatch between Si and GaAs. The fringes are bent locally and sometimes interrupted. Most of these distortions can be attributed to the presence of dislocations. In regions where the substrate has been removed completely, a high density of dislocations is observed. Stereomicroscopy reveals that they extend from the surface to the interface and back, with only short segments along the interface. The dislocations in the Moiré fringes of *fig. 12* are more pronounced for diffraction at (220) planes than for diffraction at (400) planes. This means that they are preferably located in (110) planes. When gliding on (111) planes, they have to move along $\langle 112 \rangle$ directions. As shown in *fig. 13*, these dislocations can be clearly seen in cross-sectional images. As expected, their directions are very close to the $\langle 112 \rangle$ directions.

Investigations of the interface region (*fig. 2*) indicate that thin Si films grow coherently on GaAs in spite of the 4% lattice mismatch. The defect structure observed consists of dislocations close to the interface and stacking faults in (111) planes with the density increasing towards the interface. Because LEED indicated non-pseudomorphic growth above 1.4 nm, we attribute this defect structure to the relaxation of the misfit stress. Below 1.4 nm there is a planar tensile stress in the Si

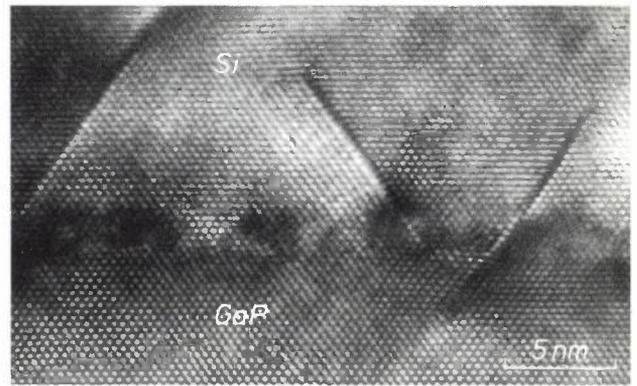


Fig. 11. High-resolution electron microscopy (HREM) cross-sectional image along the [110] direction of MBE-grown Si on (001) GaP, showing dissociated 60° dislocations close to the interface.

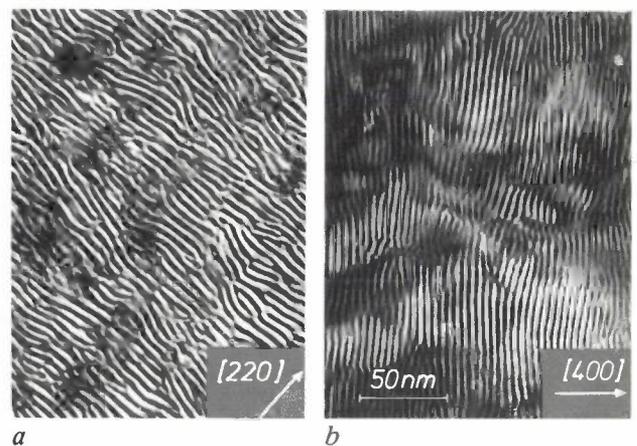


Fig. 12. Bright-field TEM images of a planar section of an Si film 200 nm thick on (001) GaAs for diffraction at (220) planes (*a*) and (400) planes (*b*). Both images show Moiré fringes, approximately perpendicular to the diffraction vector in the [220] and [400] directions. The periodicity corresponds to the 4% lattice mismatch. The distortions of the fringes are due to the presence of dislocations, which are more pronounced in (*a*) than in (*b*).

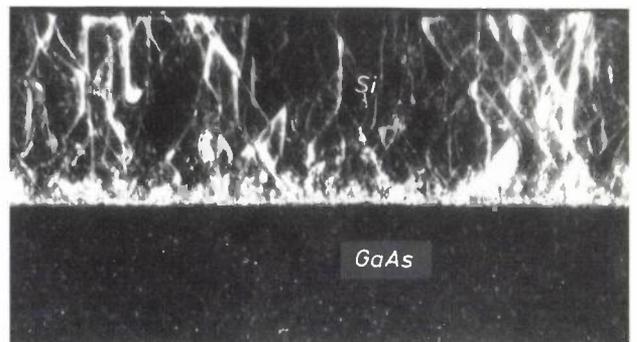


Fig. 13. Dark-field cross-sectional TEM image along the [110] direction of 200-nm MBE-grown Si on (001) GaAs. The Si film has a high density of dislocations, whereas no defects are observed in the GaAs substrate. Many dislocations run in $\langle 112 \rangle$ directions.

film, leading to a strain equal to the lattice mismatch. This stress can initially cause a homogeneous tetragonal distortion, giving rise to pseudomorphic growth. Above 1.4 nm relaxation occurs by the formation of 60° dislocations in the same way as for considerably thicker Si films on GaP. The larger lattice mismatch

gives a stronger strain field, however, and hence a higher dislocation density. In plane-view images we only see short segments in the interfacial plane, indicating that the glide process has been disturbed by interactions between the dislocations.

Films of thickness 200 nm grown at 600 °C and films of thickness 25 nm grown at 500, 550 and 600 °C are found to have about the same lateral defect density. This indicates that there is no thermodynamic equilibrium between strain and dislocations. It points either to a complete relaxation of misfit strain below 25 nm or to the existence of a barrier to dislocation formation. In the latter case extended defects can be avoided, e.g. by preventing the occurrence of nucleation sites or by reducing the growth temperature, which reduces the mobility of the dislocations.

Further characterization

For perfect epitaxial growth on GaP, the Si atoms will be displaced laterally with respect to their natural lattice positions by 0.36% at room temperature and 0.46% at the actual growth temperature of 570 °C. This parallel elongation leads to a contraction perpendicular to the interface — the Poisson effect; see *fig. 14*. A direct consequence of the resulting tetragonal deformation is a tilt $\Delta\theta$ in the non-normal crystallographic axes, e.g. the [011] axis. Measurements of $\Delta\theta$ (by ion back-scattering) or the perpendicular lattice spacing (by X-ray diffraction) give valuable information about the strain [23].

Unfortunately, X-ray diffraction can only be used with films thicker than 0.5 μm . Measurements of $\Delta\theta$ are possible, however, if ion back-scattering is used in conjunction with ion blocking. *Fig. 15* shows the principles of this technique for measuring $\Delta\theta$ for the [011] axis, for which the largest shift was expected. Light positive ions (H^+ or He^+) with an energy of 100–200 keV are incident on the surface. Upon penetration, an ion only changes direction when it passes so closely (10^{-5} nm) to an atomic nucleus that scattering occurs. The ion loses energy because of the scattering process and the interaction with electrons as it travels through the film. The spectrum of a large number of back-scattered ions therefore gives information about compositions and concentration profiles to a depth of about 100 nm. A depth resolution of 0.5 to 1 nm can be obtained with an electrostatic analysing detector with a high energy resolution [24].

Directions parallel to rows of atoms are ‘forbidden’ to outgoing ions: the atoms nearer to the surface ‘block’ the ions scattered by deeper atoms. In these directions there is a dip in the back-scattering spectrum. For Si on (001) GaP the [001] blocking direc-

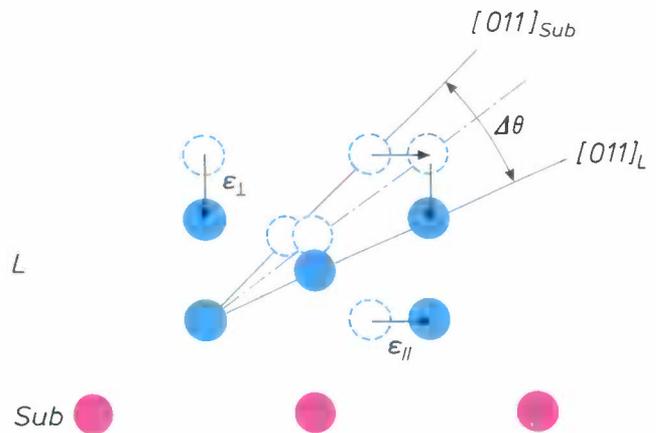


Fig. 14. Schematic illustration of the tetragonal lattice distortion in a film L with a smaller lattice constant than the substrate Sub . The parallel elongation by the coherency strain $\epsilon_{||}$ induces a perpendicular stress ϵ_{\perp} leading to a contraction (Poisson effect). The resulting tetragonal distortion corresponds to a tilt $\Delta\theta$ in the [011] direction with respect to the substrate.

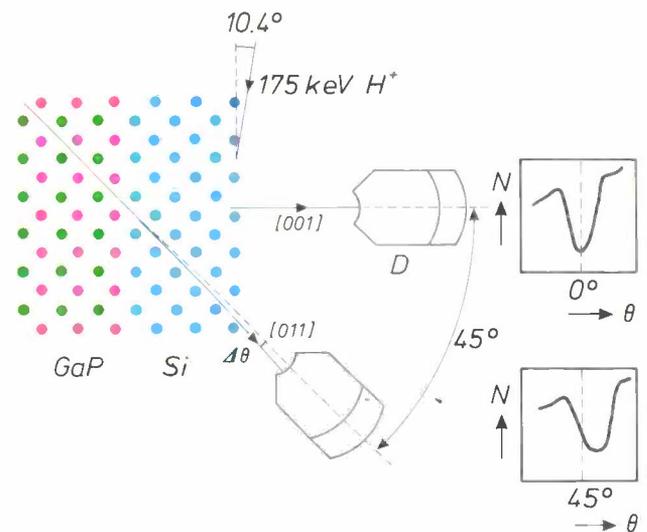


Fig. 15. Schematic view of the ion-blocking experiment, showing the scattering geometry and the measurement procedure. Hydrogen ions with an energy of 175 keV are incident on the Si surface at an angle of 10.4° . The yield of back-scattered ions is measured with a detector D at an angle θ to the normal to the surface. At $\theta = 0^\circ$, the blocking of scattered ions in the [001] direction gives a minimum in the back-scattering intensity yield N . For an undistorted lattice, the blocking in the [011] direction would occur exactly at $\theta = 45^\circ$. A tetragonal distortion with an angular change of $\Delta\theta$ will give a corresponding shift in the intensity minimum.

tions of film and substrate coincide. The tetragonal distortion, however, shifts the [011] blocking direction of the Si film with respect to the substrate [011] blocking direction, which is at an angle of 45° to the normal to the surface. The measured angular shift $\Delta\theta$ is a measure of the strain in the film.

To ensure a reliable determination of $\Delta\theta$ we first checked that a rotation of the detector through 45° corresponded to the angle between the [001] and [011] blocking directions for an unstrained Si crystal. We made sure that the surface oxide layer and interac-

tions between the ions and the surface atoms did not affect the position of the dip. This was achieved by energy selection of the ions back-scattered from depths of 3 to 11 nm. To be certain that the sample was not misaligned, we checked that the $[0\bar{1}1]$ dip, at about 45° in the opposite direction, was shifted by the same amount.

Typical results of ion-scattering measurements are given in *fig. 16*. The back-scattering spectrum of bulk Si has the expected $[001]$ and $[011]$ dips at exactly 0° and 45° . The spectrum of a 30-nm strained Si film on GaP has a shift in the $[011]$ dip: $\Delta\theta = 0.26 \pm 0.04^\circ$. For thicker films (200 and 400 nm) a smaller value of $\Delta\theta$ is found, a clear proof of strain relaxation. The results agree well with those obtained by TEM determinations of the dislocation density. However, the measured shift clearly differs from 0.19° , the theoretical value based on the elasticity constant of bulk Si. This is rather surprising, since X-ray diffraction in films thicker than $0.5 \mu\text{m}$ always confirms the elasticity-theory predictions exactly. It would appear that this theory does not apply to films as thin as those studied here.

Information about the strain was also obtained by Raman scattering of a monochromatic light beam, e.g. from the 514.5-nm line of an argon laser. The incident photons lose some of their energy by inducing a lattice vibration. The resulting shift in the frequency of the scattered light is related to the stress caused by atomic displacements. The difference from bulk Si is therefore a relative measure of the strain in epitaxial Si.

The results of strain determinations with TEM, ion back-scattering and Raman scattering all agree well; see *fig. 17*. This shows both the parallel elastic strain and the part of the lattice mismatch that is taken up by dislocations, as a function of the film thickness. It is found with all three methods that the strain is approximately equal to the lattice mismatch for films up to about 75 nm thick. This is very different from the equilibrium theory mentioned earlier [18], which gives 14 nm as the critical thickness for dislocation formation.

For Si on GaAs, information about the quality of very thin epitaxial films has been obtained by measuring the back-scattering yield as a function of energy [11]. It was found that a 5-nm Si film grew evenly, with no island formation at first.

Films have also been characterized by the back-scattering of light positive ions (He^+) incident on the surface at high energy (2 MeV) (Rutherford back-scattering, RBS). If the ions are incident in an 'open' crystallographic direction ('channelling'), the scattering yield increases as the number of ions occupying regular lattice sites decreases [25]. The yield ratio for chan-

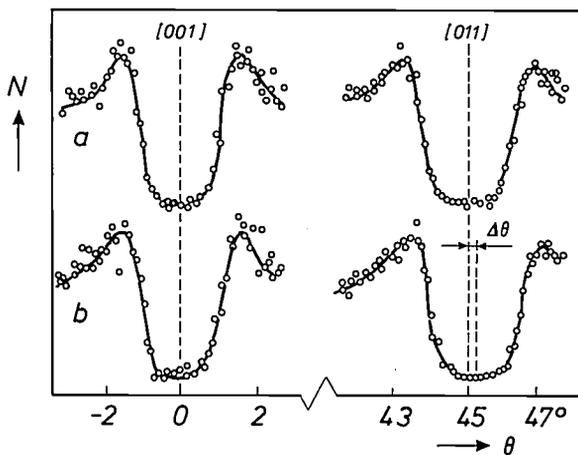


Fig. 16. Measured angular profiles near the $[001]$ and $[011]$ blocking minima for bulk Si (a) and for 30-nm MBE-grown Si on (001) GaP (b). The measured yield N is given in arbitrary units as a function of the angle θ to the normal to the surface. For bulk Si, the $[001]$ and $[011]$ minima are located at exactly $\theta = 0^\circ$ and $\theta = 45^\circ$, as would be expected in the absence of strain. For Si on GaP, a shift of $\Delta\theta = 0.26^\circ$ is measured for the $[011]$ minimum, as a result of the strain-induced tetragonal distortion.

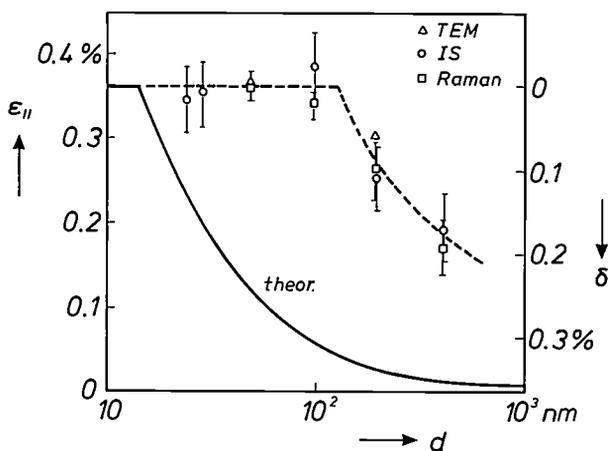


Fig. 17. Results of strain calculations from equilibrium theory and of determinations by transmission electron microscopy (TEM), ion back-scattering (IS) and Raman scattering, for MBE-grown Si on (001) GaP. The elastic strain $\epsilon_{||}$ and the part of the lattice mismatch that is taken up by dislocations (δ) are plotted against the film thickness d . Ion back-scattering and Raman spectroscopy were used to derive $\epsilon_{||}$ from the shift in the $[011]$ direction and the frequency shift, and TEM was used to determine δ from the ratio of the parallel Burgers-vector component and the mean separation between the dislocations. The sum of $\epsilon_{||}$ and δ is equal to the lattice mismatch (0.36%), so that a determination of either also gives the other. The results of the three experimental methods agree well with each other. The critical thickness for dislocation formation (75 nm) is however much larger than predicted by the equilibrium theory (14 nm).

[23] P. M. J. Marée, R. I. J. Olthof, J. W. M. Frenken, J. F. van der Veen, C. W. T. Bulle-Lieuwma, M. P. A. Vieggers and P. C. Zalm, Silicon strained layers grown on GaP(001) by molecular beam epitaxy, *J. Appl. Phys.* 58, 3097-3103, 1985.

[24] See for example J. F. van der Veen, Ion beam crystallography of surfaces and interfaces, *Surf. Sci. Rep.* 5, 199-288, 1985.

[25] See for example W. K. Hofker and J. Politiek, Ion implantation in semiconductors, *Philips Tech. Rev.* 39, 1-14, 1980.

nelling and random incidence has a lower value with increasing crystal perfection; the minimum value for pure GaP and GaAs substrates is 3%. For Si on GaP the same value has been derived from the Si contribution to the RBS spectra, which indicates perfect crystalline ordering. With Si on GaAs there is appreciable dechannelling; see *fig. 18*. Near the Si surface the amount of dechannelling is encouragingly low (9%), indicative of good but by no means perfect crystalline

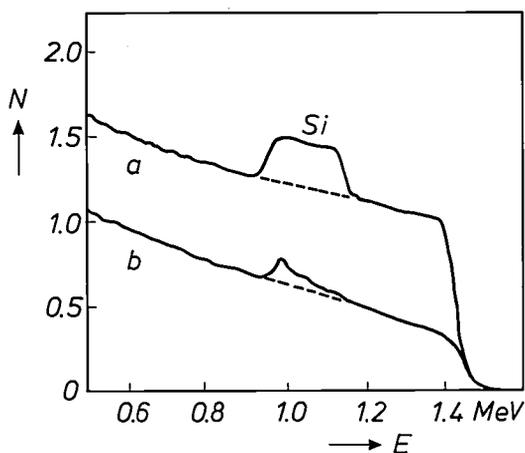


Fig. 18. Rutherford back-scattering spectra of a 500-nm Si film grown by MBE at 600 °C on (001) GaAs. The back-scattering of He⁺ ions incident with an energy of 2 MeV on the surface is measured at an angle of 165° to the incident beam. The back-scattering yield N is given in relative units as a function of the energy E of the scattered ions. Curve *a* is obtained with incidence in a random direction, curve *b* with incidence along the 'open' [001] direction. By comparing the Si contribution in the two spectra on the high- and low-energy sides, the amount of dechannelling on the Si surface or the Si-GaAs interface can be estimated.

ordering. The dechannelling increases considerably (to about 40%) however towards the Si-GaAs interface. The RBS results confirm the information gathered from our TEM investigations.

Information about the doping level in the Si films can be obtained by secondary-ion mass spectroscopy (SIMS) [9]. In this method an ion beam with an energy of a few keV sputters ionized particles from the surface. If the sputtered particles are accelerated to a well-defined energy and subjected to electric and magnetic deflection, they can be analysed by mass. The variation of the composition with depth can now be studied, since the surface is continuously 'peeled off' during the sputtering.

SIMS experiments on Si films on GaP (grown at 550 °C) and on GaAs (grown at 600 °C) reveal not only a surface region containing Ga and P or As, but also a high doping level in the bulk: 10^{18} cm⁻³ Ga and twice as much P or As. When grown on GaP at 400 °C the doping with P was considerably lower. The Ga con-

tent, on the other hand, remained the same, confirming an exchange of Ga and Si atoms suggested by LEED and Auger observations.

Prospects

The investigations have provided a better understanding of the formation of dislocations induced by lattice mismatch. We have seen that Si films can be grown on GaP to a thickness of about 75 nm, with no indication of misfit dislocations in TEM images. With GaP grown on Si it would be possible to have even thicker defect-free films, because dislocations are less easily initiated in compressed films. The occurrence of anti-phase disorder (*fig. 1*) could for example be prevented with a substrate orientation such as (211), offering two significantly different types of bonding sites [26]. The Ga and P atoms will then automatically bond to different sublattices, with the same sublattice pairing over the entire surface.

The tetragonal distortion in strained Si films on GaP has an interesting subsidiary aspect. From band-structure calculations it can be shown that the electron mobility must be twice that in bulk Si at room temperature. This offers possibilities for faster semiconductor switching devices, provided that the accompanying increase in leakage current (by a factor of about ten) poses no additional problems. The positive effect of a tensile stress on the electron mobility has already been demonstrated for thermally treated Si films on SiO₂-coated substrates and on sapphire [27]. In these films, however, the high defect density due to the large lattice mismatch prevents the mobility from exceeding that in bulk Si.

The high level of Ga and P doping in Si films on GaP may have undesirable effects on their electrical properties. Future work must therefore be directed at eliminating this 'background' doping. Some suggestions are:

- Coating the sides and base of the substrate with a material such as Si₃N₄, which acts as a barrier to the surface diffusion.
- A different method of substrate preparation that gives an As-stabilized surface.
- Reducing the temperature to 300-400 °C during growth.

[26] S. L. Wright, H. Kroemer and M. Inada, Molecular beam epitaxial growth of GaP on Si, *J. Appl. Phys.* **55**, 2916-2927, 1984.

[27] B.-Y. Tsaur, J. C. C. Fan and M. W. Geis, Stress-enhanced carrier mobility in zone melting recrystallized polycrystalline Si films on SiO₂-coated substrates, *Appl. Phys. Lett.* **40**, 322-324, 1982;

Y. Kobayashi, M. Nakamura and T. Suzuki, Effect of heat treatment on residual stress and electron Hall mobility of laser annealed silicon-on-sapphire, *Appl. Phys. Lett.* **40**, 1040-1042, 1982.

The crystal quality of Si on (001) GaAs will never be perfect. The dislocation density is enormous, especially near the interface. In thicker films the quality improves considerably towards the surface. Further investigation is required to find out if such defects could affect the operation of devices. A technology-oriented evaluation requires detailed electrical characterization as well as improved growth facilities.

On GaAs, very thin Si films (< 5 nm) can be grown with a complete coverage of the substrate surface. Such thin films are therefore potentially useful for ohmic contacts on GaAs, where a high defect density is no problem and a high background-doping level can even be an advantage.

Many colleagues at Philips Research Laboratories and the FOM Institute for Atomic and Molecular Physics have contributed to the work described here.

Summary. Thin silicon films have been epitaxially grown on GaP and GaAs by molecular beam epitaxy (MBE). Films and interfaces have been investigated by surface analysis and transmission electron microscopy. Both the surface and the bulk of the films are contaminated by atoms from the substrate; this is attributed mainly to surface diffusion and exchange between Si and Ga atoms. On GaP, films thinner than 75 nm are subject to tensile stress because of the lattice mismatch (0.36%), without the appearance of crystal imperfections. Thicker films contain dislocations and stacking faults, related to strain-induced atomic rearrangements. The critical thickness for dislocation formation is significantly larger than is predicted from elasticity theory. On GaAs, the lattice mismatch of 4% gives a greater reduction in crystal quality.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	A
Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	B
Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	E
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	R

C. Ronse	B	A strong chord property for 4-connected convex digital sets	Comput. Vision, Graphics & Image Process. 35	259-269	1986
H. C. de Graaff & W. J. Kloosterman	E	New formulation of the current and charge relations in bipolar transistor modeling for CAD purposes	IEEE Trans. ED-32	2415-2419	1985
R. L. Maresca	N	A general method for designing low-temperature drift, high-bandwidth, variable-reluctance, position sensors	IEEE Trans. MAG-22	118-123	1986
E. H. L. Aarts & J. H. M. Korst	E	Parallele netwerken voor het simuleren van leerdrag: Boltzmann-machines	Informatie 28	632-641	1986
R. Broer*, G. Aissing*, W. C. Nieuwpoort* (* Univ. Groningen) & L. F. Feiner	E	Hartree-Fock cluster study of interstitial transition metals in silicon	Int. J. Quantum Chem. 29	1059-1066	1986
A. H. van Ommen, H. J. W. van Houtum & A. M. L. Theunissen	E	Diffusion of ion implanted As in $TiSi_2$	J. Appl. Phys. 60	627-630	1986
H. A. van Sprang & P. A. Breddels	E	Numerical calculations of director patterns in highly twisted nematic configurations with nonzero pretilt angles	J. Appl. Phys. 60	968-972	1986
G. W. 't Hooft & C. van Opdorp	E	Determination of bulk minority-carrier lifetime and surface/interface recombination velocity from photoluminescence decay of a semi-infinite semiconductor slab	J. Appl. Phys. 60	1065-1070	1986
C. van Berkel & M. J. Powell	R	Photo-field effect in amorphous silicon thin-film transistors	J. Appl. Phys. 60	1521-1527	1986
C. W. J. Beenakker	E	Ewald sum of the Rotne-Prager tensor	J. Chem. Phys. 85	1581-1582	1986
H. K. Kuiken, J. J. Kelly & P. H. L. Notten	E	Etching profiles at resist edges: I. Mathematical models for diffusion-controlled cases	J. Electrochem. Soc. 133	1217-1226	1986
P. H. L. Notten, J. J. Kelly & H. K. Kuiken	E	Etching profiles at resist edges: II. Experimental confirmation of models using GaAs	J. Electrochem. Soc. 133	1226-1232	1986
W. A. P. Claassen*, H. A. J. T. van de Pol*, A. H. Goemans* (* Philips Elcoma Div., Nijmegen) & A. E. T. Kuiper	E	Characterization of silicon-oxynitride films deposited by plasma-enhanced CVD	J. Electrochem. Soc. 133	1458-1464	1986
J. W. D. Martens & W. L. Peeters	E	Anisotropy in cobalt-ferrite thin films	J. Magn. & Magn. Mater. 61	21-23	1986
R. van Mens	E	Ternary phase studies of Nd-Fe-X where X = C, Si, Ge, Pb, Sn	J. Magn. & Magn. Mater. 61	24-28	1986

J. E. Knowles	R	Magnetization reversal by flipping, in acicular particles of $\gamma\text{-Fe}_2\text{O}_3$	J. Magn. & Magn. Mater. 61	121-128	1986
D. M. Krol (<i>AT & T Bell Labs, Murray Hill, NJ</i>), C. A. M. Mulder & J. G. van Lierop	E	Raman investigation of autoclave-prepared, monolithic silica gels	J. Non-Cryst. Solids 86	241-250	1986
M. H. L. M. van den Broek	E	Calculation of the radius of the emitting area and the current in triode electron guns	J. Phys. D 19	1389-1399	1986
M. H. L. M. van den Broek	E	Experimental emittance diagrams of triode electron guns	J. Phys. D 19	1401-1419	1986
E. W. Meijer	E	A model study for coatings containing hexamethoxymethylmelamine	J. Polym. Sci. A: Polym. Chem. 24	2199-2208	1986
P. J. Courtois & P. Semal	B	Block iterative algorithms for stochastic matrices	Linear Algebra & Appl. 76	59-70	1986
R. B. Helmholtz (<i>ECN, Petten</i>), T. T. M. Palstra*, G. J. Nieuwenhuys*, J. A. Mydosh* (<i>* Univ. Leiden</i>), A. M. van der Kraan (<i>IRI, Delft</i>) & K. H. J. Buschow	E	Magnetic properties of $\text{La}(\text{Fe}_x\text{Al}_{1-x})_{13}$ determined via neutron scattering and Mössbauer spectroscopy	Phys. Rev. B 34	169-173	1986
U.ENZ	E	The confined breather: quantized states from classical field theory	Physica 21D	1-6	1986
J. G. Kloosterboer, G. F. C. M. Lijten & F. J. A. M. Greidanus	E	Structure and stability of polyacrylate radicals trapped in a network	Polym. Commun. 27	268-271	1986
R. Brehm, J. C. Driessen, P. v. Grootel & T. G. Gijsbers	E	Low thermal expansion materials for high precision measurement equipment	Precision Eng. 7	157-160	1985
W. J. J. Rey	E	Multivariate data analysis, contributions and shortcomings of robustness in practice	Proc. Compstat, Rome 1986	197-204	1986
F. Rausch*, H. Lindeman* (<i>* Philips Bipolar IC Devel., Nijmegen</i>), W. Josquin, D. de Lang & P. J. W. Jochems	E	An analog BIMOS technology	Proc. Int. Conf. on Solid state devices & materials, Tokyo 1986	65-68	1986
H. C. de Graaff, W. J. Kloosterman & T. N. Jansen	E	Compact bipolar transistor model for CACD, with accurate description of collector behaviour	Proc. Int. Conf. on Solid state devices & materials, Tokyo 1986	287-290	1986
J. W. Slotboom, K. Osinski, G. J. T. Davids, H. G. R. Maas, A. Slob, K. E. Kuijk & B. Jansen		High-density CCD video-memories: physics and technology	Proc. Int. Conf. on Solid state devices & materials, Tokyo 1986	315-318	1986
H. L. Peek & D. R. Wolters	E	Enhanced tunneling in implanted MOS structures	Proc. Int. Conf. on Solid state devices & materials, Tokyo 1986	487-490	1986
M. P. A. Vieggers, B. H. Koek & A. H. van Ommen	E	An ordered precipitate structure and the formation of coesite in oxygen implanted silicon	Proc. Int. Conf. on Solid state devices & materials, Tokyo 1986	577-580	1986
C. W. T. Bulle-Lieuwma & P. C. Zalm	E	Ion-milling procedures for cross-sectional TEM specimens	Proc. Int. Cong. on Electron Microscopy, Kyoto 1986	351-352	1986
M. P. A. Vieggers, B. H. Koek & A. H. van Ommen	E	Ordered precipitate structure in oxygen implanted silicon	Proc. Int. Cong. on Electron Microscopy, Kyoto 1986	1507-1508	1986
E. Arnold, U. Landman*, S. Ramesh, W. D. Luedtke*, R. N. Barnett*, C. L. Cleveland* (<i>* Georgia Inst. Techn., Atlanta, Ga</i>), A. Martinez, H. Baumgart & B. Kahn	N	Formation of facets at the solid-melt interface in silicon	Proc. MRS Conf., Boston, MA, 1986	21-27	1986
S. Ramesh, A. Martinez, J. Petruzello, H. Baumgart & E. Arnold	N	Addressing the problems of agglomeration, surface roughness and crystal imperfection in SOI film	Proc. MRS Conf., Boston, MA, 1986	45-51	1986
L. J. M. Esser, L. G. M. Heldens, L. J. van de Polder, H. C. G. M. van Kuijk, H. L. Peek, G. T. J. van Gaal-Vandormael & C. A. M. Jaspers	E	The PAN-imager	Proc. SPIE 591	61-66	1986

H. Baumgart	N	Silicon-on-insulator technology by crystallization on quartz substrates	Proc. SPIE 623	211-223	1986
A. W. M. van den Enden & G. A. L. Leenknecht	E	Design of optimal IIR filters with arbitrary amplitude and phase requirements	Signal processing III, I.T. Young <i>et al.</i> (eds), Elsevier Science, Amsterdam	183-186	1986
P. C. W. Sommen	E	On the convergence behaviour of a frequency-domain adaptive filter with an efficient window function	Signal processing III, I.T. Young <i>et al.</i> (eds), Elsevier Science, Amsterdam	211-214	1986
R. N. J. Veldhuis	E	A method for the restoration of burst errors in speech signals	Signal processing III, I.T. Young <i>et al.</i> (eds), Elsevier Science, Amsterdam	403-406	1986
F. L. H. M. Stumpers	E	Computers of the fifth generation and their role in communications	Software system design methods, J. K. Skwirzynski (ed.), Springer, Berlin	475-489	1986
B. Vitt	A	Characterization of a solar selective black cobalt coating	Sol. Energy Mater. 13	323-350	1986
F. M. Klaassen & W. Hes	E	On the temperature coefficient of the MOSFET threshold voltage	Solid-State Electron. 29	787-789	1986
P. R. Boudewijn, M. R. Leys (<i>Univ. Lund</i>) & F. Roozeboom	E	SIMS analysis of $Al_xGa_{1-x}As/GaAs$ layered structures grown by metal-organic vapour phase epitaxy	Surf. & Interface Anal. 9	303-308	1986
B. A. Joyce, P. J. Dobson, J. H. Neave & J. Zhang (<i>Imp. College, London</i>)	R	The determination of MBE growth mechanisms using dynamic RHEED techniques	Surf. Sci. 174	1-9	1986

Contents of Philips Telecommunication and Data Systems Review 44, No. 2, 1986

- P. H. A. Timmers & H. Zantema: MEMPHIS/4000: Software Production Automated (pp. 1-21)
 R. C. Barendregt: The P 5040 word processor as an Electronic mail terminal (pp. 22-29)
 R. Slagter & L. Kool: VPU, The multimode teletext encoder (pp. 30-37)

Contents of Electronic Components & Applications 8, No. 1, 1987

- D. Wong: High-speed 12-bit tracking ADC using programmable logic sequencers (pp. 2-14)
 W. Langenhorst & R. Waser: Ceramic-chip capacitors — high-rel products with a promising future (pp. 15-20)
 M. Roberts: Stereo sound generator for sound effects and music synthesis (pp. 21-31)
 G. Goodhue, J. Jenkins & A. Khan: Harvard architecture pushes microcontroller IC into high-speed realm (pp. 32-38)
 A. Lentzer & G. Stäcker: Integrated video programming system (VPS) decoder (pp. 39-44)
 A. Otten & J. Slakhorst: High-temperature electrolytic capacitors (pp. 45-51)
 R. Kohlman: Integrated Services Terminal (IST) Bus — the low-cost LAN and associated interface (pp. 52-61)

Contents of Valvo Berichte, November 1986

- K. Ruschmeyer: Höher ausnutzbare Spulen durch permanentmagnetische Vormagnetisierung (pp. 1-9)
 K. Ruschmeyer & H. Wieters: Stufendrosseln mit weichmagnetischen Einlagen im Luftspalt (pp. 10-15)
 H. Wieters: Induktivitätsmessung vormagnetisierter Drosseln (pp. 16-18)
 A. Petersen: Piezokeramische Aktuatoren (pp. 19-31)

An open 800-kV ion-implantation machine

H. J. Ligthart and J. Politiek

Conventional commercial ion-implantation machines are capable of implanting the usual donor and acceptor elements phosphorus, arsenic and boron in silicon. The accelerating voltage of these machines is generally no higher than 200 kV. If heavier elements are to be implanted in heavier substrates, higher accelerating voltages are required. In scientific research the ion source should also be easily interchangeable. In ion-implantation machines for accelerating voltages up to 1 MV, which are not normally commercially available, the high-voltage section is insulated from the environment by gas under pressure in a tank. Changing an ion source therefore takes a long time. The ion-implantation machine that has been developed at Philips Research Laboratories, however, has an 'open' high-voltage section, which is insulated from the environment by atmospheric air, so that the ion source can be changed very quickly.

Introduction

The object of ion implantation is to replace some of the atoms in a crystalline substance by other atoms that also fit into the lattice structure. The atoms of the element to be implanted are therefore ionized, accelerated to a high velocity in an electrostatic field and then implanted in the substrate. The resulting damage to the lattice can be removed by means of a heat treatment (annealing). Ion implantation is widely used in the manufacture of integrated circuits, for doping silicon with donors (e.g. phosphorus) or with acceptors (e.g. boron). An advantage of ion implantation as compared with thermal diffusion is that it is easier to control the concentration profile, i.e. the concentration (usually in cm^{-3}) as a function of depth. There is also much less lateral diffusion, and annealing times can be kept much shorter [1].

Fig. 1 shows a diagram of the concentration profile produced by a single implantation. (Other shapes of profile can be obtained by combining several implantations.) At a certain depth the concentration decreases sharply to zero, since the ions have completely lost their kinetic energy. Theoretically the dose (in cm^{-2}) is

equal to the integral of the ion-flux density ϕ (in $\text{cm}^{-2} \text{s}^{-1}$) over time at the substrate surface. The dose is also equal to the integral of the concentration over the depth after the implantation. In practice a narrow ion beam that describes a fine raster is used, instead of a wide beam that covers the entire surface of the sub-

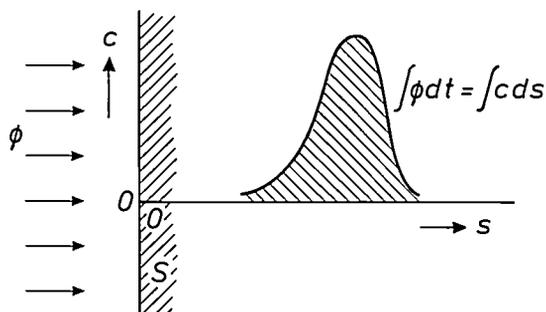


Fig. 1. Diagram of a concentration profile. *S* sample, typically a silicon slice, whose surface intersects the plane of the diagram at the vertical axis. *c* concentration of implanted ions. *s* depth from surface. ϕ ion-flux density. *t* time.

[1] W. K. Hofker and J. Politiek, Ion implantation in semiconductors, Philips Tech. Rev. 39, 1-14, 1980; S. T. Picraux and P. S. Peercy, Ion implantation of surfaces, Sci. Am. 252, No. 3, 84-92, 1985.

strate. The spacing of the lines of the raster is less than the half-width of the beam cross-section. The mean ion-flux density is therefore equal to the beam current divided by the area of the raster. It follows that the height of the concentration profile is proportional to the beam current and to the duration of the implantation process. The location of the concentration profile depends on the kinetic energy of the ions and on the atomic numbers of the implanted element and the substrate material. If one of these two atomic numbers is increased, the profile moves towards the surface of the substrate. If the energy of the ions is increased, the profile moves into the material.

In the investigation of new materials for integrated circuits, use is made of heavier substrate materials (for example gallium phosphide, gallium arsenide or indium phosphide, known as III-V compounds) and heavier donors and acceptors (e.g. tellurium and cadmium). In research on integrated circuits it is also necessary to make deeper implantations, for example of oxygen ions to produce 'buried' SiO₂ layers. This accounts for the recent trend of implanting ions at higher kinetic energy.

Fig. 2 gives an idea of the collision processes that take place when ions move through a substrate. The diagram shows the ion-energy loss per unit pathlength as a function of the square root of the ion energy; the diagram is based on the LSS model (Lindhard, Scharff and Schiøtt) [2]. The straight line corresponds to the ion-energy loss resulting from collisions with electrons. These inelastic collisions cause no damage in the substrate. The curve corresponds to the energy loss due to collisions with nuclei. These elastic collisions do cause damage in the substrate. At a kinetic

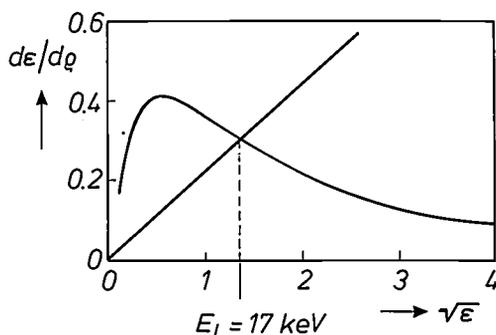


Fig. 2. The calculated energy loss per unit pathlength of the ions moving in the substrate as a result of collisions with nuclei of substrate atoms (curve) and with electrons (straight line), as a function of the square root of the energy of the ions as given by the LSS model [2]. To obtain a result valid for all elements, the 'reduced' quantities ρ and ϵ are used for the path and the energy, respectively. If the ions have a higher kinetic energy than the threshold value E_1 , the collisions with electrons predominate over those with nuclei. The slope of the straight line for the collisions with electrons is also a function of the atomic numbers of the ion and the substrate atoms. The diagram applies to the implantation of boron in silicon.

energy greater than E_1 , at the intersection of the straight line and the curve, the inelastic collisions, which cause no damage, therefore predominate. (In implantations in silicon E_1 is 17 keV for B⁺ ions, 140 keV for P⁺ ions, 800 keV for As⁺ ions and 2000 keV for Sb⁺ ions.) This means that with very deep implantations, with ions of high energy, relatively little lattice damage is caused at the substrate surface. Semiconductor structures already produced in the substrate should not therefore be damaged by implantations at greater depth.

How do we obtain ions of high kinetic energy? In the first place, of course, by increasing the strength of the electrostatic field, i.e. by increasing the potential at the ion source. (The substrate is normally at earth potential.) A higher energy can also be obtained by using multiply charged ions. A third method is to use a 'tandem' accelerator.

In a tandem accelerator the ion source is at earth potential. Single negatively charged ions from the source are accelerated and turned into positively charged ions (possibly multiply charged) in a 'stripper', consisting of a gas in a cell at positive potential. In this charge reversal one or more electrons are stripped from the ions by collisions with gas atoms. The positively charged ions then move towards the substrate, which is at earth potential. Disadvantages of tandem accelerators are that they have a low ion yield and are more complicated than single-ended accelerators.

The yield of an ion source generally decreases as the charge of the individual ions increases. Since the use of multiply charged ions instead of singly charged ions therefore results in a lower ion-flux density, the implantation time must be longer to produce a certain maximum concentration; see fig. 1. If deeper and more highly doped layers are required, implantation of multiply charged ions may not always be the right answer. Very deep implantations therefore require a high accelerating voltage, especially when heavy ions have to be implanted in heavy substrate material.

At Philips Research Laboratories in Eindhoven an ion-implantation machine has been developed, intended not only for research on new IC technologies but also for other work such as improving the hardness, wear resistance or corrosion resistance of metal surfaces. The high accelerating voltage of the machine, 800 kV, can be made even higher in the future. The machine was designed with an 'open' configuration, which means that the high-voltage section is not insulated by a gas (e.g. sulphur hexafluoride) under pressure in a tank, but is surrounded by atmospheric air at about 20 °C and at a relative humidity of less than 40%. This makes it possible to change ion sources quickly. Some of the investigations on the machine are connected with the selection of ion

sources, since the effectiveness of the ion source largely determines the time required for an implantation. Fig. 3a shows the high-voltage section of the implantation machine with the acceleration tube on the right (the lower of the two tubes projecting through the wall). Fig. 3b shows the low-voltage section; the

labeled from the environment by SF_6 gas at a slight excess pressure. The cascade circuit is supplied by an alternating voltage at 24 kHz and a maximum amplitude of 20 kV. The maximum voltage across each diode is therefore 40 kV. The generator can deliver a current of 2 mA via a copper-wire connection between



Fig. 3. Photograph of a) the high-voltage section and b) the low-voltage section of the 800-kV ion-implantation machine.

ions travel through this after they have passed through the acceleration tube.

In the following we shall first describe the construction of the implantation machine, and then we shall consider a few ion sources. The article concludes with a few examples of implantations to illustrate the usefulness of the machine.

Construction of the implantation machine

Fig. 4 is a diagram of the machine. The high-voltage section contains the high-voltage generator G , the supply unit SU and the ion-source unit IU ; see also fig. 3a. The low-voltage section contains the switching magnet SM , which directs the ion beam into one of the beam lines behind it. These beam lines terminate in the target chambers $TC1$ and $TC2$, see also fig. 3b. The high-voltage section is surrounded by concrete walls 30 cm thick. The walls provide a shield against X-radiation due to the collision of electrons moving in the opposite direction to the beam. The walls of the high-voltage section are clad with metal foil to screen the environment from interfering electromagnetic fields.

The high-voltage generator G was designed by High Voltage Engineering of Amersfoort for a voltage of 1 MV. It is a conventional cascade generator of the Cockcroft and Walton type [3]. The generator is insu-

lated from the environment by SF_6 gas at a slight excess pressure. The cascade circuit is supplied by an alternating voltage at 24 kHz and a maximum amplitude of 20 kV. The maximum voltage across each diode is therefore 40 kV. The generator can deliver a current of 2 mA via a copper-wire connection between

the corona sphere (the upper part of the generator) and the supply unit SU , giving a maximum power of 2 kW. The sum of the ion-beam current and the leakage currents cannot therefore exceed 2 mA. The supply unit SU contains two dynamos, each supplying a power of 12 kW. The dynamos are driven via insulating shafts by motors at earth potential in the basement below the high-voltage section. The dynamos provide the electrical energy for the various circuits, equipment and vacuum pumps in the units SU and IU . The unit SU also contains rectifier circuits that supply direct current to the magnet coils, and the electronic circuits for controlling the 'ion-optical' elements in IU . The control circuits communicate with the control console in the low-voltage section via optical-fibre connections. SU therefore also contains optoelectrical and electrooptical converters.

The unit IU contains the ion source, which is at a positive potential of 30 kV with respect to its immediate environment. The ions therefore arrive with a

[2] N. Bohr, The penetration of atomic particles through matter, K. Dan. Vidensk. Selsk. Mat.-Fys. Medd. 18, No. 8, 1914 (144 pages);
J. Lindhard, M. Scharff and H. E. Schiött, Range concepts and heavy ion ranges; notes on atomic collisions II, K. Dan. Vidensk. Selsk. Mat.-Fys. Medd. 33, No. 14, 1963 (39 pages).

[3] S. Gradstein, A modern high-voltage equipment, Philips Tech. Rev. 1, 6-10, 1936;
A. Kuntke, A generator for very high direct current voltage, Philips Tech. Rev. 2, 161-164, 1937.

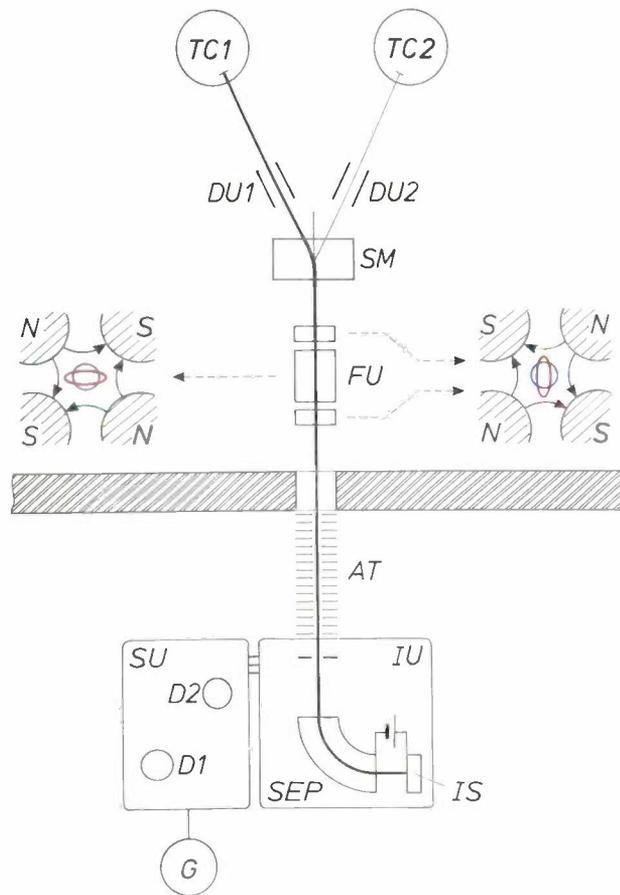


Fig. 4. Diagram of the ion-implantation machine. *G* high-voltage generator. *SU* supply unit. *D1*, *D2* dynamos. *IU* unit containing the ion source *IS*. *SEP* separator. The ion source is at a potential of 30 kV with respect to the separator. *AT* acceleration tube. *FU* focusing unit. This consists of three magnetic quadrupoles, whose action on the cross-section of the ion beam is illustrated in the separate diagrams on the left and right. *N*, *S* north and south poles of the quadrupoles. The beam cross-section in front of the quadrupole is shown in blue, the cross-section behind it is shown in red. *SM* switching magnet. *DU1* and *DU2* deflection units in the beam lines behind the switching magnet; see also fig. 6d. *TC1* and *TC2* target chambers.

certain initial velocity at the separator, which deflects them through an angle of 90° by means of a uniform magnetic field. The field-strength is adjusted so that the required ions pass through the exit slit while the unwanted ions do not. This follows from the expression for the radius R of the circular path of an ion in a magnetic field of flux density B ^[4]:

$$R = \frac{1}{B} \sqrt{\frac{2mE_{\text{kin}}}{(ze)^2}}, \quad (1)$$

where m is the mass of the ion, E_{kin} its kinetic energy, e the elementary charge and ze the charge of the ion. After the separator the ion beam enters the acceleration tube *AT*, which passes through the concrete wall between the high-voltage and low-voltage sections. The ions obtain the required kinetic energy in the acceleration tube. Turbomolecular pumps are connected to the ion source, the separator and the accelera-

tion tube; these pumps keep the pressure low enough for the mean free path of the ions to be sufficiently large.

Putting the separator at high voltage has the advantage that the magnetic field does not need to be so powerful. If the separator were located after the acceleration tube, then as eq. (1) shows the magnetic flux density would have to be five to six times higher for the same radius. The units *SU* and *IU* are mounted on insulating columns. Parts of these columns connect metal plates, each at a defined potential; see fig. 3a. This is arranged by electrically interconnecting the plates by identical resistors, which have a total resistance of about 3.85 G Ω . Resistors are used in the same way to interconnect discs that distribute the potential drop linearly along the acceleration tube. These resistors are contained in another tube above the actual acceleration tube (see fig. 3a) and have a total resistance of 10 G Ω . The leakage currents through the resistors total 0.3 mA at 800 kV, so that the ion-beam current cannot exceed about 1.7 mA. In practice, however, the ion-beam current is always smaller.

In the low-voltage section the ion beam first passes through the focusing unit *FU*. This consists of a 'magnetic triplet', i.e. three magnetic quadrupoles in a row. Each of the quadrupoles corrects the beam cross-section in the vertical and horizontal directions. This can be seen from the shape of the lines of force in the quadrupoles as shown in fig. 4. The strength of each of the outer quadrupoles is half that of the inner quadrupole. The result is that the group of three quadrupoles has a converging action. As the ions leave the separator they diverge, and the divergence increases in the acceleration tube because they repel one another, but the divergence is almost completely compensated by the magnetic-triplet focusing unit. This therefore acts as a positive lens, which produces very nearly a point image of the beam cross-section at the exit of the separator on the sample in the target chamber. Astigmatism in the beam can be compensated by increasing the strength of one of the quadrupoles.

The ion beam next passes through the switching magnet *SM*, which works on the same principle as the separator in the high-voltage section. The direction of the field in the switching magnet is reversible, so that charged particles in the beam can be deflected clockwise or anticlockwise as desired, through an angle of $\pm 20^\circ$. The switching magnet also removes neutral particles from the beam. A third and important function of the switching magnet is the elimination of unwanted ions produced by charge exchange.

The effectiveness of the double separation of particles, in the separator and in the switching magnet, can be seen in the implantation of phosphorus, for ex-

ample. At moderate temperatures, phosphorus vapour consists of P_4 molecules. However, at a temperature of 800°C , which is almost always exceeded in the ion source, the dissociation



commences. The source thus contains P^+ , P^{++} and P_2^+ ions. The P_2^+ ions may dissociate later in the reaction



If this happens at the inlet of the separator, the P^+ ions have a kinetic energy here of a quarter of that of the P^{++} ions. (We assume that the energy of a P_2^+ ion is evenly divided between the P^+ ion and the P atom.) It can be seen from eq. (1) that in this case the P^+ ions and the P^{++} ions describe a path of the same radius in the separator and are therefore not separated. It is easy to verify that these ions do in fact become separated in the switching magnet. (If an electrostatic deflection unit were to be used instead of a switching magnet, the P^+ and P^{++} ions would describe virtually the same parabolic path, so that these ions would not be separated.)

In each of the two beam lines following the switching magnet the ion beam can be moved horizontally and vertically by means of electrostatic deflection plates. At the entrance of each beam line there is a set of four short deflection plates, followed by a set of four long deflection plates, as shown in fig. 6*d*. The set of short plates is intended for deflecting ion beams of low energy, the set of long plates is for beams of high energy. The target chamber *TC1* is used for investigating various kinds of samples for other applications, unconnected with the manufacture of integrated circuits. Target chamber *TC2* is intended for implantations for IC manufacture. The slices are placed in a drum, which can accommodate 29 slices of diameter 100 mm. To comply with the clean-room conditions required in IC manufacture, target chamber *TC2* is placed in a special enclosure with a down-flow of filtered air. The vacuum is maintained by cryopumps to give a pressure of less than 10^{-7} mbar in the two beam lines and the target chambers [61].

Moving the beam over the sample

'Triangular' voltages applied to the horizontal and vertical deflection plates produce a uniform distribution of implanted ions over the specimen. The frequency of the triangular voltage for the horizontal deflection is 93 Hz; for the vertical deflection it is 800 Hz. Every 1/186 s the ion beam writes the pattern shown in blue in fig. 5. This pattern has

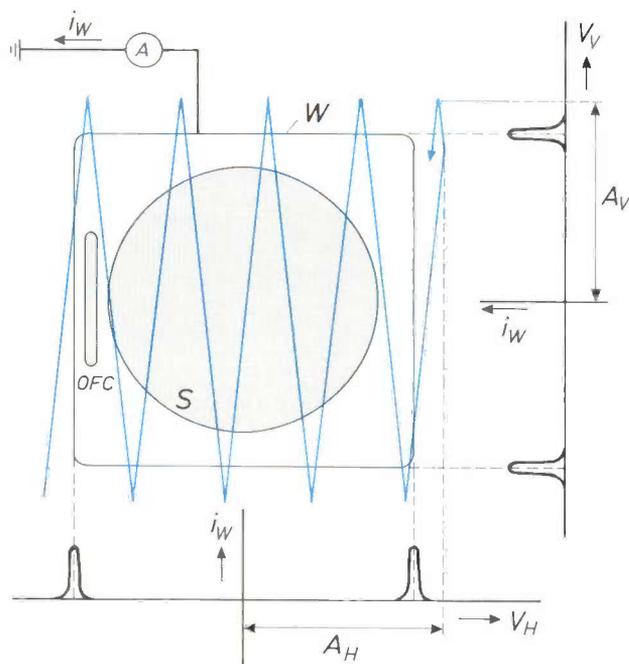


Fig. 5. The path (in blue) that the beam traces out on the sample during a half-period of the 'triangular' voltage for the horizontal deflection. *S* sample, typically a silicon slice. *OFC* opening of Faraday cup (*FC* in fig. 7). *W* wire frame, earthed via a milliammeter that measures the current i_W through the frame. The diagrams below and on the right correspond to the display on an oscilloscope screen when the vertical deflection in the oscilloscope is proportional to the current i_W and the horizontal deflection is proportional to the triangular voltage V_H or V_V on the horizontal or vertical deflection plates in the beam line. The diagram at the bottom of the figure is obtained with V_V set to 0; the diagram on the right is obtained with V_H set to 0. A_H and A_V are the amplitudes of V_H and V_V . The scales for V_H and V_V are chosen such that the distances between the peaks are equal to the width and height respectively of the wire frame.

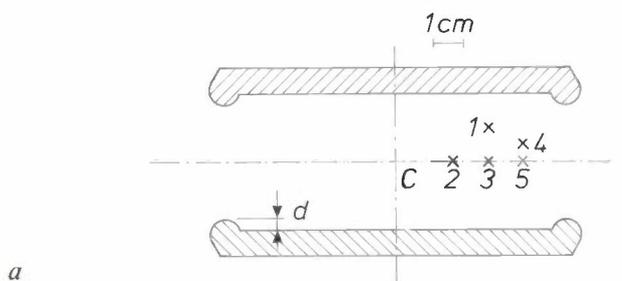
$800/2 \times 93 \approx 4.3$ periods. In each second the beam writes 186 such patterns, all adjacent. After exactly one second the writing of 186 non-overlapping patterns is repeated.

During implantation it is necessary to ensure that the beam direction does not coincide with an 'open' crystal direction (a 'channel') of the sample, otherwise the penetration depth of the ions will vary considerably if the angle of incidence is not absolutely constant. Nor should the direction of the beam lie in a 'channel plane'. The sample is therefore rotated in its own plane through a predetermined angle and the beam meets the surface of the sample at an angle of 7° to the normal.

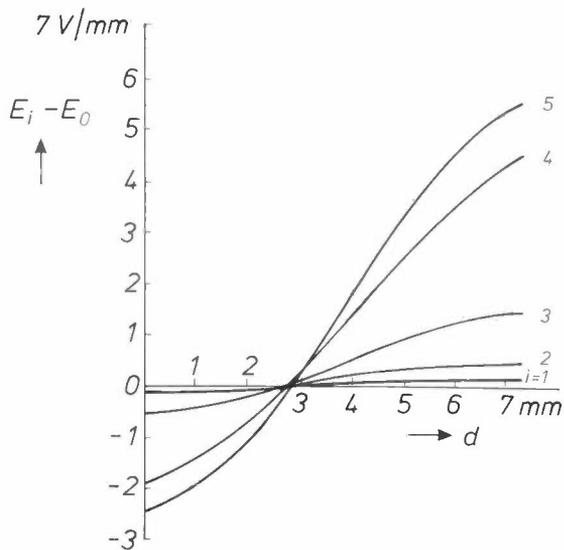
The usual requirement for the uniformity of an implantation is that the dose should not vary by more than 1% over the sample. This means that the trian-

[4] J. F. Ziegler (ed.), Ion implantation, science and technology, Academic Press, Orlando, FL, 1984; H. Rysse and H. Glawischning (ed.), Ion implantation techniques (Springer Ser. Electrophys., Vol. 10), Springer, Berlin 1982.

[61] J. J. Scheer and J. Visser, Application of cryopumps in industrial vacuum technology, Philips Tech. Rev. 39, 246-255, 1980.



a



b



c



d

gular voltage should be highly linear. In addition, the field between the deflection plates should be homogeneous over the largest possible region. If flat deflection plates were used, the strength of the electric field would decrease towards the edges. The deflection of an ion would then depend on the position where the ion entered the region between the plates. (This position depends on the deflection produced by the preceding plates and on the dimensions of the beam cross-section.) With flat deflection plates the dose can vary by a high percentage. We have solved this problem by adding a 'lip' at the edges of the plates. We determined the height of the lip with the ELOP and GELOP software packages [6]. Fig. 6 shows the shape of the deflection plates and the results of the computer calculations. For a lip 2.75 mm high the field is homogeneous over a very large region, and the dose produced by our implantation machine varies by no more than 0.5% over the wafer.

The correct amplitudes A_H and A_V for the triangular voltages on the deflection plates are set up with the aid of a wire frame around the sample; see fig. 5. The wire frame is earthed through a millimeter. While the amplitude of the voltage for the horizontal deflection is adjusted, the amplitude of the voltage for the vertical deflection is set to zero. If the vertical deflection on an oscilloscope is then made proportional to the current i_w through the wire frame and the horizontal deflection on the oscilloscope is made proportional to the voltage between the horizontal-deflection plates in the beam tube, the image that appears on the screen corresponds to the lower diagram in fig. 5. The location of the peaks on the oscilloscope screen shows how far the beam goes beyond the region inside the wire frame. The shape of the peak gives an idea of the half-width of the beam cross-section in the horizontal direction. A similar procedure is used for adjusting the voltage on the vertical deflection plates and estimating the half-width of the beam cross-section in the vertical direction.

The total number of ions incident on the sample — which is a measure of the average dose — is determined with the aid of a Faraday 'cup'. The opening in the Faraday cup is indicated by OFC in fig. 5. A diagram of the cross-section of the Faraday cup is shown

◁

Fig. 6. Calculation of the height d of the 'lip' on the deflection plates for the most uniform distribution of the electrostatic field. a) Cross-section of the deflection plates. At the points $i=1$ to 5 the field-strength E_i was determined by the ELOP and GELOP software packages [6]. E_0 field-strength at the centre C: 120 V/mm. b) The difference $E_i - E_0$ between the field-strengths at the points 1 to 5 and at the centre as a function of d . For $d=2.75$ mm the differences for these five points are virtually negligible. c) Photograph of a short deflection plate with lip. d) Photograph of the assembly of short and long deflection plates.

in *fig. 7*. The cup 'captures' all positively charged particles that pass through the opening. The sum of the charges of all these particles is measured by a current integrator. Electrode E_1 is at a low negative voltage, which ensures that secondary electrons generated on the outside wall of the cup do not enter the cup and affect the result of the measurement. Electrode E_2 captures the secondary electrons produced in the cup. Electrode E_3 captures positively charged particles of low energy, produced by ion collisions inside the cup. Electrode E_4 captures the ions from the beam that have passed through the opening of the cup. The currents through E_2 and E_3 also contribute to the measured result.

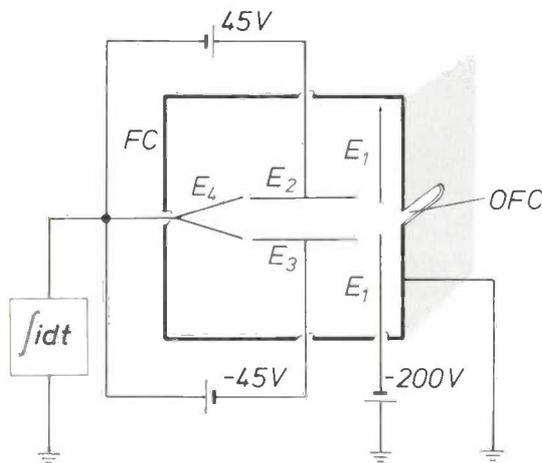


Fig. 7. Diagram of the cross-section of the Faraday cup *FC* with opening *OFC*; see *fig. 5*, and with the electrodes E_{1-4} . Electrode E_4 captures the ions from the beam. The result of the measurement with the current integrator connected to E_4 is a measure of the implantation dose. E_1 repels secondary electrons from outside the cup. The charged particles captured by E_2 and E_3 contribute to the measured result.

The ion sources

As we have seen, it is very important that the ion source in an ion-implantation machine should give a high current of the appropriate ions. We shall now discuss five ion sources that are used in our implantation machine. It will be shown that each source has its characteristic advantages and disadvantages. As noted earlier, the ion source in an implantation machine for research should be easy to change.

Ions are produced by collisions between electrons and atoms. To produce a singly charged ion the colliding electron should have an energy above a threshold value that is generally between 5 and 10 eV, depending on the element. (This threshold value expressed in volts is called the ionization potential.) To change a singly charged ion into a doubly charged ion the electron must have an energy higher than a threshold value of 15 to 30 eV. In general, doubly charged

ions can be changed into triply charged ions at an electron energy higher than a threshold value of 30 to 100 eV. In the plasma generated in ion sources the electrons are therefore made to travel the greatest possible distance at a kinetic energy sufficiently high to produce the required ions. In four of the ion sources described here the electrons acquire their kinetic energy from an electrostatic field between an anode and a cathode. The fifth source makes use of microwaves for the energy transfer. The electrons are made to travel a long path through the action of a magnetic field, whose radius is given by equation (1). The pressure in the source is a compromise. On the one hand a large number of ions is desirable, so there must be sufficient gas to supply these ions. On the other hand, it is undesirable for multiply charged ions to lose their charge through collisions with other particles. A pressure smaller than 10^{-2} mbar is normally used.

The ions produced are drawn from the source with the aid of an extraction electrode, located a short distance — about 5 mm — from the actual source and at a potential of -30 kV with respect to the source. The extraction electrode is at the same time the input diaphragm for the separator; see *fig. 4*. Since the plasma is a good conductor, the boundary of the plasma corresponds to an equipotential surface, so that the ions leave the surface of the plasma perpendicularly. The gas pressure and the general discharge conditions in the source determine the shape of the plasma at the extraction aperture. Apart from the Maxwell-Boltzmann distribution, this shape and that of the extraction electrode thus determine the divergence of the ion beam as it leaves the source.

The five ion sources will now be briefly reviewed. (The sources in their original version were a product of High Voltage Engineering.)

The hollow-cathode source

Fig. 8 shows the principle of the hollow-cathode source, also known by the name of its first manufacturer as the Danfysik source. A great advantage of the source is that both gases and solids can be ionized. These are evaporated in a ceramic oven. The maximum oven temperature is 1500°C , which is high enough for many elements to reach a sufficiently high vapour pressure. To produce ions of elements for which this does not apply, it is necessary to start from compounds. The difficulty with the use of compounds is that many unwanted ions are produced.

In the hollow-cathode source the electrons for the plasma are supplied by a filament which connects two

[6] K. J. van Oostrum, CAD in light optics and electron optics, Philips Tech. Rev. 42, 69-84, 1985.

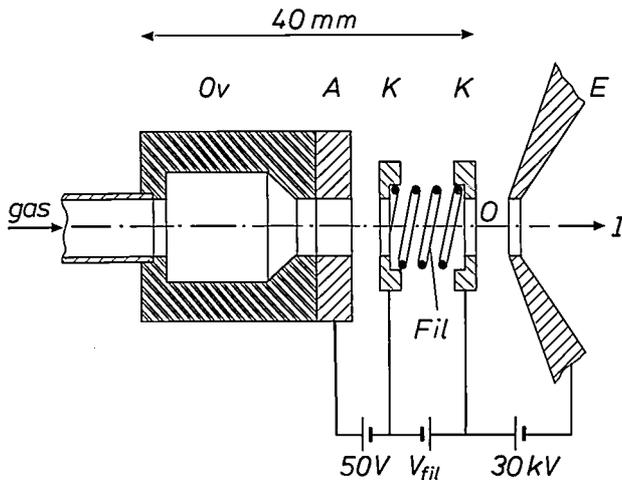


Fig. 8. The hollow-cathode source. *Ov* oven. *A* anode. *K* cathode, consisting of two parts. The two parts are connected by the filament *Fil*. V_{fil} filament voltage. *E* extraction electrode. *I* ion-beam current. The source contains a plasma that protrudes from the extraction aperture *O*.

parts of the cathode. A weak magnetic field is applied. A strong magnetic field is not necessary since the filament produces many electrons, giving rise to many ionizations. The anode-cathode voltage is low, about 50 V, to keep the anode-cathode current within bounds. Because of the low anode-cathode voltage the spread in the ion energy is small.

The versatility of the source makes it particularly suitable for research applications. There are a few disadvantages: the low anode-cathode voltage produces few doubly charged and multiply charged ions, the filament has a life of less than 24 hours, and the high temperatures of the filament and oven can cause unwanted reactions.

The Penning source

The Penning source, see fig. 9, does not have the disadvantages of the hollow-cathode source that result from the use of a filament. The operation of the Penning source, which is comparable with the well-known Penning gauge [7], is based on ionization with a high electric field-strength. The voltage between the anode and the cathodes is therefore high, 3 kV, and is more than enough to generate triply charged and multiply charged ions. A strong axial magnetic field makes the electrons move back and forth along helical paths between the two cathodes. The magnetic field is excited by a coil (not shown) mounted around the source.

The source is simple and reliable in use, requires little maintenance, and can produce a high current of multiply charged ions. It can generate, for example, a current of $0.1 \mu\text{A}$ of Ar^{4+} ions. (These ions therefore have a kinetic energy of 3 MeV at an accelerating voltage of 750 kV at the sample.) A disadvantage is that

only gaseous elements and compounds can be ionized. Also, since the ions in the source move in a strong electrostatic field, their energy varies a great deal at the extraction aperture, so that the energy spread in the beam is considerable. If this energy spread is limited by using a narrow exit slit for the separator, the ion current is reduced. If a broad exit slit is used, unwanted ions pass through and the angular spread of the ion beam also increases.

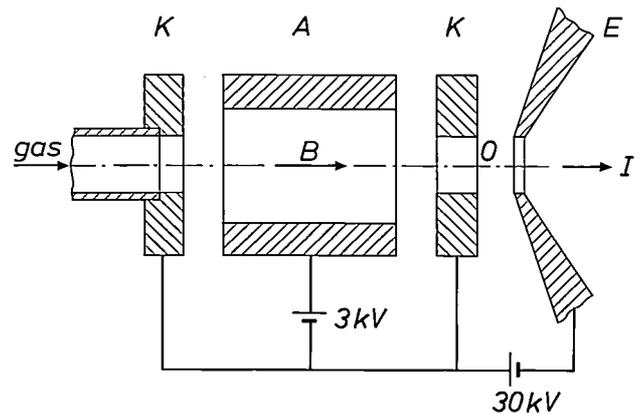


Fig. 9. The Penning source. The symbols have the same significance as in fig. 8. This source does not have a filament, but operates with a strong magnetic field of flux density *B*.

The radial Penning source

The radial-extraction Penning source, shown in fig. 10, does not have the disadvantage of a large spread in energy because the ions are extracted from one place in the plasma, where they all have about the same kinetic energy. Except for the ion extraction, this source is almost identical with the conventional Penning source. It has permanent magnets of samarium-cobalt for ease of construction.

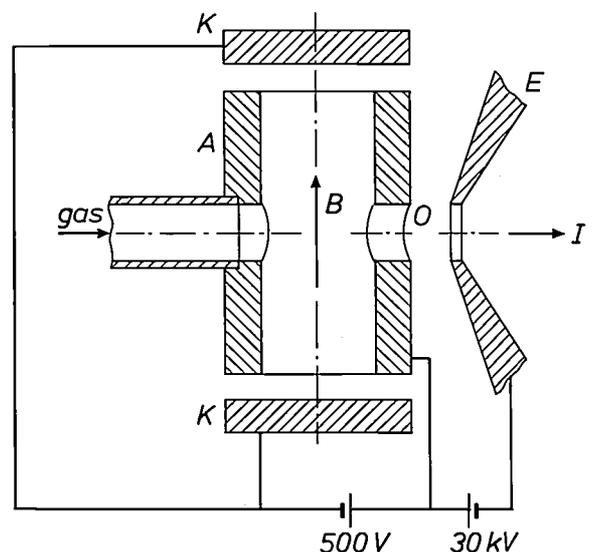


Fig. 10. The radial-extraction Penning source. The symbols have the same significance as in figs 8 and 9.

A disadvantage of this source is that it is necessary to use a high plasma density, which means that the anode-cathode current is high. The anode-cathode voltage associated with this current is about 500 V. The high current causes sputtering of cathode material, which means that the cathodes of the source have a shorter life than the cathodes of the conventional Penning source. This source is also only suitable for the ionization of gases, of course.

The sputter source

The sputter source, shown in *fig. 11*, has the special feature that it will ionize nearly all conducting solids whose melting point is not too low, including those with a very low vapour pressure (and usually a high melting point), e.g. tungsten and tantalum. This source also operates with a strong magnetic field; the anode-cathode voltage is about 180 V. The magnetic field is excited by a coil and can therefore be varied in magnitude.

The plasma is formed in argon gas. Ions of this gas collide with the cathode and sputter the cathode mater-

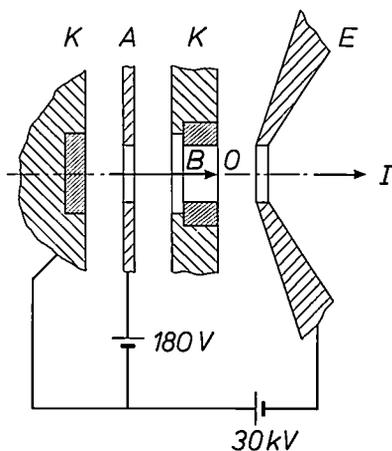


Fig. 11. The sputter source. The symbols have the same significance as in *figs 8 and 9*. The two parts of the cathode are coated with the element to be ionized.

ial. In the previous source this was a disadvantage, but here it is useful, because the two parts of the cathode are coated with the element to be ionized. The anode-cathode current must be high: 0.5 A. The ions formed at the part of the cathode near the extraction electrode contribute directly to the ion current. The ions formed at the other part of the cathode contribute much less since they cannot easily pass the anode. Indirectly, however, this part of the cathode does contribute to the ion current by generating neutral particles that are ionized later.

With the sputter source we have obtained currents of 30 μ A with Ta⁺ ions, 30 μ A with Ta²⁺ ions, 6 μ A

with Ta³⁺ ions and 2 μ A with Ta⁴⁺ ions. (A current of 1 μ A corresponds to 6.25×10^{12} singly charged particles per second.) With other ion sources it is particularly difficult to ionize tantalum.

The microwave source

In the microwave source, shown in *fig. 12*, energy is not transferred to the plasma by an electrostatic field but by a microwave electromagnetic field generated by a magnetron [8].

The frequency of the microwave circuit, which is connected to the source by a coaxial line is tuned to the cyclotron frequency of the electrons. An equation for this frequency can be found by substituting $E_{kin} = \frac{1}{2} m(2\pi fR)^2$ and $z = 1$ in equation (1):

$$f = \frac{Be}{2\pi m} \tag{2}$$

The cyclotron frequency is thus the frequency of rotation of the electrons in their helical path in a constant magnetic field; this frequency is independent of the kinetic energy of the electrons. If the energy increases, the radius of the path becomes larger, but the frequency of rotation remains the same.

In our case the frequency of the travelling waves in the microwave circuit is approximately equal to

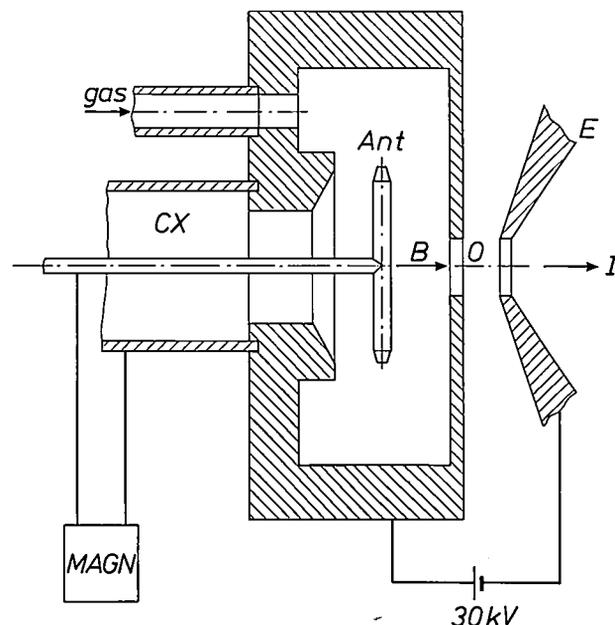


Fig. 12. The microwave source. This source is connected by a coaxial line *CX* to a magnetron *MAGN*, which supplies microwave energy to the plasma in the source. *Ant* antenna. See the captions to *figs 8 and 9*.

[7] F. M. Penning, High-vacuum gauges, Philips Tech. Rev. 2, 201-208, 1937.

[8] J. Ishikawa, Y. Takeiri and T. Takagi, Axial magnetic field extraction-type microwave ion source with a permanent magnet, Rev. Sci. Instrum. 55, 449-456, 1984.

2.45 GHz. It follows from equation (2) that the flux density B of the magnetic field in the source must be equal to about 0.0875 T (tesla). This field is excited by a permanent magnet (not shown in the drawing). This means that it is difficult to adjust the strength of the magnetic field to the microwave frequency. The energy is transferred from the microwave circuit to the plasma by a T-shaped antenna.

Because of the relatively low energy of the electrons, most of the ions produced are singly charged. An advantage is that the spread of energy in the ion beam is small, since the source does not contain an electrostatic field. The source is only suitable for the ionization of gases and gives a relatively large ion current, e.g. 400 μA of O_2^+ ions.

Examples of ion implantation

Implantation of oxygen in silicon

The implantation of a high dose of oxygen ions — up to $3 \times 10^{18} \text{ cm}^{-2}$ — at a fairly considerable depth in a single-crystal silicon slice makes it possible to produce 'buried' SiO_2 layers. If the dose is not too high, the surface layer is more or less undisturbed since E_i for O^+ ions in silicon is only about 30 keV; see fig. 2. It is therefore possible to produce integrated circuits at the surface on an insulating substrate. The transis-

tors of such a circuit on an insulating substrate interfere with one another less, are faster, and are also less sensitive to cosmic or other radiation.

We have used a microwave source to implant O_2^+ ions at an energy of 600 keV in silicon. During the implantation the sample was heated to a temperature of 500 to 600 °C. This high temperature ensured that the top layer did not become amorphous and that the surplus of oxygen ions could diffuse away upwards or downwards. Fig. 13a shows that as a result of this diffusion the relative oxygen concentration does not continue to increase after the stoichiometric value corresponding to SiO_2 , 66.7 at.%, has been reached. The total quantity of implanted oxygen ions thus determines the thickness of the SiO_2 layer produced. After the implantation an annealing treatment is necessary to remove local oxygen precipitates from the upper layer and to correct for lattice damage. Fig. 13b shows cross-sections of the sample before and after annealing at about 1200 °C for 2 hours.

Implantations in metals

As noted, the properties of a metal surface can be improved by implanting ions. Hardness, wear resistance or resistance to corrosion can be improved in this way. We have improved the corrosion resistance of copper by implanting aluminium ions. We used a hollow-cathode source that supplied Al^+ ions with a maximum dose of $5 \times 10^{17} \text{ cm}^{-2}$ at an energy of 170 keV.

The copper samples were subjected to a corrosion test to IEC 68243 KD. At a temperature of 300K treated and untreated specimens were exposed to air containing 15 ppm of H_2S , at a relative humidity of 75 %. After 8 hours the untreated copper was found to have about 50 times the corrosion of the treated samples.

Fig. 14 shows the results of Auger analyses of the treated copper surface after the corrosion test. Presumably the implanted aluminium was first oxidized to Al_2O_3 , and the copper at the surface to Cu_2O . The Cu_2O is converted to Cu_2S . The Al_2O_3 layer finally prevents further corrosion.

- [9] R. S. Muller and T. I. Kamins, Device electronics for integrated circuits, Wiley, New York 1984, p. 374; S. M. Sze (ed.), VLSI technology, McGraw-Hill, New York 1983, p. 483; W. G. Gelling and F. Valster, The new centre for submicron IC technology, Philips Tech. Rev. 42, 266-273, 1985/86.
- [10] R. D. Rung, C. J. Dell'oca and L. G. Walker, A retrograde p-well for higher density CMOS, IEEE Trans. ED-28, 1115-1119, 1981.
- [11] D. Pramanik and M. I. Current, MeV implantation for silicon device fabrication, Solid State Technol. 27, No. 5 (May), 211-216, 1984; C. McKenna, C. Russo, B. Pedersen and D. Downey, Applications for MeV ion implantation, Semicond. Int. 9, No. 4 (April), 101-107, 1986.
- [12] J. Hilibrand and R. D. Gold, Determination of the impurity distribution in junction diodes from capacitance-voltage measurements, RCA Rev. 21, 245-252, 1960.

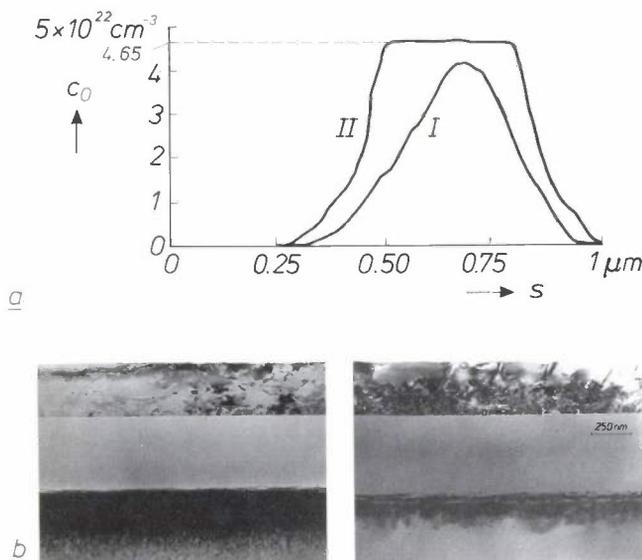


Fig. 13. a) The concentration c_O of oxygen in silicon as a function of the depth s below the substrate surface. The curves show that the implantation of oxygen in silicon produces a 'buried' SiO_2 layer as soon as the relative oxygen concentration reaches the stoichiometric value of 66.7 at.%, which corresponds to $c_O = 4.65 \times 10^{22} \text{ cm}^{-3}$. Curve I corresponds to a dose of $1.5 \times 10^{18} \text{ cm}^{-2}$ and curve II corresponds to $2.5 \times 10^{18} \text{ cm}^{-2}$. The curves are the result of measurements by the HEIS method (High-Energy Ion Scattering)^[1]. b) Micrographs made with a Philips EM400T transmission electron microscope of cross-sections of a silicon slice with a buried SiO_2 layer: on the left with no annealing, on the right after annealing for two hours at 1200 °C. The second layer is always the amorphous SiO_2 layer. The surface is at the top and is a (100) plane of the original silicon crystal.

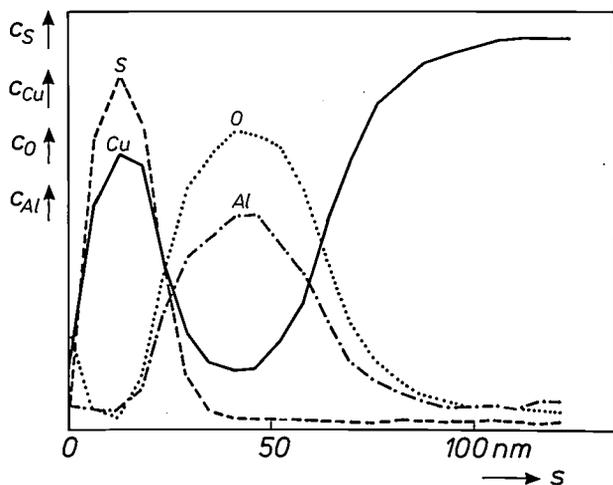


Fig. 14. Results of Auger analyses of a copper sample in which aluminium has been implanted to improve the corrosion resistance. After the implantation the sample was subjected to a corrosion test with H_2S gas and humid air. The scales for the concentrations c_S , c_{Cu} , c_O and c_{Al} of the four different elements in the sample are not comparable. The only known value is $c_{Cu} = 100\%$ for a depth s greater than 100 nm. Before the corrosion test the maximum relative value of c_{Al} was approximately 25 at. %.

High-energy implantations of phosphorus in silicon

In CMOS technology the normal practice for making 'wells' in the substrate [9] is to use low-energy ion implantation followed by drive-in annealing. The annealing enables the atoms of the implanted element to diffuse more deeply into the substrate. Disadvantages of this procedure are that lateral diffusion occurs, the heat treatment takes a good deal of time and it is only possible to obtain profiles whose concentration decreases with depth. What are wanted in the substrate, however, are wells in which the donor or acceptor concentration increases with depth, reaches a maximum and then decreases rapidly. Wells with a concentration profile of this type are called 'retrograde wells' [10]. They can be made by means of high-energy implantations. In CMOS circuits formed from retrograde wells the resistivity at the bottom of the wells is lower, so that there are fewer rejects due to 'latch-up' — a kind of short-circuiting caused by the formation of parasitic thyristors via the substrate.

Almost any shape of concentration profile can be obtained by combining implantations of different dose and energy [11], so that the concentration profiles of the separate implantations are added together. In this way the particular profiles can be made that will

ensure the best possible operation of the components of the integrated circuits. To obtain more information about such multiple implantations we have carried out a number of implantations at different energies and measured the corresponding concentration profiles. These were implantations of phosphorus in silicon at energies of 200 to 1400 keV, in steps of 200 keV. Up to 800 keV we used P^+ ions and above it we used P^{++} ions. Fig. 15 shows the results. The profiles were determined by means of the capacitance-voltage (CV) method [12].

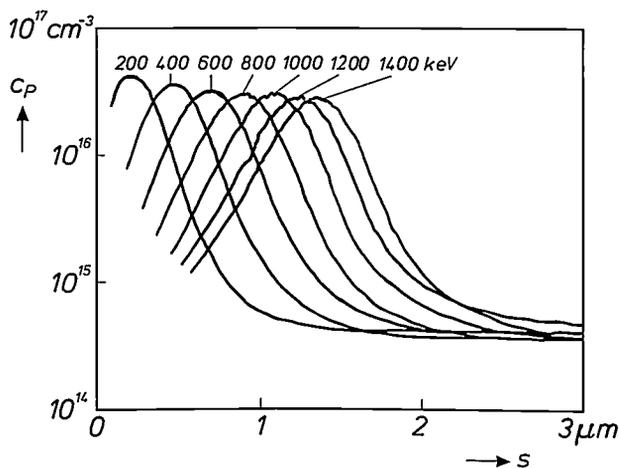


Fig. 15. Result of concentration-profile measurements made by the CV method [12]. The concentration c_p of phosphorus implanted in silicon is plotted as a function of the depth s . The different curves correspond to implantation energies from 200 to 1400 keV. Virtually any implantation profile can be obtained by combining a number of these or other curves.

Summary. Since the high-voltage section of the ion-implantation machine developed at Philips Research Laboratories is not insulated by gas under pressure in a tank, changing the ion source is relatively easy. The machine is therefore particularly suitable for research. The accelerating voltage is high, 800 kV, so that heavy elements can be implanted in heavy substrates. There are two target chambers, one designed for treating silicon slices for IC manufacture, the other designed for other samples, including metals. Implantations in metals can improve properties such as hardness or resistance to corrosion. The half-width of the beam cross-section and the location of the pattern of the beam scan on the sample are evaluated with the aid of a wire frame around the sample; the wire frame is earthed through a milliammeter. The implantation dose is measured with a Faraday cup. Five ion sources have been tested: the hollow-cathode source, the Penning source, the radial-extraction Penning source, the sputter source and the microwave source. Some examples are given to demonstrate the usefulness of the machine: implantation of buried oxygen layers in a silicon slice, implantation of aluminium to improve the corrosion resistance of copper, and the determination of a number of concentration profiles for phosphorus in silicon at different energies.

A laser-Doppler displacement meter

R. J. Asjes, C. S. Caspers and C. H. F. Velzel

In the manufacture of products such as wire, textiles or paper the moving length of material can be measured by pressing a runner wheel with a revolution counter against the surface of the moving material. This widely used method has its problems, however. The runner wheel can easily slip, resulting in an inaccurate measurement, and the material may be damaged by the wheel, which has to be in contact with it.

At the Philips Centre for Manufacturing Technology a displacement meter based on the laser-Doppler principle has now been developed. The instrument can be used for measuring displacements of widely diverse materials without being in contact with the material and with an error of less than 0.2%. Before looking more closely at the instrument itself, we shall first touch briefly on the laser-Doppler principle.

When a laser beam is split into two beams that are made to intersect, an interference pattern is produced in the intersection volume; see *fig. 1*. This interference pattern consists of light and dark planes, parallel to the 'optical axis' and perpendicular to the plane of the drawing. A small particle moving through the intersection volume passes through alternate light and dark planes. When it passes through a light plane it will scatter light, but not when it passes through a dark plane. If δ is the distance between the successive light planes, and v_p is the velocity of the particle perpendicular to the planes, intensity variations in the scattered light can be observed, with a frequency

$$f_D = \frac{v_p}{\delta}. \quad (1)$$

The principle of a laser-Doppler velocity measurement is that the frequency f_D is measured and the velocity v_p is calculated: the distance δ can be calculated from the geometry of the configuration and the wavelength of the laser light.

In the situation shown in *fig. 1* the detector has to have a very wide frequency range, since a high velocity implies a high frequency and a low velocity a low frequency. At very low velocities this can present

problems. The difficulty can be overcome, however, by arranging for the planes of the interference pattern to move at a constant velocity v_g . If we assume that the direction of v_g is opposite to that of v_p , then the frequency of the intensity variations is

$$f_m = \frac{v_p + v_g}{\delta} = f_D + f_g, \quad (2)$$

where f_m is the modulation frequency observed, v_p is the velocity of the particle, f_D is the Doppler frequency corresponding to the velocity v_p given by Eq. (1), and f_g is the frequency observed when $v_p = 0$. In this way, intensity variations (at a frequency f_g) are also observed when the particle itself is stationary, as the interference planes 'pass over it'. The interference pattern can be made to move by using a circular grating as the beam splitter. (The lines of such a grating are arranged radially around the circumference of a disc.) If the grating is rotated at a constant angular velocity, the planes in the interference pattern will also move at a constant velocity^[1]. When the interference planes move in the opposite direction to the particle, the frequency observed increases in proportion to the increase in the velocity of the particle. In reality, of

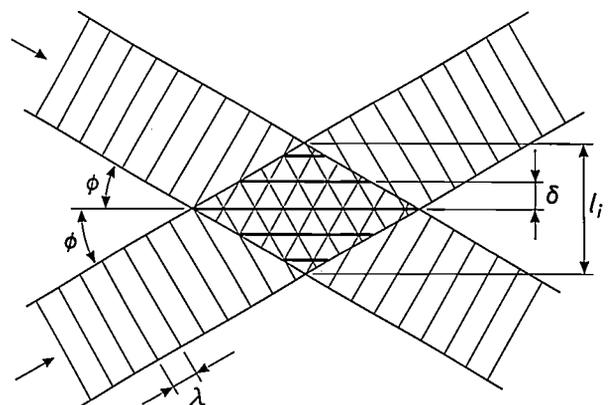


Fig. 1. The interference pattern obtained by splitting a laser beam into two beams and then making them intersect. The pattern consists of light and dark planes perpendicular to the plane through the two beams. If λ is the wavelength of the laser light and ϕ the angle between either beam and the optical axis, then the distance δ between the two maxima in the interference pattern is $\delta = \lambda / (2 \sin \phi)$. The length of the interference pattern in the direction of the velocity and length measurement is l_i .

course, there will not be just a single particle, but the surface of something like a cable or a wire, for example. Such a surface will almost invariably contain enough small asperities to represent a stream of particles whose mean velocity is measured.

As appears from Eq. (2) either f_m or f_g can be kept constant. Usually it is preferable to keep f_g constant^[1]. This can be done by rotating the circular grating at a constant angular velocity. When the velocity v_p changes, so too does the observed frequency f_m . We have preferred, however, to keep f_m constant (at least approximately)^[2]. This is done by controlling the angular velocity of the grating — and hence f_g — so that f_m always has about the same value. In this way, velocity and displacement can be determined in principle from the angular velocity of the grating, which is continuously measured by tachogenerator (we shall return to this presently). Some calculation is still necessary, however, since we do not keep f_m exactly constant. If we did, it would be only necessary to measure f_g to determine v_p . However, keeping f_m accurately constant would require an extremely refined control system. Keeping f_m only approximately constant means that the bandwidth of the processing electronics can be much smaller, and the velocity can be calculated very accurately from the two measured values of f_m and f_g .

We shall now look briefly at the mechanical part of the instrument (the 'measuring head') before dealing with the electronics in more detail.

The light source of the instrument (see *fig. 2*) is an HeNe laser with a radiated power of about 2 mW at a wavelength of 632.8 nm. The laser beam is focused on the rotating grating, and an image of the grating is produced on the moving surface of the object by the objective. Between the grating and the objective there is a diaphragm with two apertures, which passes only the two first-order beams and stops those of higher order. The light scattered by the object is focused on the detector, a silicon P-I-N diode, by a condenser, which consists of a combination of two plastic Fresnel lenses. These have a large diameter (so that they intercept much of the scattered light), and they are flat, light and inexpensive. Each lens has two holes for the two beams to pass through; the detector is mounted between the two beams. The optical system described here contains a minimum of optical components. The focusing lens and the objective must have good imaging characteristics and can be relatively small (about 2 cm). The condenser, on the other hand, does not need to have such good imaging characteristics but it does have to be relatively large (about 10 cm), as we noted earlier. It is therefore convenient to use different lenses for objective and condenser. (The same

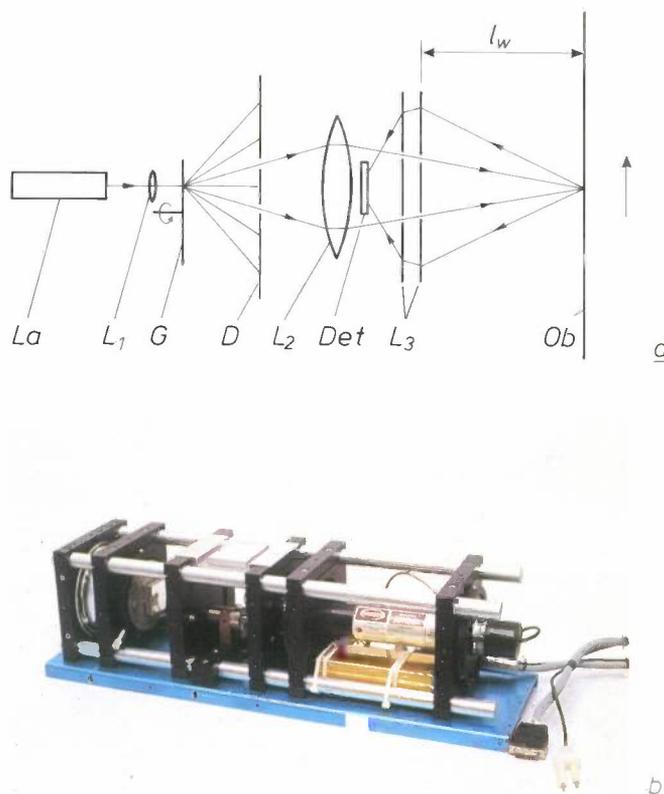


Fig. 2. a) Diagram of the 'measuring head' of the displacement meter. *La* HeNe laser. *L₁* focusing lens. *G* rotating grating. *D* diaphragm that passes only first-order beams. *L₂* objective lens. *Det* detector. *L₃* condenser, consisting of two Fresnel lenses. *Ob* object whose displacement is to be measured. *l_w* distance between the front of the measuring head and the object (working distance). b) Photograph of the measuring head. All the components are contained in square mounts connected to each other by four steel rods. The cylindrical laser is on the right between the rods; the Fresnel lens can just be seen on the left.

lens is often used for both functions, which makes the optics fairly expensive since a large lens with good imaging characteristics is required.) All the components are contained in a 'measuring head', made from components of a commercially available assembly system.

The electronic processing of the signal will be explained with the aid of the simplified diagram in *fig. 3*. The detector signal is fed to a preamplifier and a filter that passes frequencies between 420 kHz and 485 kHz. The signal passed by the filter is converted by a pulse shaper into a binary signal that drives a phase-locked loop consisting of a mixer (in our case an exclusive OR gate), a switch, an integrator and a voltage-controlled oscillator. The amplitude of the input signal of the pulse shaper is used for keeping the switch in the

[1] J. Oldengarm, Development of rotating diffraction gratings and their use in laser anemometry, *Opt. & Laser Technol.* 9, 69-71, 1977.

[2] The idea of keeping the modulation frequency constant is due to M. P. Weistra of Philips Research Laboratories, Eindhoven; see Netherlands patent application No. 8301917, 31 May 1983.

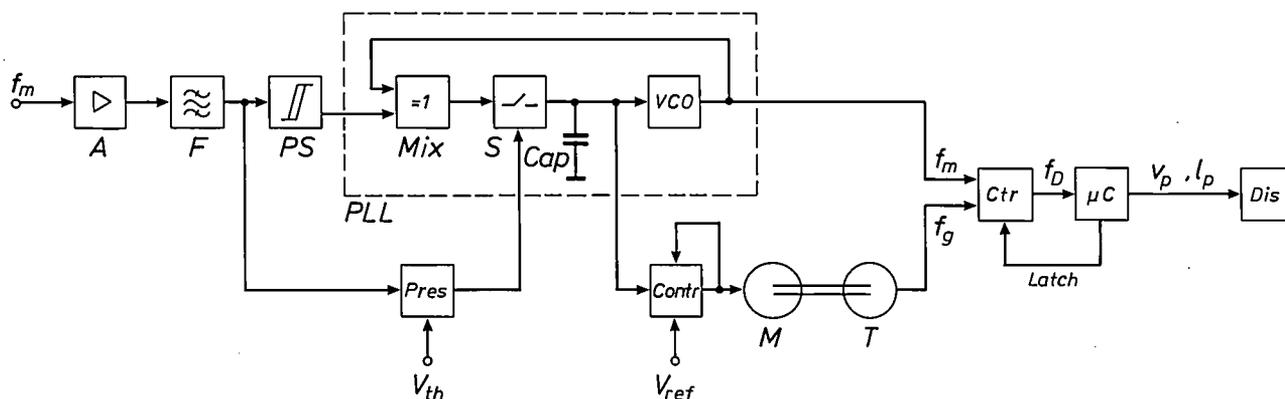


Fig. 3. Simplified block diagram of the signal-processing electronics, described in the text. f_m frequency of the input signal. A preamplifier. F bandpass filter. PS pulse shaper. Mix mixer (here an exclusive OR gate). S switch. Cap capacitor (integrator). VCO voltage-controlled oscillator. Mix , S , Cap and VCO form the phase-locked loop PLL . $Pres$ circuit for checking whether the amplitude of the input signal is smaller than the threshold voltage V_{th} . If it is, S is opened. $Contr$ control circuit for the motor that drives the grating. V_{ref} reference voltage. M motor for the grating. T tachometer. f_g frequency of the interference-pattern movement. Ctr difference counter. $latch$ counter-latching device. f_D Doppler frequency corresponding to the velocity v_p (given by Eq. 1). μC microcomputer. Dis display panels that display the measured velocity v_p and the displacement l_p .

phase-locked loop closed. When the amplitude of the detector signal falls below a threshold value V_{th} , e.g. if the signal is interrupted, the switch opens and the voltage-controlled oscillator continues to oscillate for about a second. In this way the system is protected from short interruptions. The input voltage of the oscillator is compared with a reference voltage V_{ref} ; the difference voltage is used for driving the grating motor. The effect of this is to limit the frequency of the detector signal to a band of a few kHz about a mean value of 455 kHz, which makes the analog part of the circuit considerably simpler^[2]. The output signal of the phase-locked loop consists of pulses at the frequency f_m . The frequency f_g at which the interference pattern moves is obtained from a tachometer mounted on the shaft of the motor that drives the grating. The pulse trains from the phase-locked loop and the tachometer are fed to a latched difference counter. A microcomputer determines the sampling time t_s (in our case 1/800 s) in which the frequencies f_m and f_g are subtracted one from another. After a time t_s the contents of the counter amount to $(f_m - f_g)t_s = f_D \times t_s$. As this value is latched after each sampling interval, it is available for the microcomputer during the next sampling interval. The computer applies a calibrating factor, to give $f_D \times t_s$ directly as the mean velocity over the sampling interval that has just elapsed. It also counts the successive values of

Table I. Some numerical data for the displacement meter. The symbols are explained in figs 1 and 2.

Maximum measurable displacement	16 km
Maximum measurable velocity	12 m/s
δ	0.02 mm
l_i	1 mm
l_w (minimum)	60 mm
l_w (maximum)	75 mm

$f_D \times t_s$ during a period of time, to give the displacement that has occurred during that time.

It can be shown by error analysis that the inaccuracy of the result of the measurement is about 0.03% when the velocity of the object is constant^[3]. If it is not constant, an absence of signal (e.g. because of a brief interruption of one of the two light beams) can introduce an error in the measured displacement. However, if the signal is present for 90% of the time, the standard error is still less than 0.2% — much less than with conventional methods of displacement measurement. The principal data of the instrument are listed in Table I.

^[3] C. H. F. Velzel, Laser Doppler displacement meter with controlled grating speed, Proc. Int. Conf. on Optical techniques in process control, The Hague 1983, pp. 289-294.

1937

THEN AND NOW

1987

Measuring instruments

By 1937 there was already a need for an easily operated test set that could be used for various electrical measurements on devices such as 'wireless valves' (photo lower right) [*]. The operation of the test set was substantially automated by using a contact bridge with 140 contacts. Inserting the 'code card' for a particular valve (photo lower centre) set the instrument up correctly in a single operation. Certain effects such as short-circuits between electrodes or a broken filament were indicated by signal lamps; other characteristics could be read from a milliammeter with its scale divided into two regions: blue for 'pass', red for 'fail'.

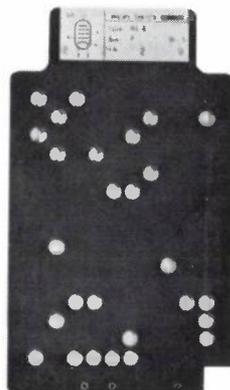
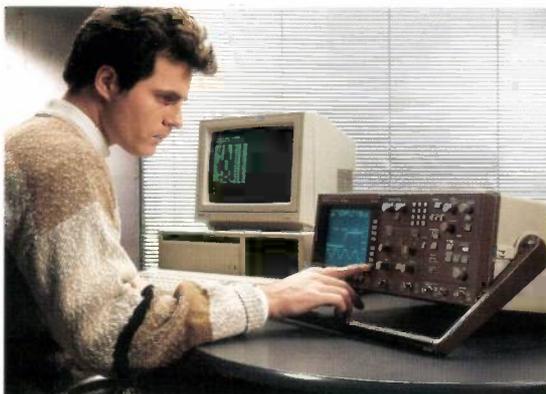
In 1987 instruments generally contain the most modern electronic components such as memories and microprocessors. Our example is the highly advanced PM 3320 digital storage oscilloscope (large photo). This instrument can be used for measuring analog signals with a bandwidth of no less than 200 MHz. It does so with the aid of P²CCD memories [**] and analog-to-digital conversion with an accuracy of 10 bits. The oscilloscope also includes a memory that will store more than 4 × 4000 measured values. Mathematical manipulations, such as the multiplication of two signals, can be performed before display on the screen (10 cm × 12 cm). If the two signals represent current



and voltage, for example, a curve representing the instantaneous power is obtained. The number of controls is kept as small as possible by introducing eight variable-function 'soft keys' next to the screen (photo lower left). The actual functions in use are shown on the screen. Without this facility the PM 3320 would need more than 100 controls. The complete setting-up routines for certain frequently made measurements can also be stored in the memory of the instrument, so that it can be set up instantaneously at the touch of a key.

[*] From Philips Technical Review, February 1937.

[**] See also H. Dollekamp, L. J. M. Esser and H. de Jong, P²CCD in 60 MHz oscilloscope with digital image storage, Philips Tech. Rev. 40, 56-68, 1982.



Quantitative measurements by the Schlieren method

G. Prast

In many different areas of technology the tolerances for the dimensions of parts and components are very strict. This means that inaccuracy of form and surface roughness must be as low as possible. Checking workpieces for these requirements can be difficult and sometimes it can only be done in the laboratory. In production conditions the Schlieren method may provide the answer. The combination of digital computation with this originally qualitative method from the last century allows us to use it for quantitative measurements.

Roughness and shape of surfaces

With the growing importance of optics — in the storage and exchange of digital information, for example — optical components such as mirrors and lenses will make an even greater contribution to Philips products than in the past. The specifications for the shape and surface roughness of lenses and mirrors are extremely demanding. The shape of the aspheric objective lens in the tracking mechanism of a LaserVision or Compact Disc system — a lens that can take the place of a system of four spherical lenses for reading out the information on a Compact Disc — must not differ from the standard by more than $0.1\ \mu\text{m}$. Obviously, it only makes sense to manufacture these extremely precise optical components if we can find out whether the product meets the specifications. The moulds for pressing the lenses can be made with the COLATH high-precision lathe^[1]. The normal practice so far has been to inspect the optical components afterwards. A number of optical components are then assembled to form a unit that must satisfy some easily checked criteria. This 'post'-inspection can mean that the moulds have to be machined once again. It is better to have a method that can be used to assess the quality of the mould when it is being made. If such a method can also be performed on the lathe itself, any

errors can be corrected immediately, without the need for repeated setting up, which can introduce new errors. Making optical components is only one use of this high-precision machining technique, of course. The specifications for the shape and roughness of machined surfaces for other purposes may be just as strict.

Two types of measurements are widely used for determining the shape and roughness of surfaces: mechanical methods and optical methods. In the mechanical methods of measurement a thin stylus is moved at constant velocity over the surface. During the movement the stylus is pressed lightly against the surface so that it follows its contours. The movement of the stylus perpendicular to the surface is measured and gives the roughness and shape of the surface. The mechanical methods have their disadvantages: they are time-consuming, sensitive to vibrations and cannot be used for all materials because of the direct contact between stylus and surface. Nor are they very practical for inspection during production, because the workpiece usually has to be moved elsewhere for the measurement.

The optical methods of measurement do not have the disadvantage of direct contact between workpiece and measuring device. This category includes light-scattering methods, interference methods and the Schlieren method. The Schlieren method is an example of an originally qualitative method that has since

been made quantitative by using modern equipment and digital signal processing, so that it can easily be used in production conditions. In this article we shall look first at the principle of the method and then discuss its application in a special type of microscope.

Principle of the Schlieren method

The method was introduced in the middle of the last century by J. B. L. Foucault as a way of checking the quality of large astronomical mirrors. A few years later A. Töpler discovered a new application: studying inhomogeneities in the flow of gases, e.g. in wind tunnels. It was this application that gave the method its name. A 'Schliere' or 'streak' is a thin zone of locally different refractive index in a transparent object. Light incident on the Schliere emerges in a different direction from the direction it would take in the absence of any such inhomogeneity. The change in direction is a measure of the inhomogeneity (or of the irregularity in the surface of an object if the light is reflected). A purely qualitative observation of a number of inhomogeneities in the refractive index of air can be seen in *fig. 1*.



Fig. 1. This picture, made by the Schlieren method, shows a bullet travelling through the air above two candle flames. This photograph is taken from the book of note [2] (page 560).

Differences in the refractive index of parts of a transparent object are not always easy to observe. It is possible, however, to enhance the contrast between light beams that pass through parts of the object with different refractive indices. This is done by screening off a part of each light beam emerging from the object [2]. *Fig. 2* shows the principle of the method. Lenses L_1 and L_2 produce an image of the light source B , and lenses L_2 and L_3 produce an image of the object O on the screen S . Part of the beam is screened off by the knife-edge K at the point where the image of the source is formed. The 'Schliere' is therefore shown on the screen as a streak that is more brightly illuminated than the background, since a smaller part of the beam from the 'Schliere' is intercepted by the knife-edge, as compared with neighbouring locations. In this example the object is transparent, but the same measurement method can of course be used with a reflecting object. Light source and screen are then situated on the same side of the object.

By measuring the light intensity (illuminance) on the screen the location of any inhomogeneity in the object or of any irregularity on its surface can be determined. Originally this method was mainly qualitative, but now that accurate photodiodes are available (in our configuration we use a row of 1024 CCD photodiodes with dimensions of $13\ \mu\text{m} \times 13\ \mu\text{m}$ on a single chip), and large amounts of measured data can be processed rapidly by computer, quantitative measurements are possible.

In the remainder of this section we shall derive these quantitative results. We shall initially assume that we have a configuration with ideal optical components. Later we shall see that errors that arise because these components are not ideal can be eliminated by using differential measurements.

Configuration with ideal optical components

In the arrangement shown in *fig. 2*, L_1 is an optically perfect lens and B is a light source that illuminates the object plane uniformly. A planar wavefront is incident on the object. If we assume that the gradient in the optical pathlength s_{opt} through the object at position P is $\partial s_{\text{opt}}/\partial x$, then the wavefront, which was truly planar before it met the object, will have a corresponding local deviation. There can be two reasons for a gradient in the optical path through the object. The refractive index in the object may differ from point to

[1] J. Haisma, W. Mesman, J. M. Oomen and J. C. Wijn, *Aspherics, III. Fabrication, testing and application of highly accurate aspheric optical elements*, Philips Tech. Rev. 41, 296-303, 1983/84.

[2] H. Wolter, *Schlieren-, Phasenkontrast- und Lichtschnittverfahren*, in: S. Flügge, *Handbuch der Physik*, Bd. 24, Springer, Berlin 1956, pp. 555-645.

point, or the object may be thicker or thinner in some places. If it is thicker or thinner the relation between the gradient in the object and the deviation in the shape of the wavefront can be explained in the following way. Consider an object of the shape shown in

of the wavefront) is equal to $d/(c/n)$ and also to $\Delta z/c + (d - \Delta z)/(c/n) + \Delta d/c$. The distance Δd by which the wavefront at the position Δx is 'ahead of' the wavefront at $x=0$ is therefore $\Delta d = \Delta z(n - 1)$. The tangent of the angle α between the wavefront and

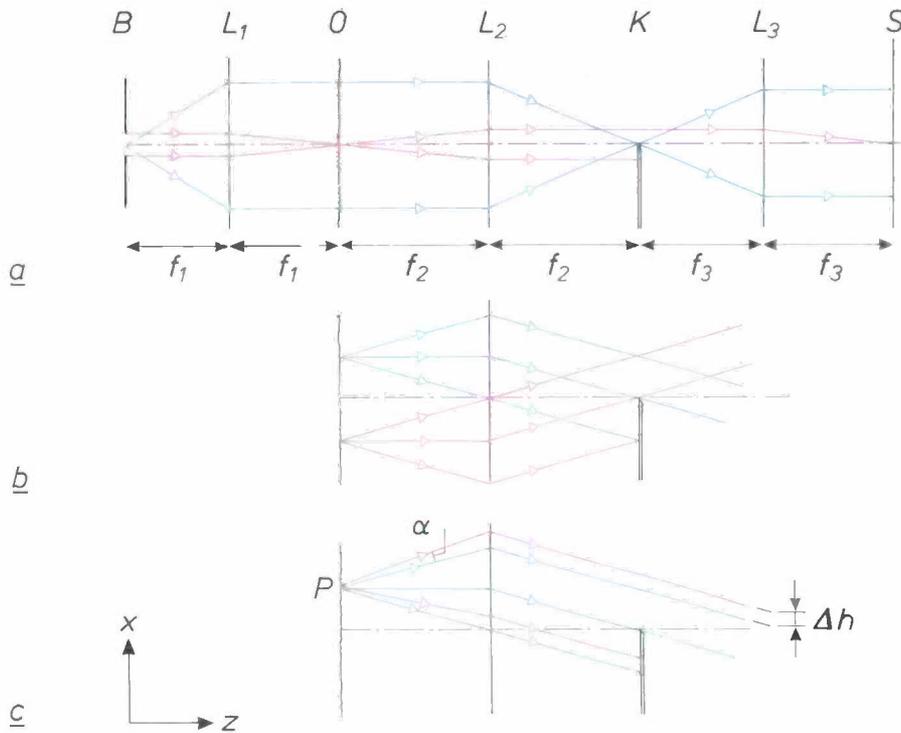


Fig. 2. a) Example of a configuration for the Schlieren method. An image of a rectangular light source *B* is formed by two lenses (*L*₁ and *L*₂ with focal lengths *f*₁ and *f*₂) at the focal plane of the second lens (blue beam). Between *L*₁ and *L*₂ there is a transparent object *O*. An image of the object is produced on the screen *S* by the lenses *L*₂ and *L*₃ (red beam). At the focus of the second lens there is a knife-edge *K*. The coordinate perpendicular to the knife-edge is *x*, and *z* indicates the direction perpendicular to the object. b) Half of a light beam originating from an arbitrary point *P* of the object is screened off by the knife-edge. c) If there is an inhomogeneity in the object at the position *P*, a light beam is deflected through an angle α . The result is that less than half of the beam originating from *P* is screened off by the knife-edge. The image of the light source produced at the focal plane of the second lens by the light incident on the object through point *P* is displaced by a distance Δh . The inhomogeneity is visible on the screen as a streak brighter than its surroundings.

fig. 3, with a planar wavefront *WF* incident on it at time *t*. The position of the wavefront is indicated at a number of times separated by a fixed interval Δt . In equal intervals Δt the wavefront travels identical optical pathlengths. (The optical pathlength is the product of distance and refractive index.) The object is made of a material with a refractive index *n* higher than that of the surrounding air. The velocity of the wavefront in the object is $1/n$ times the velocity in air. Consequently, at positions where the object is thinner the wavefront reaches the position $z = d$ earlier. The distance the wavefront has travelled at $x = 0$ between $t + \Delta t$ and $t + 3\Delta t$ is exactly equal to the thickness *d* of the object. At $x = \Delta x$ this distance is $d + \Delta d$. The time $2\Delta t$ (the distance travelled divided by the velocity

of the wavefront) is equal to $d/(c/n)$ and also to $\Delta z/c + (d - \Delta z)/(c/n) + \Delta d/c$. The distance Δd by which the wavefront at the position Δx is 'ahead of' the wavefront at $x=0$ is therefore $\Delta d = \Delta z(n - 1)$. The tangent of the angle α between the wavefront and

$$\tan \alpha = \frac{\Delta d}{\Delta x} = \frac{\Delta z(n - 1)}{\Delta x}$$

If α is small — the Schlieren method can only be used if it is — then α is also small and we therefore have:

$$\alpha = \frac{\Delta z(n - 1)}{\Delta x}$$

The deflection of the light rays, which are perpendicular to the wavefronts, is therefore proportional to the ratio $\Delta z/\Delta x$, which is equal to the gradient $\partial z/\partial x$ in the limit for small distances.

At the focal plane of *L*₂ the image formed of the light source by light rays passing through point *P* is

displaced by a distance $\Delta h = \alpha f_2$. This implies that the intensity of the light passing the knife-edge increases because of this gradient. With no gradient the light intensity (I_0) is proportional to $\frac{1}{2}h$, whereas with a gradient ($I_0 + \Delta I$) it is proportional to $\frac{1}{2}h + \Delta h$, where h is the dimension of the image of the light source. The relative increase in intensity at the screen as a result of the gradient is

$$\frac{\Delta I}{I_0} = \frac{\Delta h}{\frac{1}{2}h} = \frac{\alpha f_2}{\frac{1}{2}h} \quad \text{or} \quad \alpha = \frac{h \Delta I}{2f_2 I_0} \quad (1)$$

The gradient can thus be directly determined from the measured relative increase in intensity.

The sensitivity of the method is

$$\frac{d\Delta(I/I_0)}{d\alpha} = \frac{2f_2}{h} = \frac{2f_1}{l} \quad (2)$$

where $h = (f_2/f_1)l$, with l the dimension of the light source perpendicular to the optical axis and perpendicular to the edge. By increasing f_1 or reducing l the sensitivity of the method can be increased, within certain limits. The sensitivity cannot be increased indefinitely, however, because the measurement range also depends on these two quantities. If $\Delta h = \alpha f_2 > \frac{1}{2}h$, the image is either fully illuminated or fully blacked out and the image intensity is no longer linearly dependent on the gradients. In addition, diffraction effects become significant if the dimensions of the light source become too small.

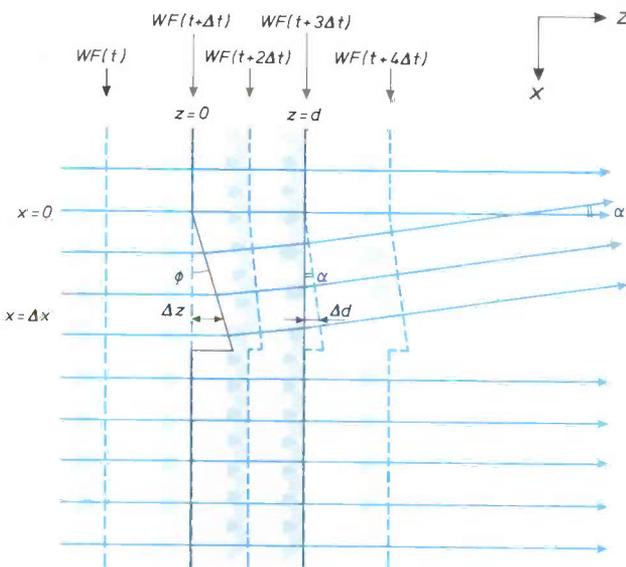


Fig. 3. The relation between the gradient in the optical path and the perturbation of the wavefront. The grey area represents part of a transparent object of thickness d , in which a wedge-shaped piece (with an apex angle ϕ) is missing. A plane wavefront WF is incident on the object from the left at time t_1 . The position and shape of the wavefront are shown at a number of times t_1 to t_5 at intervals of Δt . The blue continuous lines represent the light rays, which are perpendicular to the wavefronts. At the height of the wedge the light rays leave the object at an angle α .

Configuration with real optical components

The above argument only applies for perfect optical components and uniform illumination of the object. In reality, neither condition can be satisfied. Nevertheless the Schlieren method remains useful because these drawbacks can be overcome by making four measurements of the light intensity on the screen S : with and without a knife-edge and with and without an object.

The intensity without knife-edge and without object is given by:

$$I'_0 = c_{x,y}hb, \quad (3)$$

where $c_{x,y}$ is the position-dependent illumination of the plane where the object will be located, and h and b are the height and width of the image of the source that will be produced at the knife-edge. Next, the knife-edge is introduced in such a way that the actual edge is situated on the optical axis of the configuration. Owing to imperfections in the positioning of the lenses, however, the centre of the image does not lie exactly on the edge but is separated from it by a distance $\xi_{x,y}$. With a knife-edge and without object the intensity is:

$$I_0 = c_{x,y}(\frac{1}{2}h + \xi_{x,y})b. \quad (4)$$

The knife-edge is then removed and the object put in its place. The intensity with object and without knife-edge is:

$$I' = c_{x,y}hbT_{x,y}, \quad (5)$$

where $T_{x,y}$ is the transmittance of the object. Finally the intensity is measured with both knife-edge and object present:

$$I = c_{x,y}(\frac{1}{2}h + \xi_{x,y} + \alpha_{x,y}f)bT_{x,y}, \quad (6)$$

where f is the focal length of the second lens. We now have two relative measurements of the light intensity: I_0/I'_0 measured without an object and I/I' measured with an object. The difference $M_{x,y}$ between these measurements is:

$$M_{x,y} = \frac{I}{I'} - \frac{I_0}{I'_0} = \frac{c_{x,y}(\frac{1}{2}h + \xi_{x,y} + \alpha_{x,y}f)bT_{x,y}}{c_{x,y}hbT_{x,y}} - \frac{c_{x,y}(\frac{1}{2}h + \xi_{x,y})b}{c_{x,y}hb} = \frac{\alpha_{x,y}f}{h} \quad (7)$$

Between $M_{x,y}$ and $\alpha_{x,y}$ we therefore find a simple relation in which a correction is made not only for the imperfections of the lenses and the source but also for different diode sensitivities. The only requirement for the diodes is that they should produce a signal that is proportional to the incident light intensities. The measurements without an object are in effect a cali-

bration of the measurement configuration and only have to be made once.

The Schlieren microscope

The Schlieren method is a useful means of measuring thickness differences and refractive-index variations in transparent objects. The principle can also be applied successfully in a microscope. This is not self-evident, because in deriving the quantitative method we used geometrical optics, which cannot always be used directly for the theory of image formation in a microscope.

It is possible to make f_2 and f_3 in fig. 2 large so as to obtain a high magnification. For example, with $f_2 = 30$ mm and $f_3 = 400$ mm we obtain a magnification of 13 times. With the diode dimensions of $13 \mu\text{m} \times 13 \mu\text{m}$ we thus look at 'points' in the object of $1 \mu\text{m} \times 1 \mu\text{m}$. We have then in fact made a microscope. However, we shall have to show in another way that the linear relation between the light intensity in the image and the gradient at the corresponding point in the object, as derived by geometrical optics, is preserved. This can be done by applying Abbe's calculation of image formation in the microscope to the case where half the diffraction pattern formed at the back focal plane of the objective lens (the Fourier plane) is covered and the image is formed by the uncovered half. From the calculation given below it will be seen that the Schlieren theory remains valid as long as the diffraction orders in the Fourier plane overlap [3].

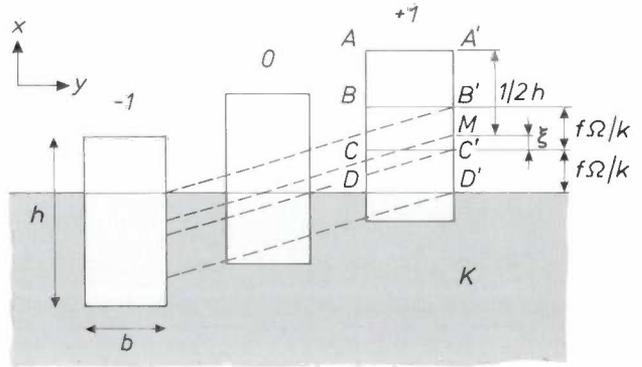


Fig. 4. The location of the diffraction images of order $-1, 0$ and $+1$ relative to the knife-edge. In fact the three images are only displaced at right angles to the edge, but for clarity they are shown here displaced along the edge as well. M indicates the centre of the image of height h , and C is the position of the knife-edge in the diffraction image. There are three regions: $AA'BB'$, where all three diffraction images pass the knife-edge, $BB'CC'$, where only the zero-order and 1st-order diffraction images pass the knife-edge, and $CC'DD'$, where only the 1st-order diffraction image passes the knife-edge. The images are spaced at a distance $f\Omega/k$.

objective three images of the light source are formed, as shown in fig. 4. Light waves originating from the same point in the light source and reaching the same point in the image via different paths (the different orders) will interfere with each other. The amplitudes of the interfering light waves have to be added. After squaring and integrating over the light source we obtain the intensity pattern of the image.

Since one, two or three diffraction orders may pass the knife-edge, there are three different regions in the diffraction pattern:

- In the region between AA' and BB' , where all three orders of the diffraction pattern pass the knife-edge, the amplitude E_{AB} of the wavefront is:

$$E_{AB} = E_0[\cos(\omega t + kz_0) - k\Delta z \cos \Omega x \sin(\omega t + kz_0)],$$

or

$$E_{AB}^2 = E_0^2 [\cos^2(\omega t + kz_0) - 2k\Delta z \cos \Omega x \sin(\omega t + kz_0) \cos(\omega t + kz_0)],$$

with the condition that $k\Delta z \ll 1$. The surface is given by

$$S_{AB} = [\frac{1}{2}h + \xi - f\Omega/k]b,$$

and the total intensity I_{AB} in this region is:

$$I_{AB} = S_{AB} \int E_{AB}^2 dt = \frac{1}{2} \left[\frac{1}{2}h + \xi - \frac{f\Omega}{k} \right] E_0^2 b. \tag{9}$$

- Next we consider the region between BB' and CC' , where only the zero-order and 1st-order diffraction images pass the knife-edge. Here we have:

$$E_{BC} = E_0[\cos(\omega t + kz_0) - \frac{1}{2}k\Delta z \sin(\omega t + kz_0 + \Omega x)],$$

or

$$E_{BC}^2 = E_0^2 [\cos^2(\omega t + kz_0) - k\Delta z (\cos^2(\omega t + kz_0) \sin \Omega x - k\Delta z \sin(\omega t + kz_0) \cos(\omega t + kz_0) \cos \Omega x)].$$

For the surface of this region $S_{BC} = (f\Omega/k)b$. The total light intensity on this area is:

$$I_{BC} = \frac{1}{2} \frac{f\Omega b}{k} E_0^2 [1 - k\Delta z \sin \Omega x]. \tag{10}$$

This condition can be derived as follows. A microscope forms an image of a reflecting object with a surface whose height follows a sine function: $z = z_0 + \frac{1}{2} \Delta z \cos \Omega x$, where Ω is the spatial frequency of the irregularity, z the coordinate parallel to the optical axis and x is the coordinate perpendicular to the optical axis and perpendicular to the knife-edge. An incident plane wave $E = E_0 \cos(\omega t - kz)$, where E is the amplitude and k the circular wave number, is reflected:

$$\begin{aligned} E &= E_0 \cos(\omega t + kz_0 + k\Delta z \cos \Omega x) \\ &= E_0 [\cos(\omega t + kz_0) \cos(k\Delta z \cos \Omega x) \\ &\quad - \sin(\omega t + kz_0) \sin(k\Delta z \cos \Omega x)]. \end{aligned}$$

For $k\Delta z \ll 1$ this becomes

$$E = E_0 [\cos(\omega t + kz_0) - \frac{1}{2}k\Delta z \sin(\omega t + kz_0 + \Omega x) - \frac{1}{2}k\Delta z \sin(\omega t + kz_0 - \Omega x)]. \tag{8}$$

The different terms in this expression represent the orders 0, 1 and -1 of the diffraction pattern. At the focal plane of the objective they are imaged as three points spaced by $f\Omega/k$, where f is the focal length of the objective. In reality there is not just a single plane wave incident on the object, but a beam originating from the light source, which has non-zero dimensions. At the focal plane of the

• Finally, for the region between CC' and DD' , where only the 1st-order image passes the knife-edge:

$$E_{CD} = E_0 \left[-\frac{1}{2} k \Delta z \sin(\omega t + kz_0 + \Omega x) \right],$$

or $E_{CD}^2 = 0$ and hence also $I_{CD} = 0$.

The total intensity of the light that passes the knife-edge is:

$$\begin{aligned} I &= I_{AB} + I_{BC} + I_{CD} \\ &= \frac{1}{2} \left[\frac{1}{2} h - \frac{f\Omega}{k} + \xi + \frac{f\Omega}{k} (1 - k\Delta z \sin \Omega x) \right] E_0^2 b \\ &= \frac{1}{2} E_0^2 b \left(\frac{1}{2} h + \xi - f \frac{\partial z}{\partial x} \right), \end{aligned} \tag{11}$$

from which it appears that the image intensity is indeed proportional to the gradient in the refractive index in the object, except for a constant term.

However, if $\frac{1}{2} h + \xi$ is smaller than $(f\Omega)/k$, then light of the diffraction image of order -1 does not pass the knife-edge. We then have $AB = 0$, so that $I_{AB} = 0$ and $BC = \frac{1}{2} h + \xi$. The total intensity is then:

$$I = I_{BC} = \frac{1}{2} E_0^2 b \left(\frac{1}{2} h + \xi \right) (1 - k\Delta z \sin \Omega x). \tag{12}$$

If we define a critical frequency Ω_{cr} for the spatial frequency of the irregularities on the surface such that light of the image of order -1 only just fails to pass the knife-edge

$$\Omega_{cr} = \left(\frac{1}{2} h + \xi \right) k / f,$$

or

$$k = \frac{\Omega_{cr} f}{\frac{1}{2} h + \xi}.$$

Then

$$I_{BC} = \frac{1}{2} E_0^2 b \left(\frac{1}{2} h + \xi - \frac{\Omega_{cr}}{\Omega} f \frac{\partial z}{\partial x} \right). \tag{13}$$

If the spatial frequency of the irregularities Ω is smaller than Ω_{cr} , then equation (11) applies and the gradient is proportional to the intensity of the image. If this condition is satisfied, the three diffraction images overlap. It is therefore immediately clear that, if the light source is made larger, the condition for ordinary Schlieren calculations is satisfied up to higher frequencies. The critical frequency Ω_{cr} also increases if ξ is made larger. This means that the knife-edge should be adjusted so that part of the image of order -1 passes the knife-edge. If the spatial frequency of the irregularities is greater than Ω_{cr} , then equation (13) applies. There are still intensity differences in the image, but these are no longer proportional to the gradient, since they are attenuated by a factor of Ω_{cr}/Ω with respect to the gradient.

The distance between the zero-order and 1st-order diffraction images of the source because of periodic irregularities in the object is $d = (\lambda/w)f$, where w is the period of the irregularity and λ the wavelength of the light. The dimension of the image of the source perpendicular to the optical axis is h , so that the Schlieren theory applies when $h > (\lambda/w)f$ or $w/\lambda > f/h$. Irregularities at a small spacing (w small) can be seen if f/h is small. Unfortunately, this ratio also sets a

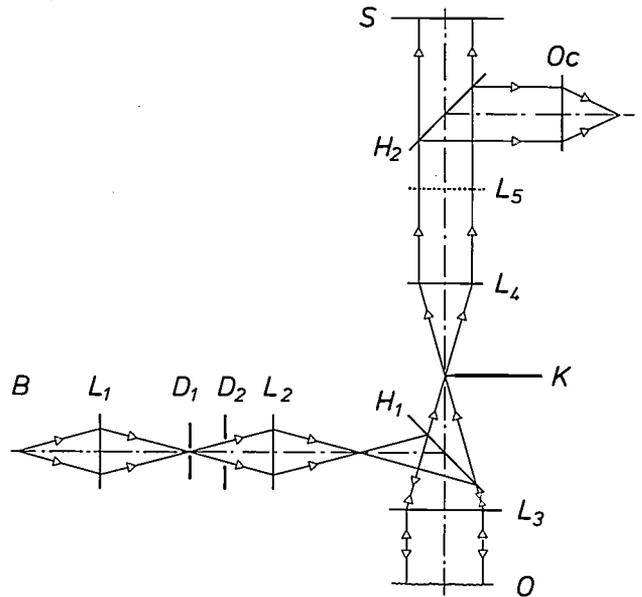


Fig. 5. Principle of the Schlieren microscope. Light from a light source is projected on to the object O by the lenses L_1 and L_2 , the diaphragms D_1 and D_2 and the beam-splitter H_1 . An image of the source is produced by the objective L_3 at the focal plane, where the knife-edge is located. An image of the object is produced on the diode array S by lens L_4 . The second beam-splitter H_2 and the eyepiece Oc are used for focusing at the object. An auxiliary lens L_5 can be used for focusing at the back focal plane of the objective, where the diffraction image can be observed together with the knife-edge.

limit to the sensitivity of the method for determining the magnitude of the gradient (see equation 2) and we cannot therefore go on reducing f/h indefinitely. Practical values in such a measurement are $f = 30$ mm, $h = 3$ mm. The intensity $M_{x,y}$ on the screen can be produced with present-day diode arrays to an accuracy $(\Delta M_{x,y})$ of 5×10^{-3} . The intensity is proportional to the gradient in the optical path $\partial s_{opt}/\partial x$, which, for reflecting surfaces, is twice the roughness $\partial z/\partial x$ on the surface:

$$M_{x,y} = \frac{f}{h} \left(\frac{\partial s_{opt}}{\partial x} \right) = \frac{2f}{h} \left(\frac{\partial z}{\partial x} \right).$$

The accuracy to which the height differences on the surface $[\Delta(\partial z/\partial x)]$ can be determined is therefore directly related to the accuracy to which the intensities can be determined:

$$\begin{aligned} \frac{\partial z}{\partial x} &= \frac{h}{2f} M_{x,y}, \\ \Delta \frac{\partial z}{\partial x} &= \frac{h}{2f} \Delta M_{x,y} = \frac{3}{2 \times 30} 5 \times 10^{-3} = 2.5 \times 10^{-4} \text{ rad}. \end{aligned} \tag{14}$$

The height differences Δz can be found by integrating the gradient $\partial z/\partial x$ in the x -direction over the surface.

[1] These calculations were made by G. Bouwhuis of Philips Research Laboratories.

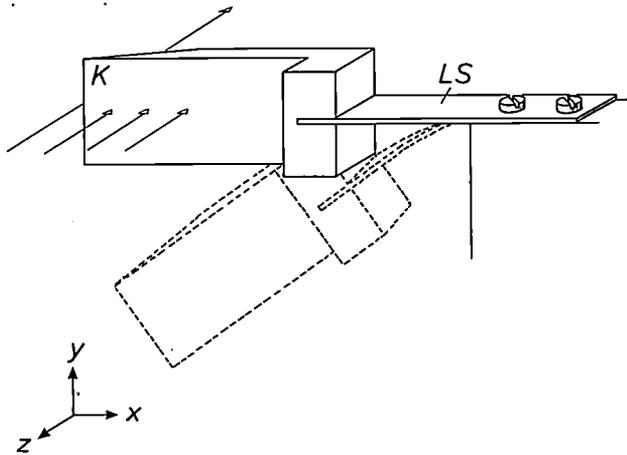


Fig. 6. The knife-edge K , which is connected to the system by a leaf spring LS , is correctly positioned by energizing a magnet. The movement takes place in the direction of the edge (y -direction), perpendicular to the optical axis (z -direction) and takes 50 ms. The edge of the knife is ground to prevent unwanted reflections in the direction of the diode array.

The inaccuracy of measurements over a distance of $1 \mu\text{m}$ (integration over one diode) is

$$\Delta z = \Delta \left(\frac{\partial z}{\partial x} \right) \Delta x.$$

In integration over a number of diodes the inaccuracy increases with the square of the number of measured points (statistically independent measurements), so that

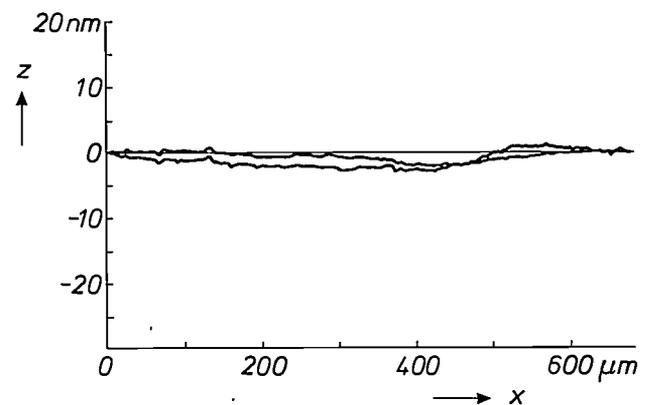
$$\Delta z = \sqrt{n} \Delta \left(\frac{\partial z}{\partial x} \right) \Delta x.$$

Roughness measurements over a distance of $10 \mu\text{m}$ can be made to an accuracy of 0.75 nm . In profile measurements over a distance of 500 nm height differences can be determined to an accuracy of 5 nm .

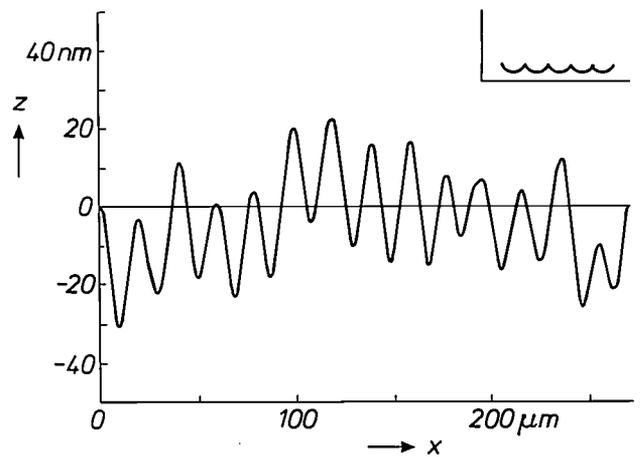
The smallest distance in the object over which full contrast can be observed is $w = (f/h)\lambda = (30/3)0.8 \mu\text{m} = 8 \mu\text{m}$. Since the contrast decreases very slowly with w , the resolution is of the order of $5 \mu\text{m}$.

A Schlieren microscope is illustrated schematically in *fig. 5*. The microscope has a conventional illumination system, consisting of a lamp, two lenses, an aperture diaphragm and a field diaphragm. The only difference from a conventional microscope is that the objective is replaced by an achromatic lens with a focal length of 30 mm and a numerical aperture of 0.2 . This is necessary because it is not possible to place the knife-edge in the focal plane of a conventional microscope objective. The total magnification produced by the objective and the second lens of the imaging system is 13 times. The first beam-splitter, H_1 , is used for illuminating the object, the second one, H_2 , is required for viewing the image through an

eyepiece and for focusing the microscope on the object. Since an experiment consists of four measurements, two with the knife-edge and two without it, it must be possible to reposition the knife-edge at exactly the same place after it has been removed from the microscope. The light intensity is measured to an accuracy of 0.1% , so that the reproducibility of the position of the knife-edge must also be 0.1% . This means that the position of the knife-edge must be adjustable within $1 \mu\text{m}$ of the optical axis, i.e. 0.1% of the dimension of the image of the source. The design of the knife-edge can be seen in *fig. 6*. It is attached to the optical system by two screws and a leaf spring. This meets the requirement that the knife-edge must



a



b

Fig. 7. a) A measurement of the height differences (z) of a highly polished surface as a function of position (x) on a line across the surface. The two curves represent two measurements made with the Schlieren microscope along the same line on the surface. In this case the surface asperities are of interest. The results of the measurements are digitally filtered with a highpass filter. This means that height variations of long wavelength (profiles) cannot be accurately reproduced. This appears from the fact that the two curves do not coincide, whereas small variations (roughness) in the two curves are indicated in the same way. b) A measurement of the height differences (z) on a surface that has been machined on a high-precision lathe by a diamond tool of radius 1 mm . This operation produces a pattern of grooves on the surface (see inset), with a pitch of $20 \mu\text{m}$. The sharp edges found in such patterns are not visible because a lowpass filter cuts off the high spatial frequencies in these measurements.

always coincide with the optical axis to an accuracy of better than $1\ \mu\text{m}$.

The optical system alone as described above is not sufficient for making quantitative measurements. The light intensities are determined by a linear array of 1024 diodes, with a dynamic range better than 0.1%. The measured values are converted by an analog/digital converter into words of 10 bits and stored in the memory of a microcomputer. The microcomputer controls the measurements, checks the state of the components of the system and performs simple arithmetical operations (e.g. averaging). The actual calculations, the conversion of light intensities into gradients and height differences, and the numerical processing of the profile data, e.g. filtering or the determination of autocorrelation diagrams or Fourier spectra of the profile, is performed on a minicomputer. *Fig. 7* gives two examples of results obtained with the Schlieren microscope. *Fig. 7a* shows the roughness of a very smoothly polished surface, and *fig. 7b* the pattern of grooves produced on a surface by a tool in a high-precision lathe.

The employment of the Schlieren method in a microscope is only one of the possible applications. Similar methods (e.g. the Foucault method) can be used for quantitative measurements in the same way as described above. The method can also prove useful in cases where it is necessary to determine whether products meet certain criteria. In such a situation a perfect workpiece can be used as a reference surface. Any variations from the standard can then be immediately identified.

Summary. The Schlieren method, originally only used as a qualitative method for displaying the presence of gradients in the refractive index of transparent objects, is used for determining the shape and roughness of surfaces. A gradient in the refractive index is made visible by screening off a part of a light beam incident on an object. Accurate diodes for measuring the light intensity and the use of computers now make it possible to employ this method for quantitative determinations as well. The optical components do not have to satisfy any special conditions. Differential measurements are used, to correct for any imperfection in these components. The principle of the method has been successfully applied in a microscope. It is easy to use and insensitive to vibration, and therefore suitable for use in production conditions.



Television test decor

Picture quality with existing television technology is usually good. However, certain types of scene are notoriously difficult to capture and reproduce truly faithfully. Difficulties can arise if there are very large contrasts in brightness, strongly saturated colours, finely detailed patterns or highly reflecting surfaces. Research on new and better television systems does not only consist of theoretical studies, there are extensive experimental tests as well. To give the clearest possible

comparison with existing systems, special studio decors are used. These contain as much 'difficult' scene material as possible, and the visual effect can be rather startling. The photograph shows some experimental test decor recently in use in the television studio at Philips Research Laboratories, Eindhoven. (In the next issue of this journal we shall give a detailed account of a new high-definition television system, HD-MAC, which gives greatly improved picture quality.)

Scientific publications

These publications are contributed by staff of laboratories and plants that form part of or cooperate with enterprises of the Philips group of companies, particularly by staff of the research laboratories mentioned below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	<i>A</i>
Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	<i>B</i>
Philips Natuurkundig Laboratorium, Postbus 80000, 5600 JA Eindhoven, The Netherlands	<i>E</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	<i>R</i>

W. J. Bartels, J. Hornstra & D. J. W. Lobeek <i>E</i>	X-ray diffraction of multilayers and superlattices	Acta Crystallogr. A42	539-545	1986
R. G. Gossink <i>E</i>	Polymers for audio and video equipment	Angew. Makromol. Chem. 145/146	365-389	1986
U. Enz <i>E</i>	Bloch walls, solitons, particles: an analogy	Ann. Fond. Louis de Broglie 11 (no. 2)	87-100	1986
M. J. Verkerk & I. J. M. M. Raaymakers <i>E</i>	Characterization of the topography of vacuum-deposited films. 1: Light scattering	Appl. Opt. 25	3602-3609	1986
I. J. M. M. Raaymakers & M. J. Verkerk <i>E</i>	Characterization of the topography of vacuum-deposited films. 2: Ellipsometry	Appl. Opt. 25	3610-3615	1986
Luisa Gonzalez (<i>Univ. Madrid</i>), J. B. Clegg, D. Hilton, J. P. Gowers, C. T. Foxon & B. A. Joyce <i>R</i>	Silicon migration during MBE growth of doped (Al, Ga)As films	Appl. Phys. A 41	237-241	1986
J. Samitier*, J. R. Morante* (<i>Univ. de Barcelona</i>), L. Giraudet & S. Gourrier <i>L</i>	Optical behavior of the U band in relation to EL2 and EL6 levels in boron-implanted GaAs	Appl. Phys. Lett. 48	1138-1140	1986
G. G. P. van Gorkom, A. van Oostrom, J. E. Crombeen & A. M. E. Hoeberechts <i>E</i>	Cesium migration and equilibrium in a strong electric field on the surface of silicon	Appl. Phys. Lett. 49	507-509	1986
A. H. van Ommen, B. H. Koek & M. P. A. Vieggers <i>E</i>	Amorphous and crystalline oxide precipitates in oxygen implanted silicon	Appl. Phys. Lett. 49	628-630	1986
A. H. van Ommen, B. H. Koek & M. P. A. Vieggers <i>E</i>	Ordering of oxide precipitates in oxygen implanted silicon	Appl. Phys. Lett. 49	1062-1064	1986
J. J. P. Bruines, R. P. M. van Hal, H. M. J. Boots, A. Polman* & F. W. Saris* (<i>FOM, Amsterdam</i>) <i>E</i>	Time-resolved reflectivity measurements during explosive crystallization of amorphous silicon	Appl. Phys. Lett. 49	1160-1162	1986
C. Ronse <i>B</i>	Definitions of convexity and convex hulls in digital images	Bull. Soc. Math. Belg. B 37	71-85	1985
J. A. Pals <i>E</i>	Basic properties in semiconductor devices	Crystalline semiconducting materials and devices, P.N. Butcher <i>et al.</i> (eds), Plenum, New York	507-548	1986
C. Niessen <i>E</i>	Abstraction requirements in hierarchical design methods	Design Methodologies, S. Goto (ed.), Elsevier Science, Amsterdam	151-182	1986

- | | | | | |
|---|--|--|-----------|------|
| H. J. Schouwenaars, E. C. Dijkmans, B. M. J. Kup* & E. J. M. van Tuijl*
(*Philips Elcoma Div., Nijmegen) E | A monolithic dual 16-bit D/A converter | IEEE J. SC-21 | 424-429 | 1986 |
| M. Verhaegen & P. van Dooren B | Numerical aspects of different Kalman filter implementations | IEEE Trans. AC-31 | 907-917 | 1986 |
| T. A. C. M. Claasen & W. F. G. Meclenbräucker E | Authors' reply to "Comments on 'Comparison of the convergence of two algorithms for adaptive FIR digital filters'" | IEEE Trans. ASSP-34 | 202-203 | 1986 |
| A. J. E. M. Janssen, R. N. J. Veldhuis & L. B. Vries E | Adaptive interpolation of discrete-time signals that can be modeled as autoregressive processes | IEEE Trans. ASSP-34 | 317-330 | 1986 |
| T. Krol E | (N,K) concept fault tolerance | IEEE Trans. C-35 | 339-349 | 1986 |
| W. J. van Gils E | A triple modular redundancy technique providing multiple-bit error protection without using extra redundancy | IEEE Trans. C-35 | 623-631 | 1986 |
| J. W. M. Bergmans E | Density improvements in digital magnetic recording by decision feedback equalization | IEEE Trans. MAG-22 | 157-162 | 1986 |
| F. M. Dekking (Univ. of Technol., Delft) & P. J. van Otterloo E | Fourier coding and reconstruction of complicated contours | IEEE Trans. SMC-16 | 395-404 | 1986 |
| A. Huizing E | Bereiding van zuiver kwartsglas door superkritisch drogen | Ing. Inf. Procestech-nol. No. 10, 1986 | 19-23 | 1986 |
| E. H. L. Aarts, F. M. J. de Bont, E. H. A. Habers & P. J. M. van Laarhoven E | Parallel implementations of the statistical cooling algorithm | Integration 4 | 209-238 | 1986 |
| J. G. Beerends* & A. J. M. Houtsma*
(*Inst. Perception Res., Eindhoven) | Pitch identification of simultaneous dichotic two-tone complexes | J. Acoust. Soc. Am. 80 | 1048-1056 | 1986 |
| J. 't Hart (Inst. Perception Res., Eindhoven) | Declination has not been defeated — A reply to Lieberman <i>et al.</i> | J. Acoust. Soc. Am. 80 | 1838-1840 | 1986 |
| J. A. Geurst E | Momentum-flux condition for Landau-Squire jet flow | J. Appl. Math. & Phys. 37 | 666-672 | 1986 |
| R. Schäfer (Philips Lighting Div., Eindhoven) & H.-P. Stormberg A | Particle number and pressure determination in high-pressure lamps | J. Appl. Phys. 60 | 1263-1268 | 1986 |
| F. J. A. den Broeder E | Layered magnetic domains in Fe-Cu multilayer films imaged by Schlieren-Lorentz microscopy | J. Appl. Phys. 60 | 3381-3383 | 1986 |
| M. H. L. M. van den Broek E | Electron-optical simulation of rotationally symmetric triode electron guns | J. Appl. Phys. 60 | 3825-3835 | 1986 |
| E. T. J. M. Smeets & A. M. W. Cox E | Influence of alkyl substituents of OMs and operating pressure on the quality of $\text{In}_x\text{Ga}_{1-x}\text{As}/\text{InP}$ heterostructures grown by OMVPE | J. Cryst. Growth 77 | 347-353 | 1986 |
| J. P. André, E. Dupont-Nivet, D. Moroni, J. N. Patillon, M. Erman & T. Ngo L | GaInP-AlGaInP-GaAs heterostructures grown by MOVPE at atmospheric pressure. | J. Cryst. Growth 77 | 354-359 | 1986 |
| B. J. Fitzpatrick, T. F. McGee III & P. M. Harnack N | Self-sealing and self-releasing technique (SSSR) for the crystal growth of II-VI compounds | J. Cryst. Growth 78 | 242-248 | 1986 |
| A. G. Tangena & G. A. M. Hurkx E | The determination of stress-strain curves of thin layers using indentation tests | J. Eng. Mater. & Technol. 108 | 230-232 | 1986 |
| R. B. Helmholtz (ECN, Petten), K. H. J. Buschow E | A neutron diffraction and magnetization study of PdMnTe | J. Less-Common Met. 123 | 169-173 | 1986 |
| J. J. van den Broek, J. L. C. Daams & R. Coehoorn E | Early proof of fivefold symmetry | J. Less-Common Met. 123 | L5-L7 | 1986 |
| P. J. Severin, A. P. Severijns & C. H. van Bommel E | Passive components for multimode fiber-optic networks | J. Lightwave Technol. LT-4 | 490-495 | 1986 |
| C. M. G. Jochem & J. W. C. van der Ligt E | Cooling and bubble-free coating of optical fibers at a high drawing rate | J. Lightwave Technol. LT-4 | 739-742 | 1986 |
| G.-D. Khoe & A. H. Dieleman E | TTOSS: an integrated subscriber system for direct and coherent detection | J. Lightwave Technol. LT-4 | 778-784 | 1986 |

D. J. Broer & G. N. Mol	E	Fast curing primary buffer coatings for high strength optical fibers	J. Lightwave Technol. LT-4	938-941	1986
G.-D. Khoe, J. A. Luijendijk & L. J. C. Vroomen	E	Arc-welded monomode fiber splices made with the aid of local injection and detection of blue light	J. Lightwave Technol. LT-4	1219-1222	1986
P. Schobinger-Papamantellos (<i>Inst. für Kristallogr. und Petrogr., Zürich</i>) & K. H. J. Buschow	E	A neutron diffraction and magnetic study of the first-order phase transition in $TbGe_{1-x}Si_x$ ($0 \leq x \leq 0.4$)	J. Magn. & Magn. Mater. 62	15-28	1986
M. J. Verkerk & W. A. M. C. Brankaert	E	Optical properties of reactively evaporated aluminium films	J. Mater. Sci. Lett. 6	115-117	1987
J. W. M. Jacobs	E	Photochemical nucleation and growth of Pd on TiO_2 films studied with electron microscopy and quantitative analytical techniques	J. Phys. Chem. 90	6507-6517	1986
A. A. van Gorkum & L. C. M. Beirens	E	Experiments on eliminating spherical aberration in electron guns using aspherical mesh lenses	J. Vac. Sci. & Technol. A 4	2297-2306	1986
M. J. Verkerk & G. J. van der Kolk	E	Effects of oxygen on the growth of vapor-deposited aluminium films	J. Vac. Sci. & Technol. A 4	3101-3105	1986
R. A. de Groot (<i>Univ. Nijmegen</i>), A. M. van der Kraan (<i>Interuniv. Reactor Inst., Delft</i>) & K. H. J. Buschow	E	FeMnSb: A half-metallic ferrimagnet	J. Magn. & Magn. Mater. 61	330-336	1986
D. Washington	R	Colour display using the channel multiplier CRT	Jpn. Disp.	218-221	1986
W. J. A. Goossens	E	Tilt and temperature dependence of the pitch in the chiral smectic C phase	Liquid Cryst. 1	521-528	1986
H. Baumgart, F. Philipp (<i>Max Planck Inst., Stuttgart</i>) S. Ramesh, B. Khan, A. Martinez & E. Arnold	N	Twin stabilized planar growth of SOI films	Mater. Res. Soc. Symp. Proc. 53	65-70	1986
B. A. Khan, T. Marshall, E. Arnold & R. Pandya	N	Leakage currents in p-channel accumulation mode TFT's	Mater. Res. Soc. Symp. Proc. 53	435-439	1986
G. J. van der Kolk & M. J. Verkerk	E	Recrystallization of oxygen contaminated Al films	Mater. Res. Soc. Symp. Proc. 71	357-362	1986
K. H. J. Buschow	E	New permanent magnet materials	Mater. Sci. Rep. 1	1-63	1986
R. L. Bronnes & R. C. Sweet	N	The metallographic examination of an aluminium-titanium heat exchanger	Microstruct. Sci., Vol. 14, M.R. Louthan <i>et al.</i> (eds), Am. Soc. Met., Metals Park, OH	181-187	1986
L. J. van der Pauw	E	Light absorption in a bent dielectric slab with a lossy coating	Philips J. Res. 41	431-444	1986
C. Colinet (<i>Lab. Thermodyn. & Physico-Chimie Metallurgique, St. Martin d'Hères</i>) & K. H. J. Buschow	E	Formation enthalpies of Gd-Pd and Gd-Pt compounds	Philips J. Res. 41	445-451	1986
B. J. Fitzpatrick, P. M. Harnack & S. Cherin	N	Crystal growth of cadmium chalcogenides by the SSSR-zone melting method	Philips J. Res. 41	452-459	1986
D. Snyers & A. Thayse	B	Theorem proving techniques and P-functions for logic design and logic programming	Philips J. Res. 41	460-505	1986
P. J. Severin	E	Self-routing fibre-optic networks and switches using multi-tailed receiver/transmitter units	Philips J. Res. 41	507-530	1986
J. Bergmans	E	Discrete-time models for digital magnetic recording	Philips J. Res. 41	531-558	1986
T. N. Saadawi (<i>City Univ., New York</i>), N. Jain & M. Schwartz (<i>Columbia Univ., New York</i>)	N	Protocol for a distributed switching interactive CATV network	Philips J. Res. 41	559-575	1986
R. S. Prodan	N	Multidimensional digital signal processing for television scan conversion	Philips J. Res. 41	576-603	1986
C. B. Marshall	R	A radio-paging receiver architecture and demodulator	Philips J. Res. 41, Suppl. no. 2	128 pp.	1986

D. C. Rogers*, J. Singleton*, R. J. Nicholas* (*Univ. Oxford), C. T. Foxon & K. Woodbridge	R	Magneto-optics in GaAs-Ga _{1-x} Al _x As quantum wells	Phys. Rev. B 34	4002-4009	1986
G. N. A. van Veen, F. H. M. Sanders, J. Dieleman, A. van Veen (Interuniv. Reactor Inst., Delft) D. J. Oostra* & A. E. de Vries* (*FOM, Amsterdam)	E	Anomalous time-of-flight distributions observed for argon implanted in silicon and resputtered by Ar ⁺ ion bombardment	Phys. Rev. Lett. 57	739-742	1986
F. J. C. M. Toolenaar & G. de With	E	Translucent Y ₃ Al ₅ O ₁₂ . The influence of process parameters	Proc. Br. Ceram. Soc. 37	241-246	1986
H. A. van Sprang & J. L. M. van de Venne	E	Switching times for display application of the cholesteric to nematic phase transition	Proc. Eurodisplay Conf., Paris 1984	207-210	1984
H. Sari, L. Desperben & S. Moridi	L	Optimum timing recovery for digital equalizers	Proc. GLOBECOM '85, New Orleans 1985	1460-1465	1985
A. F. de Jong, W. Coene* & D. van Dijk* (*Univ. Antwerpen)	E	<i>New methods for the construction of the phase-object function, to be used in dynamical electron diffraction calculations</i>	Proc. Int. Cong. on Electron Microscopy, Kyoto 1986	503-504	1986
B. H. Verbeek	E	Coherence properties and l.f. noise in AlGaAs lasers with optical feedback	Proc. SPIE 587	93-98	1986
A. A. Turnbull	R	The application of heat-collector fins to reticulated pyroelectric arrays	Proc. SPIE 588	38-43	1986
C. H. L. Weijters & W. C. Keur	E	Reduction of reflection losses in solid-state image sensors	Proc. SPIE 591	75-79	1986
J. O. Voorman	E	Analog integrated filters or continuous-time filters for LSI and VLSI	Rev. Phys. Appl. 22	3-14	1987
A. J. Kil*, R. J. J. Zijlstra*, P. M. Koenraad* (*Univ. Utrecht), J. A. Pals, J. P. André	E,L	Noise due to localized states in the quantum Hall regime	Solid State Commun. 60	831-834	1986
R. Houdré*, C. Hermann*, G. Lampel* (*Ecole Polytech. Palaiseau) & P. M. Frijlink	L	Photoemission study of a single GaAlAs/GaAs/GaAlAs quantum well	Surf. Sci. 168	538-545	1986
E. E. Havinga & L. W. van Horsen	E	Dependence of the electrical conductivity of heavily-doped poly-p-phenylenes on the chain length	Synth. Met. 16	55-70	1986
F. Baaijens	E	On a numerical method to solve contact problems	Thesis, Eindhoven	121 pp.	1987
M. van den Broek	E	The design of rotationally symmetric triode electron guns	Thesis, Eindhoven	127 pp.	1986
B. Greenberg, W. K. Zwicker & I. Cadoff	N	ZnS epitaxy on sapphire (110)	Thin Solid Films 141	89-97	1986
P. C. Zalm	E	Ion-beam assisted etching of semiconductors	Vacuum 36	787-797	1986
J. H. Waszink & M. J. Piena	E	Experimental investigation of drop detachment and drop velocity in GMAW	Weld. J. 65 (Weld. Res. Suppl.)	289s-298s	1986
R. E. van de Leest	E	On the atmospheric corrosion of thin copper films	Werkstoffe & Korrosion 37	629-632	1986

Contents of Philips Telecommunication and Data Systems Review 44, No. 3, 1986

R. F. Wielinga & J. Preston: The new medium of Compact Discs, and introduction to Optical Disc Recording (pp. 2-6)

M. G. Noordenbos: CD-ROM; User aspects and application fields (pp. 7-17)

K. Meissner: Device, system integration and standardization (pp. 18-31)

J. P. H. Maat: Call processing in a fully integrated network of SOPHO-S2500 PABX's (pp. 32-45)

HD-MAC: a step forward in the evolution of television technology

M. J. J. C. Annegarn, J. P. Arragon, G. de Haan,
J. H. C. van Heuven and R. N. Jackson

The first experiments with television were made about a hundred years ago, and were mainly based on mechanical devices such as the Nipkow disc. Then in the thirties of this century purely electronic solutions became available for all the elementary problems of image transfer. From that time, in its triumphal progress, television has marched on, and on. Now there are some 500 000 000 television receivers throughout the world. Yet the end of the road is still over the horizon: with broadcasting satellites, chips and advanced signal-processing methods, new perspectives have opened before us — literally and figuratively. From a purely technical point of view the possibilities for the future seem to be virtually unbounded. However, for quite other reasons (economic ones especially) it is vitally important that the development should be evolutionary rather than revolutionary. In the article below a new television system of such an evolutionary nature is described.

Introduction

A map of the world showing the television system used in each country looks something like a patchwork quilt. To start with, some systems have 30 pictures (frames) per second, with 525 lines per picture, whereas others have 25 pictures per second, with 625 lines per picture. Then there are three different standards at present for transmitting colour pictures: NTSC, PAL and SECAM. Recently, for use in television broadcasting satellites and cable systems, a 'family' of MAC systems (MAC stands for Multiplexed Analog Components) has been added^[1]. There are also many smaller differences, e.g. in the methods used for including the sound, for encoding the stereo sound information and for handling teletext and videotex. Clearly, in spite of all the international agreements and recommendations, the situation is very varied,

and this is far from ideal. The undesirable consequences of this are becoming more and more noticeable as the exchange of picture information by electronic methods becomes more common. Although conversion of one type of television signal into another is always possible in principle (and is often used in practice), complicated professional equipment is generally required, and loss of quality cannot always be avoided.

But another, and in fact more fundamental, criticism can be made of the situation just described. Television is sometimes called 'the window on the world'. However, it is a very small window giving only a limited experience of reality. It can be shown that larger, sharper pictures greatly improve the viewing experience — with a difference that is of the same order of magnitude as was the addition of colour to monochrome TV. Unfortunately the existing TV systems were not designed for producing adequate pic-

Ir M. J. J. C. Annegarn and Ir G. de Haan are with Philips Research Laboratories, Eindhoven; J. P. Arragon, Ingénieur I.S.E.P., is with Laboratoires d'Electronique et de Physique Appliquée (LEP), Limeil-Brévannes, France; Ir J. H. C. van Heuven is with the Consumer Electronics Division, Philips NPB, Eindhoven and R. N. Jackson, F.I.E.E., is with Philips Research Laboratories (PRL), Redhill, Surrey, England.

^[1] H. Mertens and D. Wood, Standards proposed by the EBU for satellite broadcasting and cable distribution, J. IERE 56, 53-61, 1986.

ture quality when used with large-screen displays. What is needed, therefore, is a new and improved television signal that satisfies a 'high-definition' *broadcast* standard and is known as an HDTV broadcast signal.

There is yet another important point: much of the picture material intended for transmission by television is initially recorded ('produced') on photographic film, because the quality on film is superior to that of the present television standards. It will only be possible to change this situation by using a considerably improved and as yet not agreed (electronic) 'high-definition' *production* standard.

There is therefore a clear need for a new television standard that will make HDTV possible, and it is most important to achieve the best possible world-wide consensus here. However, it is at present a matter of debate whether the HDTV broadcast standard should in fact be exactly the same as the HDTV production standard, in spite of the close relationships between the two. For an HDTV broadcast standard it is important that the vast quantities of existing consumer equipment should not suddenly be made obsolete: existing equipment must be capable of processing any improved signals to give the quality now attainable. In other words, this standard should permit a gradual ('evolutionary') progress towards equipment with improved picture quality [2][3]. This could perhaps imply that a single world standard for HDTV broadcasting might not be sufficient: for example in relation to the picture frequency it might be better to consciously retain the existing 25-Hz/30-Hz frontier. The HDTV production standard, on the other hand, should be so universal that it will permit conversion without loss of picture quality to any other existing standard (thus including any accepted HDTV broadcast standard).

In this article we should like to show, as an example, how an HDTV broadcast standard could be based on the standardized MAC family. First, however, we shall pause briefly to look at HDTV in general. Next, we shall give a short description of the existing MAC family, which we shall usually call the 'MAC/packet' system from now on. The remaining — and longest — part of this article will be devoted to the equipment and the special signal-processing operations that are required to obtain the desired high-definition quality with this system.

High-definition television (HDTV)

The term 'high-definition television' is not new; it was used some fifty years ago by a British government committee to refer to all television systems that used more than 240 lines. Today we take it to mean something quite different. HDTV means satisfying the need

for an enhanced viewing experience by providing large, high-quality pictures. To do this we must improve present television quality in the following respects:

- an improvement in the resolution in the vertical and horizontal directions by a factor of about two;
- suppression of the undesirable interaction between chrominance and luminance information ('cross-colour' and 'cross-luminance');
- an increase in the aspect ratio (this is the ratio of picture width to picture height, and is now 4:3) to 16:9 to give greater compatibility with cine images and the characteristics of the human eye;
- stereophonic sound with 'hi-fi' quality.

Experiments have shown that the first two improvements allow a considerably larger picture to be used for the display, and in combination with the change in aspect ratio this gives a greater involvement of the viewer. None of the three older colour-television standards PAL, NTSC and SECAM is suitable for the full attainment of the above improvements. Fortunately, however, this does seem to be possible within the newer MAC television standard. This opens the way for a gradual transition to HDTV, with existing equipment remaining fully serviceable and providing the quality for which it was originally designed [4].

Multiplexed Analog Components (MAC)

With the prospect of geostationary terrestrial satellites as 'suspended' television transmitters, which transmit the signal directly to the individual television viewers ('direct-broadcasting satellites'), the need arose a few years ago for a new television broadcast standard. As far back as 1977 international agreements were concluded at the World Administrative Radio Conference (WARC) about the available frequency band and the permitted signal level for each television channel. The European countries were each allocated five channels with a bandwidth of 27 MHz in the frequency range around 12 GHz. For these channels, mainly because of the maximum permissible transmitted power, frequency modulation (FM) is the best modulation method. This means, however, that the noise appearing on reception has the characteristic triangular spectral shape associated with FM, so that the noise power increases linearly with the frequency. In the existing PAL/NTSC/SECAM systems, however, it is the higher frequencies that carry the chrominance and sound information, so that these are relatively the most affected by this FM noise (*fig. 1*).

Quite apart from this, the long range of satellite transmitters creates almost automatically a need for more sound channels and more arrangements for sub-

titling than exist in conventional transmitters. Moreover, it is highly desirable to be able to 'encode' television signals in a simple and reliable way so that certain programmes can only be obtained after extra payment (pay TV and subscriber TV). A solution for all these problems and special features has been found in the form of the MAC/packet system. This system has the same scanning structure as conventional TV systems, i.e. 625 lines per picture, 25 pictures per second and 2:1 interlacing [6]. The most fundamental characteristic of a MAC television signal, however, is that all the types of information — such as luminance, chrominance and sound — occur alternately (in 'time-division multiplex') and only become simultaneously available again on display. To make this possible the conventional luminance signal Y and the conven-

tional chrominance signals U and V are compressed in time, so that they then fit together, one after the other, into the same time interval as before ('the line period'; fig. 2). Each line period, moreover, contains only U or V , in alternate sequence. Because of the time compression the bandwidth increases proportionately. Two possible combinations of compression factors have been selected: a compression factor of $3/2$ for Y combined with a compression factor of 3 for U and V or a factor of $5/4$ for Y and a factor of 5 for U and V . In this latter case, a greater bandwidth for the Y information is available for a given bandwidth of the compressed signal, at the expense of a reduced bandwidth for U and V .

One of the great advantages of a MAC television signal is that in principle there is now no question of cross-colour and cross-luminance, because luminance information and chrominance information are no longer present simultaneously in the signal. Moreover, it is also easier to take advantage of the full bandwidth of the luminance and chrominance signals after reception, so that a MAC signal alone (i.e. without specific HDTV measures) can give a better picture quality than existing systems.

During each line period (in the line-flyback time) there is also room for sound information (4 to 8 mono channels), synchronization and various forms of digital information ('data'). The sound is digitally encoded to 'hi-fi' quality. The sound and the data are divided up into 'packets' (whence the name MAC/packet system) of 751 bits. Of these, 31 bits represent auxiliary information (the opening or 'header') and the remaining 720 bits serve as the actual information (see fig. 3). The header mainly indicates the kind of information transmitted in the packet. Each packet extends over more than one line-flyback time. In the field-flyback times more space is reserved for

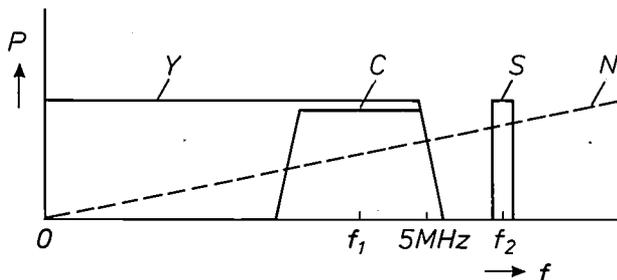


Fig. 1. General diagram of the frequency spectrum P of a colour-television signal in one of the current systems PAL, NTSC and SECAM. Y luminance information. C colour information. S sound information. N triangular noise resulting from frequency modulation. f_1, f_2 subcarrier frequencies for C and S respectively.

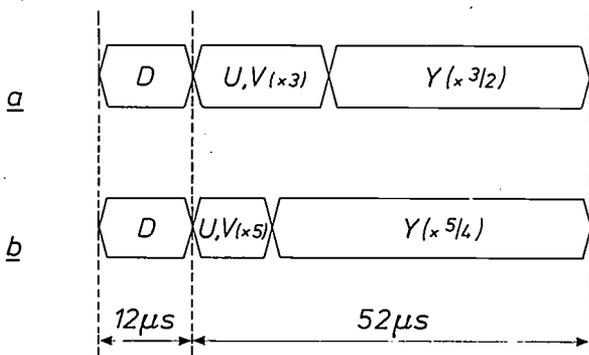


Fig. 2. In the MAC system each line period of $64 \mu s$ contains one of the two colour-difference signals U or V and the luminance signal Y in a compressed form, one after the other. Two combinations of compression factors are possible: a) 3 and $3/2$, or b) 5 and $5/4$. The rest of the line period (the flyback time D) is used for sound signals, synchronization signals, etc.

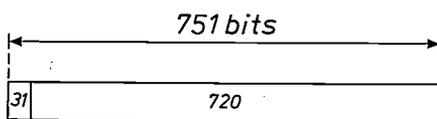


Fig. 3. The sound and other digital information are included in the MAC signal as packets of 751 bits, with 720 bits representing actual information. The bit rate is $10\frac{1}{2}$ or $20\frac{1}{2}$ MHz. One packet is distributed over about 8 or about 4 line-flyback times respectively.

[2] H. Mertens, The future evolution of broadcasting standards (The 1984 Shoenberg Memorial Lecture of the Royal Television Society), *Television (J.R. Telev. Soc.)* 22, 4-12, 1985; C. J. van der Klugt, New television standards: revolution or evolution (The 1985 Shoenberg Memorial Lecture of the Royal Television Society), *Television (J.R. Telev. Soc.)* 23, 6-12, 1986.

[3] C. P. Sandbank and I. Childs, The evolution towards high-definition television, *Proc. IEEE* 73, 638-645, 1985.

[4] M. J. J. C. Annegarn, J. P. Arragon and R. N. Jackson, High definition MAC: the compatible route to HDTV, paper presented at the International Broadcasting Convention, Brighton 1986; G. M. X. Fernando and D. W. Parker, Improved display conversion for high definition MAC, paper presented at the International Broadcasting Convention, Brighton 1986; G. de Haan and W. Crooymans, Subsampling techniques for high definition MAC, paper presented at the International Broadcasting Convention, Brighton 1986; F. Fonsalas and J. Y. Lejard, A method for chrominance contour enhancement applied to HD-MAC television pictures, paper presented at the International Broadcasting Convention, Brighton 1986.

[5] A similar system also exists for 525 lines per picture and 30 pictures per second.

other extra information (such as teletext and subtitles) and — perhaps at least as important — information about the precise way in which the MAC signal has been built up from various signal components (*fig. 4*). For example, in addition to the conventional aspect ratio of 4:3 an aspect ratio of 16:9 is permitted in the MAC/packet system. Since this is possible with each of the two combinations of time-compression factors mentioned earlier, there are already a total of four different kinds of picture coding in the MAC/packet system. The extra information transmitted creates a large measure of flexibility, which at the same time opens the way for a gradual evolution to HDTV display of suitably encoded MAC/packet signals. This will be explained further in the following section.

width of 20 MHz and an aspect ratio of 16:9. A signal that can be transmitted via the standard channel is derived from the camera via an HD-MAC encoder. At the receiving end a picture of HDTV quality can be displayed with the aid of a special HD-MAC decoder and a special receiver; the signal that originates from the standard MAC channel can however also be processed to give the conventional quality by standard MAC equipment. *Fig. 5b* shows clearly the 'downward compatibility' of the HD-MAC approach. The key to proper operation is to be found in the special HD-MAC encoder and decoder. In the following sections we shall therefore look more closely at these circuits and the signal-processing theory on which they are based.

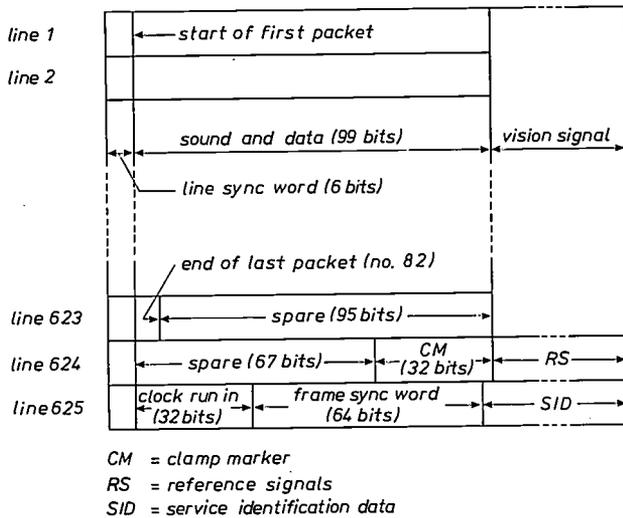


Fig. 4. General arrangement of a complete frame for a member of the MAC/packet family ('D2-MAC'), corresponding to a complete television picture of 625 lines. The intervals that correspond to one line period are shown here one above the other; in fact the picture information occupies about 5/6 of the line period (see *fig. 2*).

High-Definition MAC (HD-MAC)

A schematic, but almost self-explanatory presentation of the possibilities of the MAC signal standard can be found in *fig. 5*. *Fig. 5a* shows a simple MAC television system. A television picture, consisting of 625 lines per picture, with 2:1 interlacing, 50 fields (rasters) per second, a bandwidth of about 5 MHz for the luminance signal and an aspect ratio of 4:3, is generated in a camera. The signals pass through a MAC encoder, a standard MAC channel (e.g. a broadcast satellite) and a MAC decoder, to produce eventually a good picture on a standard receiver. *Fig. 5b* shows an HDTV system in which the same standard MAC channel is used. At the transmitting end, however, an HDTV camera is now used, with 1250 lines per picture, 2:1 interlacing, 50 fields per second, a Y band-

Television signal processing

The essence of generating a television signal is the conversion of an image, generally moving, into an electrical signal. For the moment we shall confine ourselves to the luminance information from the picture, which in both monochrome and colour television is represented by a single signal $s(t)$.

The picture is usually considered as an intensity function $I(x,y,t)$, where x is an independent variable in the horizontal direction, y is an independent variable in the vertical direction and t is time. The conversion of $I(x,y,t)$ into $s(t)$ is made by means of a line scan, in much the same way as our eyes scan when we read the successive pages of a book. In the acquisition of the image (and later when it is displayed) this scan corresponds essentially to a discretization of the variables y and t , or in other words, to a double 'sampling' by vertical position and time. This leads to certain peculiarities that are not fully taken into account in existing television systems. The subdivision of a picture into horizontal lines, for example, means that the television camera cannot distinguish between a uniform white surface and a fine pattern of horizontal black and white lines if it so happens that only the white lines are scanned. Also, the subdivision of a moving scene into a series of successive separate images means that changes that occur in the interval between two images are not observed, whereas certain parts of stationary images tend to show a periodically varying intensity on display ('large-area flicker'). By taking better account of the peculiarities of the scanning process a considerable improvement in the picture quality of television can be obtained [6]. We shall describe this in the section 'Improvement of the resolution in the vertical direction'.

As long as most of the signal processing in television is 'analog', there is no discretization or sampling

in the x -direction. As soon as analog-to-digital conversion comes into play, there is sampling in that direction as well. Then, if the sampling theorem is not satisfied, we run into the danger of certain additional unwanted effects occurring. However, by deliberately applying 'sub-Nyquist sampling', we can again obtain a considerable improvement in picture quality, particularly for stationary images, at practically the same sampling rate. We shall consider this more closely in the section 'Improvement of the resolution in the horizontal direction'.

To describe the various types of signal processing of images and their consequences, we can make good use of the concept of frequency. However, there are now three kinds of frequencies. First of all, there is the vertical frequency f_y , which we express in *cycles per picture height* or c/ph and the horizontal frequency f_x , which we express in *cycles per picture width* or c/pw . In addition to these there is the temporal frequency f_t , expressed in Hz, which indicates how fast the image intensity changes as a function of time. We can illustrate the concepts of horizontal and vertical frequency

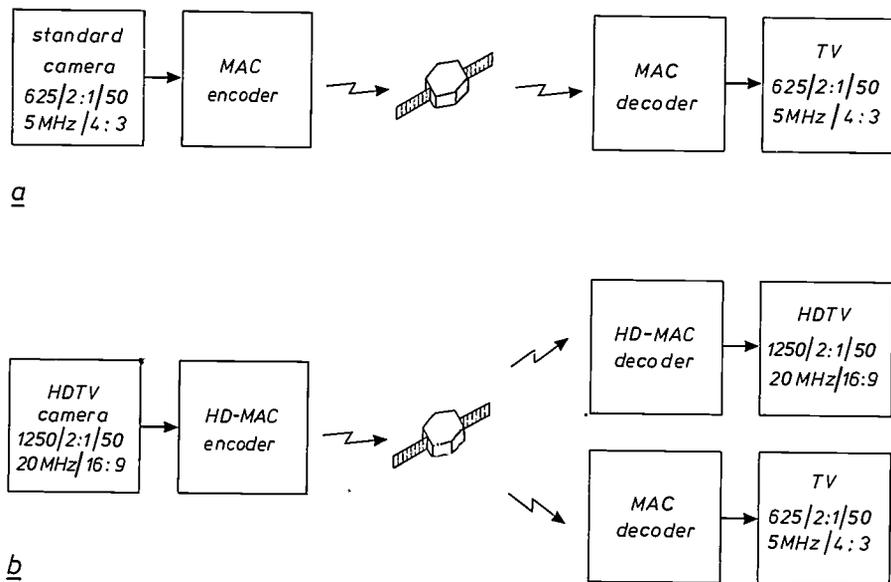


Fig. 5. a) With standard equipment for encoding and decoding MAC signals good television quality can be obtained via broadcasting satellites. b) By using an HDTV camera and an HD-MAC encoder at the transmitting end it is possible — depending on the receiving equipment — to obtain both a standard-quality display and the much better 'high-definition' quality.

		$f_x [c/pw]$		
		0	1	2
$f_y [c/ph]$	0			
	1			

Fig. 6. The concepts of 'horizontal frequency' f_x and 'vertical frequency' f_y , expressed in number of cycles per picture width (c/pw) and number of cycles per picture height (c/ph) can easily be recognized from these simple stationary images. For simplicity the sinusoidal variation in luminance between maximum white and maximum black is replaced here by a stepwise variation.

with the aid of a number of stationary images (so that $f_t = 0$); see fig. 6. Until further notice we shall limit our considerations mainly to stationary and — as we stated earlier — monochrome images.

Television signal spectra

The concept of 'sampling' is of central importance in a detailed analysis of signal processing in television. The sampling theorem for simple one-dimensional signals (such as a mono audio signal) is well known: no information is lost on sampling provided that the sampling frequency is at least twice that of the highest

[6] B. Wendland, High definition television studies on compatible basis with present standards, in: Television technology in the 80's, SMPTE, New York 1981, pp. 151-165;
 B. Wendland, Extended definition television with high picture quality, SMPTE J. 92, 1028-1035, 1983;
 G. J. Tonge, The television scanning process, SMPTE J. 93, 657-666, 1984.

frequency in the sampled signal. The spectrum of the sampled signal then consists of an infinitely extending periodic repetition of the original signal spectrum. The original signal can be reconstructed exactly from this with an ideal lowpass filter *LP*. This is illustrated in *fig. 7a* and *b*. *Fig. 7c* shows a situation in which the sampling theorem is *not* satisfied; now *aliasing* or folding distortion occurs and the original signal cannot be reconstructed. It is also interesting to note here that a shift in the sampling times by half the sampling interval (*fig. 7d*) multiplies alternate periodic repetitions by a factor of -1 . By adding $S_1^*(f)$ and $S_2^*(f)$ from *fig. 7c* and *d* we can however recover $S^*(f)$ from *fig. 7b* exactly, without aliasing. This is because

$$s^*(t) = s_1^*(t) + s_2^*(t),$$

and therefore

$$S^*(f) = S_1^*(f) + S_2^*(f).$$

The same principle can be very usefully applied in television. However, things are rather more complicated because we are dealing with *three* frequencies, f_x , f_y and f_t . These frequencies can be plotted along three axes perpendicular to one another. Any moving scene can thus be represented as a three-dimensional body (a 'three-dimensional spectrum') located around the origin (*fig. 8a*). The bounding planes of this body indicate the limiting values which we wish to consider for each of the three frequencies. A stationary image with limited detail, for example, is characterized by a finite part of the (f_x, f_y) -plane, and a moving image that only consists of infinitely long vertical lines is bounded by the (f_x, f_t) -plane. The different stages in the processing of television signals can now be clearly demonstrated within this context:

- Scanning in a horizontal line pattern is sampling in the vertical direction and gives a periodic repetition of the spectrum along the f_y -axis, with a period proportional to the number of lines f_v .
- Breaking a scene down into successive images is sampling in time and gives a periodic repetition of the spectrum along the f_t -axis, with a period proportional to the picture frequency f_p .
- Any analog-to-digital conversion leads to repetition of the spectrum along the f_x -axis, with the period determined by the number of samples on each line.

If more than one of these three types of sampling occur at the same time, then there is periodic repetition of the spectrum in more than one direction, of course; see *fig. 8b*. In each of these spectral repetitions aliasing can arise (especially in the f_x - and f_y -directions) and this degrades the final picture. If there was no kind of aliasing at all on display, and all the periodic spectral repetitions could be completely re-

moved, then we would be able to recover the original scene absolutely faithfully without line structure or flicker, for example. The objective in HDTV is to approach this situation as closely as possible. Here we make extensive use of operations that we can characterize as manipulations of the spectrum — especially

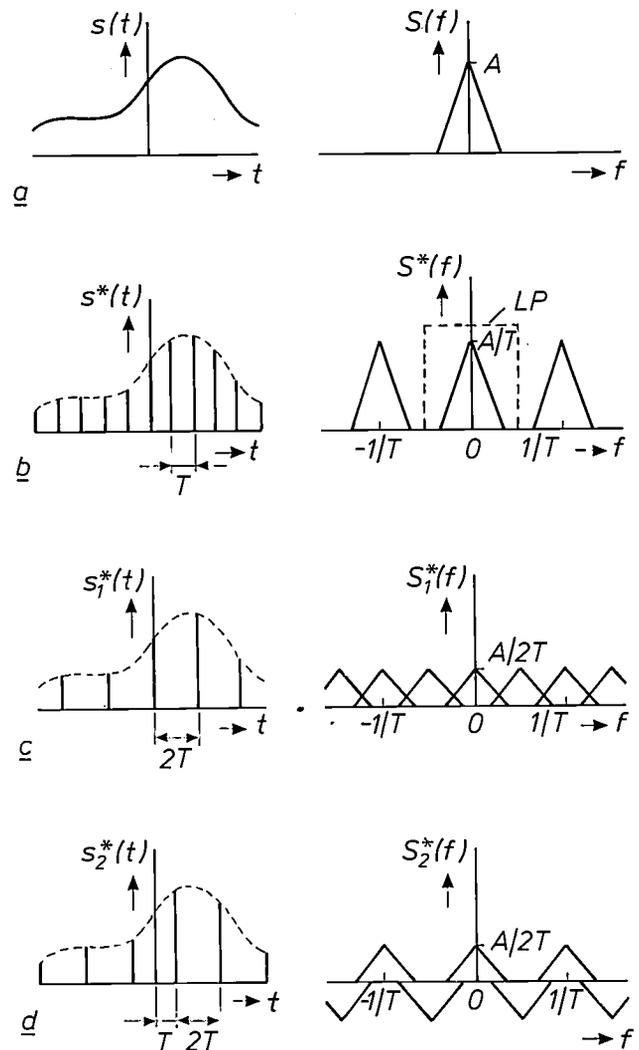


Fig. 7. a) Analog signal $s(t)$ and the associated frequency spectrum $S(f)$. b) After sampling the signal $s^*(t)$ is obtained with frequency spectrum $S^*(f)$. If the sampling theorem is satisfied, the original signal $s(t)$ can be reconstructed with an ideal lowpass filter *LP*. c) and d) If the sampling theorem is not satisfied, there is aliasing in the frequency spectrum: parts of the spectrum overlap. By adding $S_1^*(f)$ and $S_2^*(f)$, however, $S^*(f)$ can be recovered accurately.

as filtering. Therefore it is useful to distinguish between a number of types of filtering, and we should remember that filtering (in the classical sense as well) amounts to the combination (e.g. averaging) of the information from a number of different pixels. If we only combine pixels on the same line, we refer to *horizontal filtering*. We can represent this in the (f_x, f_y, f_t) -space by means of planes that are parallel to the

(f_x, f_y) -plane. The effect is shown in fig. 9a of an ideal horizontal lowpass filter with a cut-off frequency of f_{x1} on the spectrum $S(f_x, f_y, f_t)$ of fig. 8a. In a similar way we refer to *vertical filtering* if the combined pixels

we can obtain various kinds of combined filtering, such as horizontal-temporal or vertical-temporal filtering. The general term *spatio-temporal filtering* is also used in these cases.

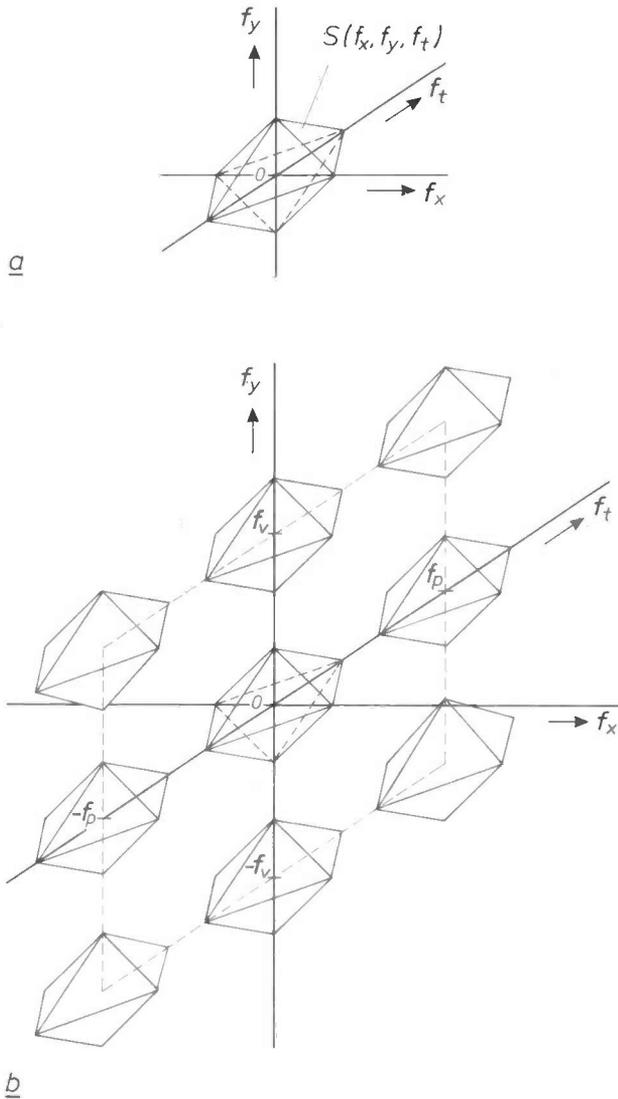


Fig. 8. a) By making use of the horizontal frequency f_x , the vertical frequency f_y and the temporal frequency f_t we can produce a three-dimensional representation of any scene. The surface $S(f_x, f_y, f_t)$ indicates the limiting values in all directions for the frequencies to be considered. b) On sampling there is periodic repetition of the spectrum along the frequency axis corresponding to the dimension (horizontal position, vertical position or time) in which the discretization occurs. Aliasing can occur here, just as in fig. 7. This diagram shows the first spectral repetitions that occur in the discretization of the scene in terms of vertical position and time; it is assumed here that we have f_v lines per picture and f_p pictures per second.

in successive lines lie directly one above or below the other.

The term *temporal filtering* indicates that the only pixels combined in the filtering are those that relate to the same position in successive images. By including more than one kind of pixel in the same filtering pro-

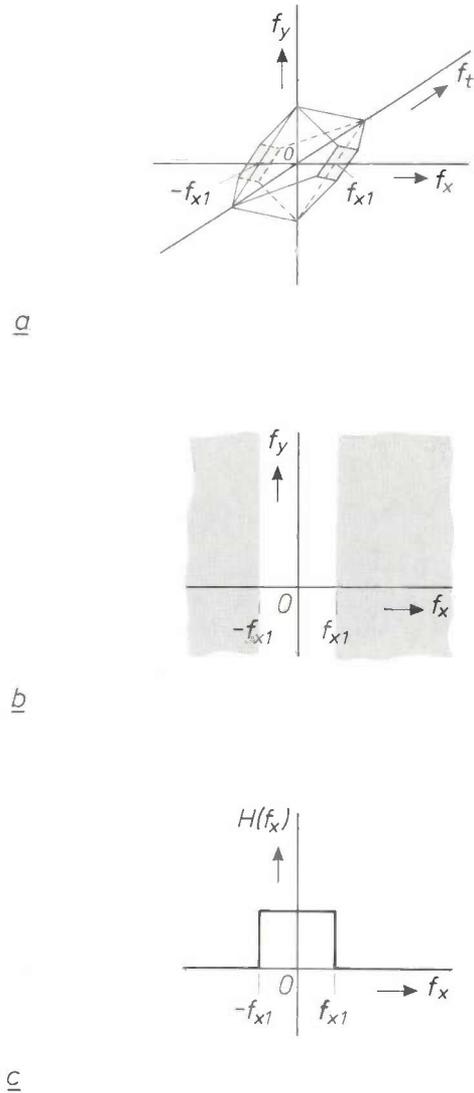


Fig. 9. a) An ideal horizontal lowpass filter with cut-off frequency f_{x1} restricts the spectrum of fig. 8a in the x -direction to the interval $(-f_{x1}, f_{x1})$. b) Two-dimensional representation of the ideal horizontal lowpass filter. c) Conventional one-dimensional representation $H(f_x)$ of the same filter characteristic.

Fortunately we do not always have to deal with three-dimensional spectra; we can frequently make do with a two-dimensional cross-section, since all the relevant information can be derived from this. The ideal horizontal lowpass filter is shown in fig. 9b as a two-dimensional filter. (In this case, in fact, the conventional one-dimensional representation alone is sufficient; fig. 9c.)

Improvement of the resolution in the vertical direction

It has been known since 1934 [7] that a television picture built up from horizontal lines has a much worse resolution in the vertical direction than might at first sight be expected. For example, it has been found from subjective tests that with 100 scanning lines a scene consisting of a maximum of 32 black horizontal lines and 32 white horizontal lines (and not the expected 2×50 lines) can be displayed successfully. This loss of 'sharpness' in the vertical direction we call the 'Kell effect'. It can be expressed by the 'Kell factor', which can vary from 0.5 to 0.7 depending on various circumstances. If the Kell factor could be made equal to 1, it would be possible to achieve a considerable improvement in the quality of the existing television pictures and hence make a significant step in the direction of HDTV. What exactly causes the Kell effect? Let us look at the spectrum in the (f_t, f_y) -plane of a television image consisting of $f_v = 625$ lines, which are scanned sequentially at a picture frequency of 50 Hz. This spectrum looks like the diagram of *fig. 10a*.

It can happen that in considering vertical frequencies account is only taken of the number of *active* lines instead of the *total number* of lines in the picture, as we shall usually do in this article. In a system with 625 lines, for example, the number of active lines is 575; the change to this other approach means that all vertical frequencies should be multiplied by a constant factor of $575/625 = 0.92$.

In yet another approach reference is made to the number of pixels N_v that can be distinguished in the vertical direction; this number corresponds to twice the highest possible vertical frequency $f_{y,\max}$ that can be displayed with the given number of active lines N_{act} . In the most favourable case (Kell factor = 1):

$$f_{y,\max} = N_{\text{act}}/2$$

and therefore

$$N_v = N_{\text{act}}.$$

In general, however, N_v will have a smaller value.

In scanning, interlacing is generally used, to limit the bandwidth of the final television signal. This gives a spectrum like the one shown in *fig. 10b*. In both cases the frequencies that can be observed with the human eye fall within the dashed circle drawn around the origin O . We therefore 'see' not only the desired spectrum (centred on O), but also parts of the spectral repetitions around the points A , B and C . The spectra around A give rise to large-area flicker; the spectra around B are responsible for the visibility of the (static) line structure and the spectra around C lead to line flicker and an apparent vertical movement of the line structure (line crawl). And in fact the situation is worse than it appears from *fig. 10*. In the television pictures now generally used the original spectrum

extends beyond $312\frac{1}{2}$ c/ph in the vertical direction. This means that the 'repeat spectra' also overlap and therefore give rise to aliasing. It is in fact these spectral repetitions and aliasing in the vertical direction that cause the Kell effect. (In passing, we should mention that for very rapidly changing images comparable

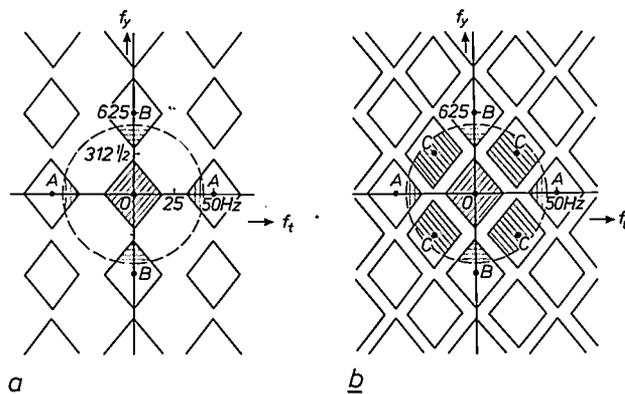


Fig. 10. a) Cross-section at $f_x = 0$ of the three-dimensional spectrum of a television signal that has been built up from 625 lines per picture and 50 pictures per second (without interlacing). *b)* As in *(a)* but now with 2:1 interlacing; there are now 50 fields and 25 complete pictures per second. The spectrum of the original scene is concentrated around the origin O ; the spectral repetitions at A cause large-area flicker, the ones at B are responsible for the visibility of the (static) line structure and those at C lead to line flicker and an apparent movement of the line pattern. The dashed circles give a rough indication of the limits of perception of the human eye under normal viewing conditions.

effects arise in the f_t -direction; however, these are much less of a nuisance.)

By starting with a camera with $2f_v$ lines instead of the usual number f_v , but still with 2:1 interlacing and using a number of processing operations, we can obtain a very interesting television signal with f_v lines and 2:1 interlacing. This signal has *no* vertical aliasing and can be transmitted and displayed with standard equipment for f_v lines. Yet the quality of the displayed picture is improved, because some of the Kell effect has been eliminated. However, the same signal can also be used, after a number of additional processing operations in the receiver, for display with $2f_v$ lines, and then the Kell effect can be completely eliminated.

The successive stages in this process are shown in the block diagram of *fig. 11* and the spectra of *fig. 12* for a system starting with a camera signal with 1250 lines, 2:1 interlacing and a field frequency of 50 Hz (designated as 1250/2:1/50 for brevity). With the aid of a number of memories a non-interlaced (i.e. 'sequential') signal with 1250 lines is first generated from this in a 'scan converter'. In this conversion

[7] R. D. Kell, A. V. Bedford and M. A. Trainer, An experimental television system, *Proc. I.R.E.* 22, 1246-1265, 1934.

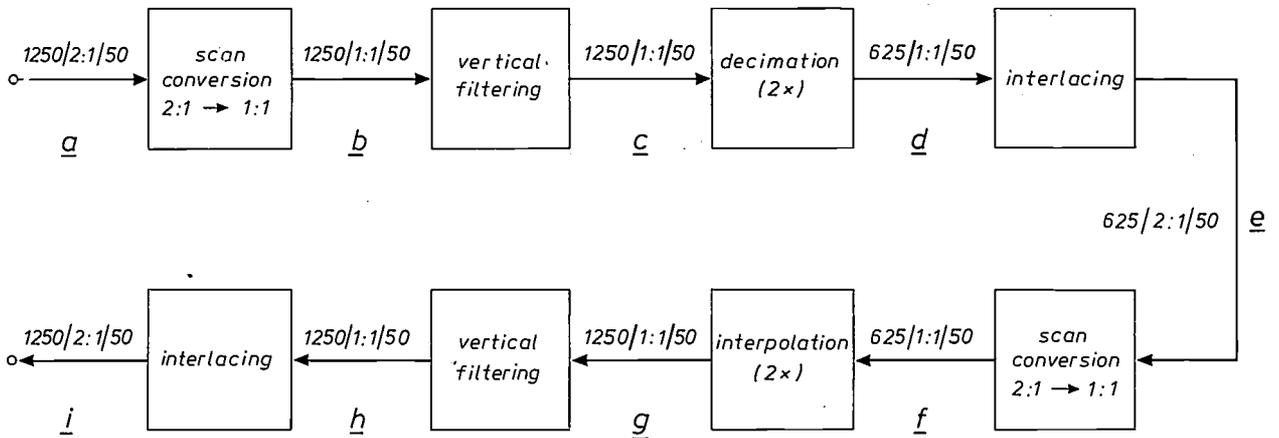


Fig. 11. By starting from a television signal of type 1250/2:1/50, followed by scan conversion (changing from an interlaced picture to a sequential one), vertical filtering, decimation (halving the number of lines) and then interlacing again a signal of type 625/2:1/50 can be obtained in which there is no aliasing in the vertical direction. This signal can be displayed on a standard television set. However, by carrying out a number of operations complementary to the ones we have just mentioned, it is possible to obtain another signal of type 1250/2:1/50 that is completely free from Kell effects and therefore represents a considerable improvement in quality. The spectra of the signals *a*, *b*, ... are shown in fig. 12.

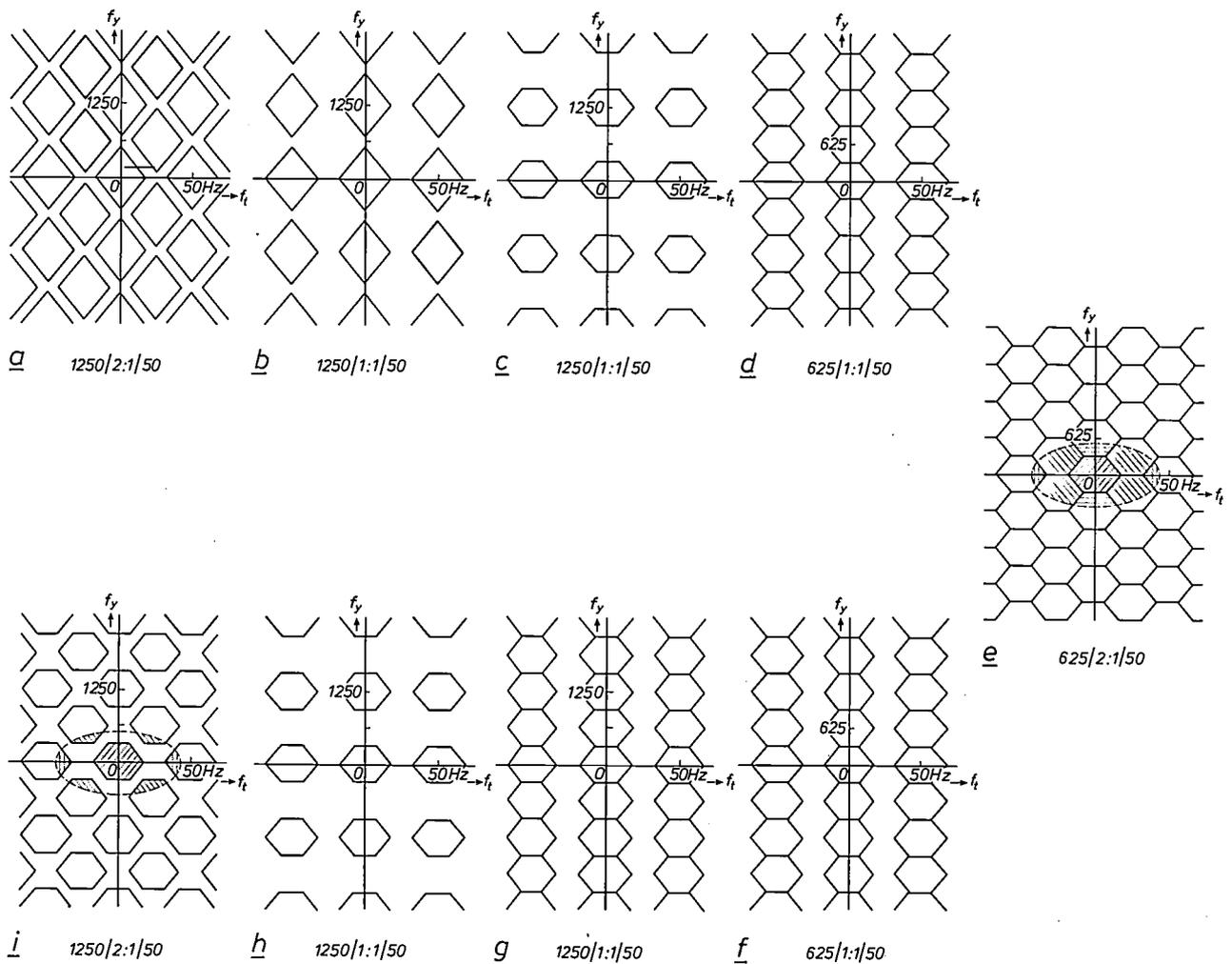


Fig. 12. Two-dimensional spectra in the (f_t, f_y) -plane at various points in the block diagram of fig. 11. The letters *a*, *b*, ... correspond in these two figures. (In this figure the limits of human visual perception appear as an ellipse instead of a circle as a result of the altered scale division in the vertical direction.)

something similar takes place to what we saw for one-dimensional spectra in fig. 7: by combining two discrete signals — representing odd fields and even fields respectively — that contain essentially the same information, certain spectral repetitions (at least for stationary images) can completely compensate one another. Next, the original spectrum is limited to frequencies below $312\frac{1}{2}$ c/ph by purely vertical filtering. After this the number of lines can be reduced to 625 by simply omitting every other line (625/1:1/50). From this sequential signal an interlaced signal can be obtained (625/2:1/50) that is compatible with conventional television signals. Because of the vertical filtering there is no longer any aliasing in this signal in the f_y -direction, but there are still vertical spectral repetitions within the perceptual range of the eye (dashed ellipse).

In the form now reached the signal can be transmitted or stored by all the conventional methods. For optimum display we carry out the following extra operations. First we go from an interlaced signal (625/2:1/50) to a sequential signal (625/1:1/50). Then, with the aid of a vertical 'interpolating' filter [8], we produce in two steps a sequential signal with twice the number of lines (1250/1:1/50) and from which half of the total of spectral repetitions are removed. Finally, we derive an interlaced signal again from this (1250/2:1/50). Because of the interlacing new spectral repetitions appear, but these contain no purely vertical frequencies. There are now no spectral repetitions or overlaps in the vertical direction within the perceptual range of the eye: the Kell effect has therefore been eliminated.

Improvement of the resolution in the horizontal direction

In the standard MAC/packet system the complete television signal after time compression (and hence before the FM processing) can have a bandwidth of about 8.4 MHz. We should now like to concentrate our attention on systems with a time-compression factor of 5/4 for the luminance signal; this means that before time compression a bandwidth of $4/5 \times 8.4$ MHz \approx 6.75 MHz is available. How can we obtain from this the best possible resolution in the horizontal direction? Let us consider a television signal of the type 625/2:1/50. To permit the time compression to be carried out this signal is sampled; a possibility for this is indicated in fig. 13a (crosses for the even fields and squares for the odd fields); if (as here) the pixels to be sampled are in a rectangular pattern in each field, this is called 'orthogonal sampling'. For stationary images we can combine the samples from dif-

ferent fields with one another in the receiver and hence halve the actual vertical distance between the samples.

We can now use the (f_x, f_y) -plane to express the maximum resolution in both horizontal and vertical direction: this is exactly equal to the maximum frequency range to which we have to limit the spectrum of the television signal if we wish to avoid aliasing as a result of spectral repetitions. (Strictly speaking, we are of course dealing with a three-dimensional spectrum [9] again; but now we are confining ourselves to the plane $f_t = 0$.) We consider a picture of height h and width w . At a line spacing (within one field) of Δy and a horizontal spacing of Δx between the samples we find that the maximum resolution for stationary images is given by a rectangle of height $2h/\Delta y$ and width $w/\Delta x$ around the origin (fig. 13b). In the display of moving images, we cannot simply combine samples from different fields with one another, since we then get movement blur. This means that without special measures the maximum resolution in the vertical direction for moving images is only half that for stationary images. (See also the section 'Moving images'.)

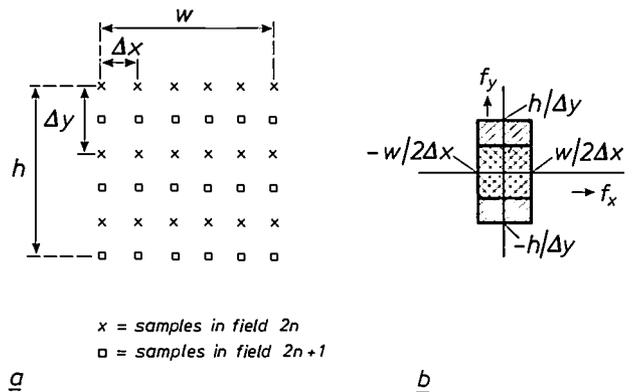


Fig. 13. a) Orthogonal sampling of a television picture. The crosses represent the samples taken in an even field and the squares represent the samples taken from the odd fields. The relative line spacing within one field is $\Delta y/h$ and the relative horizontal spacing is $\Delta x/w$, where h and w represent the height and width of the picture respectively. b) The maximum attainable resolution can be indicated in the (f_x, f_y) -plane. It corresponds to the maximum dimensions of the region in which we can prevent aliasing by performing the appropriate operations. Since image signals from successive fields can be combined for stationary images, the resolution in the vertical direction in this case (////) is twice as large as for moving images (||||).

There is of course a very direct relation between the horizontal sampling distance and the sampling rate $1/\tau$ (in Hz) used in taking the samples. This is

$$\frac{w}{\Delta x} \text{ [c/pw]} \triangleq \frac{1}{\tau} \text{ [Hz].}$$

We can therefore think of the f_x -axis as having a scale in hertz. This

can sometimes be extremely useful: for example a horizontal low-pass filter with a passband up to K/τ (in Hz) has — expressed in another way — a passband up to $Kw/\Delta x$ (in c/pw).

In some treatments the starting point is not the total picture height h and total picture width w , but only the active portion of the picture is considered (i.e. excluding the flyback times). We have already come across this for the vertical direction in the small print of p. 8. By analogy with the concept of active lines we can also refer to the active picture width w_{act} and relate the horizontal frequency to this. Such an approach is particularly useful if we want to refer to the number of pixels N_H in the horizontal direction, since this has a maximum value of

$$N_H = \frac{w_{act}}{\Delta x}$$

In general, however, N_H will have a smaller value, e.g. as a result of horizontal filtering.

A much more interesting situation arises when a sampling pattern like the number 5 on a dice is used in sampling each field; fig. 14a. This is called 'quincunx

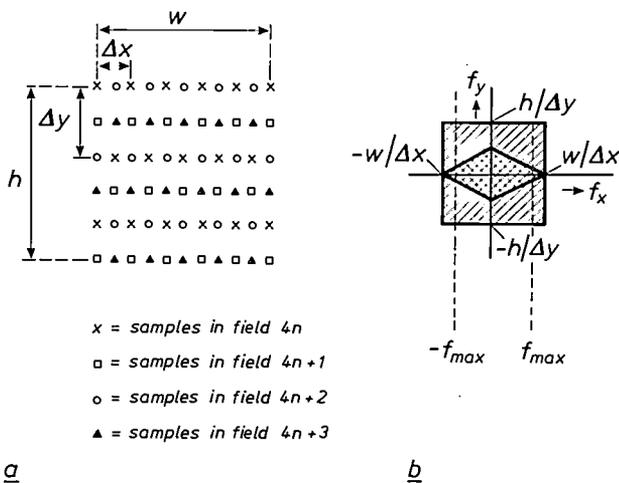


Fig. 14. a) In quincunx sampling the samples in two successive lines in the same field are displaced by half a sampling period with respect to one another. The sampling rate is little different from that in the previous figure (half the line frequency, to be precise), but now four fields must elapse, not two, before the same pixels are sampled again. b) The maximum resolution is quite different, however: for stationary images it is increased by a factor of two in the horizontal direction (////) and for moving images it has a diamond shape (\\\\). For moving images with $f_y = 0$ the maximum horizontal resolution is also increased by a factor of two. If the television picture is limited to a bandwidth of $f_{max} < w/\Delta x$ before sampling, this gives a loss in horizontal resolution as indicated by the dashed lines.

sampling'. Now four field periods elapse before the sampling cycles repeat, i.e. before exactly the same pixels are sampled again. The associated maximum resolution for stationary and moving images [10] is shown in fig. 14b. In the horizontal direction the resolution for stationary images has been increased by a

factor of two and also for moving images containing no vertical frequencies. We have to remember, however, that to avoid all aliasing at moving parts of the picture we should use two-dimensional filtering, as indicated by the cross-hatched ('diamond shaped') area. By using quincunx sampling at a sampling rate of $1/\tau$ we can thus transmit a maximum horizontal frequency corresponding to $1/\tau$. Since (seemingly) we are not satisfying the sampling theorem here, but are really creating a situation like that of fig. 7c and d, we call this *sub-Nyquist sampling* or *subsampling*.

With the bandwidth of 6.75 MHz available to us (before time compression) in the MAC system we can accept a maximum sampling rate of 13.5 MHz. We further ensure that the television signal is limited to 10 MHz before sampling [11]. After sampling and bandwidth-limiting to 6.75 MHz we then have a one-dimensional spectrum as shown in fig. 15a; in the range between about 3.5 and 6.75 MHz a special kind of

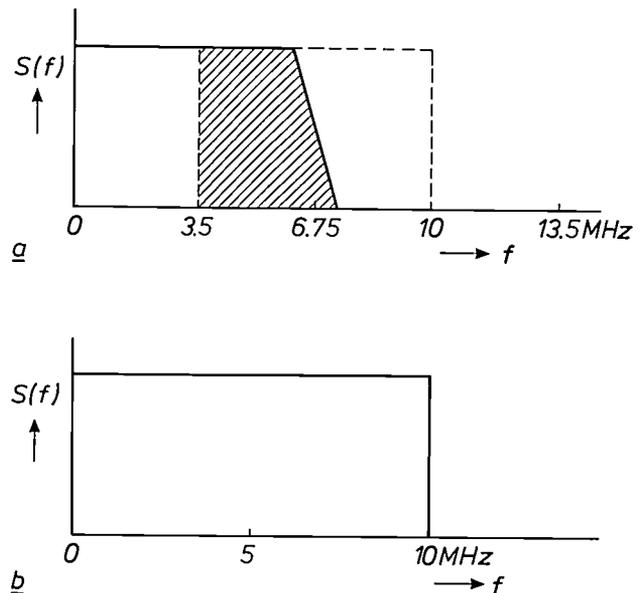


Fig. 15. a) If sub-Nyquist sampling ('subsampling') is used, a 'constructive' form of aliasing is obtained. Here this takes place in the frequency range between about 3.5 and 6.75 MHz. b) By applying the correct 'interpolating' filter operations we can reconstruct the original frequency components from this. The original frequency band from 0 to 10 MHz thus becomes available again without distortion.

[8] A. W. M. van den Eenden and N. A. M. Verhoeckx, Digital signal processing: theoretical background, pp. 110-144 in the special issue 'Digital Signal Processing I, background', Philips Tech. Rev. 42, 101-148, 1985.
 [9] G. J. Tonge, The sampling of television images, IBA report 112/81, Independent Broadcasting Authority, Winchester, Hants, England, 1981 (34 pp.).
 [10] The diamond-shaped figure that corresponds to moving images can be derived directly from the theory described in [9].
 [11] Because of this restriction the movement detection is simplified, since there is no aliasing in the frequency range below 3.5 MHz; see fig. 15a.

aliasing occurs. Because of the quincunx sampling and the two-dimensional filtering of moving parts of the image before sampling, we can remove this aliasing and at the same time reconstitute the frequency spectrum in the range between 6.75 and 10 MHz. At the receiving end we need the same two-dimensional filtering operation for this (again at moving parts of the image) as before sampling. This requires a sampling rate of at least $2 \times 1/\tau = 27$ MHz. We then obtain the one-dimensional spectrum of fig. 15*b* as the result ^[12]. In this way we can obtain the maximum horizontal resolution corresponding to fig. 14*b*. Because of the bandwidth limitation of the television signal to 10 MHz before sampling, as mentioned earlier, however, we remain within the strip $(-f_{\max}, f_{\max})$ in the horizontal direction. This means that in the horizontal direction we can resolve a maximum of about a thousand pixels per line.

So far in this section we have been dealing with the luminance signal of a television picture consisting of 625 lines. If we use the techniques of the previous section to display this as a signal of type 1250/2:1/50, we in fact double the total equivalent bandwidth to 20 MHz. This is a great improvement in comparison with the bandwidths of the luminance signal in the existing PAL and MAC systems, which are about 4 and 6 MHz respectively (for 625 lines per picture).

Moving images

In the foregoing we have seen how by applying the appropriate signal processing with a MAC/packet signal of type 625/2:1/50 we can transmit HDTV information that is suitable for display as a 1250/2:1/50 picture with an equivalent bandwidth of 20 MHz.

Strictly speaking, the maximum improvement that we can achieve in this way only applies to *stationary* images. This is because the information from different fields of stationary images can be combined without disadvantage ('inter-field processing'); this is one of the prerequisites in converting from an interlaced (2:1) picture to a sequential (1:1) picture and in reconstituting a subsampled signal. When we have to deal with moving images, however, we have to make sure that the inter-field processing does not cause movement blur. Since we subdivide the television image into pixels at both transmitting and receiving ends, we can in principle decide for each pixel whether there is any movement (this is done in a 'movement detector') and adapt the signal processing accordingly.

Various kinds of adaptation are possible, and at the present it is not yet clear which is the best (and most economic) solution or combination of solutions. We

should like to mention a few possibilities here in connection with scan conversion.

In this conversion we want to calculate missing picture lines. This can be done by combining information from a higher and a lower line from the same field (intra-field processing) or by combining a preceding line and a following line, seen in the time direction (inter-field processing); see fig. 16. The first solution works better for moving images, the second one gives the best result for stationary images. We can base our choice between the two on movement detection. An alternative possibility is to take the same combination of all four adjacent lines in both situations (movement

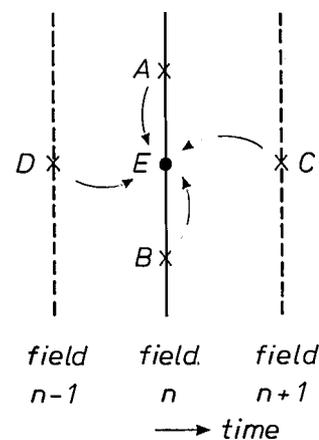


Fig. 16. Calculation of the missing line *E* in scan conversion (the picture lines are perpendicular to the plane of the diagram). If *E* is calculated from *A* and *B* we have intra-field processing. If *E* is calculated from *C* and *D* we have inter-field processing. If all four of *A*, *B*, *C* and *D* are used, we are performing spatio-temporal processing.

or not). In that case we perform a permanent vertical-temporal filtering operation of a lowpass nature in the f_y - and f_t -directions. An advantage of this is that we never have to switch from one kind of processing to the other.

Yet another method of minimizing movement blur is to introduce a movement vector; for example, if a television camera performs a movement, it effects all pixels in the same way. In principle it is possible to pass this information to the receiver in a very efficient way (especially in the MAC/packet system) where it can be accounted for in the display.

The colour signals in HD-MAC

So far we have confined ourselves to a discussion of various processing operations that are carried out on the luminance signal *Y* in the HD-MAC system. To display a colour picture, however, we also require two

colour-difference signals U and V . In general terms these signals are comparable in their characteristics with the signal Y . There is however one important difference: colour information can be displayed in much less detail than the associated luminance information. This feature is put to use in various ways, as we shall explain later. Otherwise the processing of the signals U and V in HD-MAC is much the same as for Y , to obtain the best possible resolution in the vertical and horizontal directions for the colour as well. We shall not go into further detail about this here.

It is however important that in the HD-MAC system on which we shall now mainly concentrate our attention the bandwidth of the signals U and V in the horizontal direction is limited to a quarter of the bandwidth of the signal Y . In the time compression the signals U and V can be compressed four times more than the signal Y (i.e. 5 times instead of 5/4 times). Because of the smaller bandwidth the signals U and V each represent 1/4 of the number of pixels of the signal Y in the horizontal direction.

In addition, in the standard MAC/packet system the Y signal is transmitted during each line period, but only one of the colour-difference signal is transmitted. The missing colour-difference signal is then reconstituted as well as possible by combining colour information from adjacent lines. This gives a smaller resolution for the colour than for the luminance in the

vertical direction. In HD-MAC as well, because of the compatibility, only one colour-difference signal can be transmitted during each line period. To obtain improved vertical resolution as compared with standard MAC, nevertheless, for display as an HDTV picture of type 1250/2:1/50, a number of signal manipulations are required at the transmitting end in the conversion from the 1250/2:1/50 camera signal to the 625/2:1/50 transmitter signal and in the reverse conversion at the receiving end; this is shown in fig. 17. The picture lines of the camera signal are designated from top to bottom by the numbers 1, 2, 3, The continuous lines refer to the odd fields, the dashed lines to the even fields. The colour-difference signals with one asterisk are not directly available at the transmitting end, but are calculated [13] there from signals that are available. For example:

$$U_5^* = \frac{U_4 + U_6}{2} \quad \text{and} \quad V_9^* = \frac{V_8 + V_{10}}{2}.$$

In a similar way the signals with two asterisks are calculated at the receiving end from the received HD-MAC signals. For example:

$$U_2^{**} = \frac{3U_1 + U_5^*}{4}, \quad V_2^{**} = \frac{3V_1^* + V_5}{4},$$

$$U_3^{**} = \frac{U_1 + U_5^*}{2}, \quad V_3^{**} = \frac{V_1^* + V_5}{2},$$

$$U_4^{**} = \frac{U_1 + 3U_5^*}{4}, \quad V_4^{**} = \frac{V_1^* + 3V_5}{4}.$$

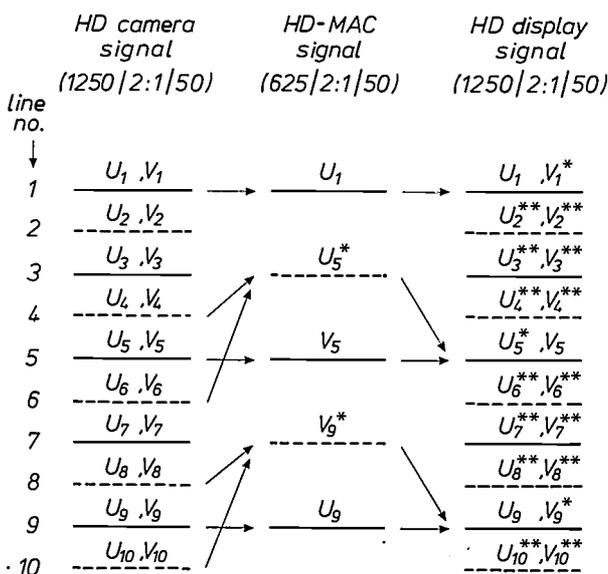


Fig. 17. Diagram illustrating the conversion of the colour-difference signals U and V of an HD camera signal of type 1250/2:1/50 into an HD-MAC signal of type 625/2:1/50 and of the conversion of this signal into an HD display signal. The colour-difference signals with one asterisk are calculated in the first conversion, the colour-difference signals with two asterisks are calculated in the second conversion. The lines from the odd fields are continuous, those from the even fields are dashed.

Since in signal processing in the HD-MAC system it is



Fig. 18. Quincunx sampling is also used for each of the two colour-difference signals in HD-MAC. Because only one of the two difference signals is transmitted in each line, the sampling pattern shown here is obtained. (The symbols have the same significance as in fig. 14 and fig. 17.)

[12] It is very important here that the overall frequency characteristic $H(f)$ of the MAC channel, including any transmitter or receiver filters has 'Nyquist shaping' around $f_0 = 6.75$ MHz, i.e. $H(f_0 + f) = 1 - H(f_0 - f)$.

[13] The calculations given here as an example ('linear interpolation') are about the simplest that can be devised for this purpose; those used in practice are more complicated.

required that the *U* and the *V* signals should both be sampled in the quincunx pattern described earlier, a cycle of sampling patterns like that shown in *fig. 18* occurs during the transmission; the same symbols are used to designate the different fields as in *fig. 14*.

In HD display, as a result of all the measures described above, a maximum vertical frequency of $156\frac{1}{4}$ c/ph is possible for both colour-information sig-

The final result can be displayed on both a special HDTV set and a standard television set. The HDTV set is suitable for display of the same type of signal as that originally provided by the HDTV camera. The sampling rate f_s used is indicated at various places in the block diagram of *fig. 19*.

The most noticeable difference from the foregoing is the position where the horizontal processing takes

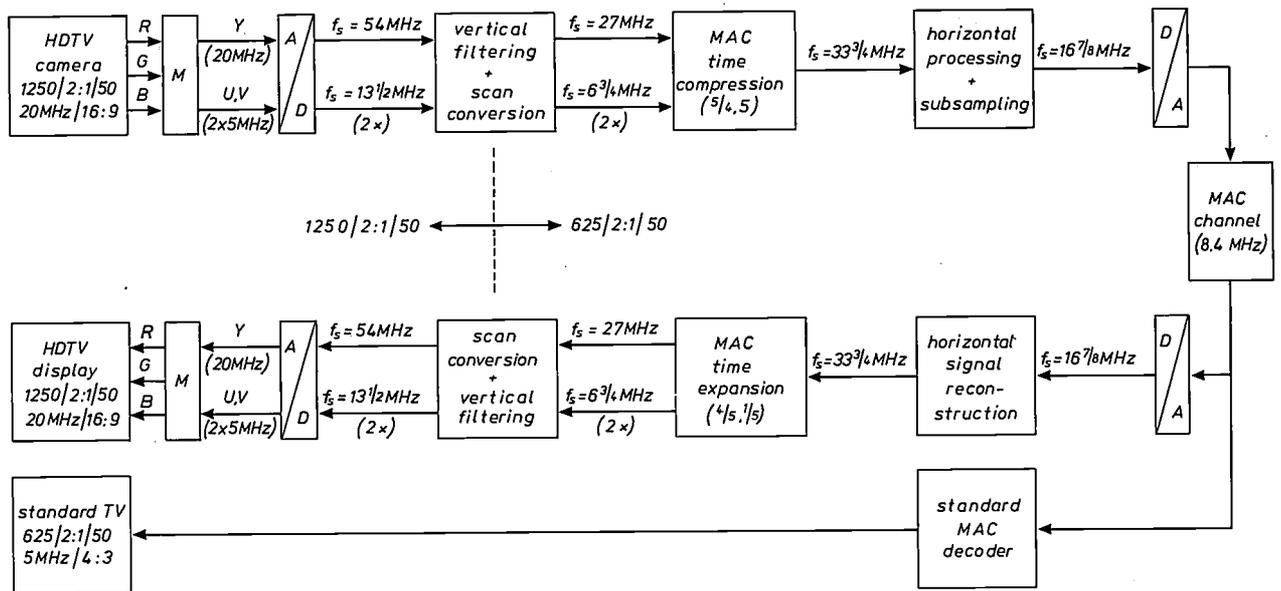


Fig. 19. A complete HD-MAC system corresponding to the simplified diagram of *fig. 5b*. Most of the signal-processing operations we have discussed can be found in this diagram. It is interesting to note that all the horizontal operations take place between time compression and time expansion. Since most of the operations are performed digitally, the system contains A/D and D/A converters. The sampling rate used f_s is indicated at various points. *M* is a matrix circuit for the conversion of the three original colour signals *R*, *G* and *B* into the signals *Y*, *U* and *V* (or vice versa).

nals with stationary images. This is half the maximum vertical frequency for the luminance signal, but twice the value permitted by the standard MAC/packet system (see also p. 211).

A complete HD-MAC television system

To demonstrate the practicability of HDTV based on the MAC/packet system and the compatibility between HD-MAC and standard MAC we have constructed a complete television system in which most of the signal processing described earlier is applied; see *fig. 19*. The starting point is an HDTV camera that provides a signal of type 1250/2:1/50 with an aspect ratio of 16:9 and a *Y* bandwidth of 20 MHz. After the appropriate processing the signal is transmitted on a MAC channel with an effective bandwidth of 8.4 MHz.

place: after time compression at the transmitting end, and before time expansion at the receiving end. The reason for this is that the compression gives the signals *Y*, *U* and *V* the same bandwidth, separates them in time and enables them to be processed in a similar way.

At the receiving end the signals *U* and *V*, after time expansion, are further processed in a number of special circuits (further detail will be omitted here) to give extra noise reduction with stationary images and to improve the colour transients [4]. This gives an additional improvement in the quality of the displayed picture.

We should mention that all of the television signals referred to in *fig. 19* are interlaced. The non-interlaced signals that we encountered in the theoretical treatment earlier were only introduced to simplify the des-

cription. In practical signal processing there is no explicit need for non-interlaced signals.

Compatibility of HD-MAC with standard MAC

As we can see from the block diagram of fig. 19, an HD-MAC signal can be presented directly to a standard MAC receiver on reception. Because of the vertical filtering that has been performed on the HD-MAC signal at the transmitting end, the picture quality in the vertical direction will be better than for reception of an 'ordinary' MAC signal, for there can now be no aliasing in the vertical direction. On the other hand, there are a few slight degradations in quality. First, there is an increase in line flicker. This is caused by the fact that in HD-MAC there is less attenuation at the high vertical frequencies (up to $312\frac{1}{2}$ c/ph). Because of this, however, the corresponding frequency components in the periodic repetitions that are responsible for line flicker (see the points C in fig. 10b) are stronger. Also, because of the subsampling of the HD-MAC signal in the horizontal direction, aliasing is introduced into the luminance signal Y

play. This can be explained with the aid of fig. 17. If the HD camera signal shown there were to be converted in accordance with the specifications that apply for the *standard* MAC system into a signal of type 625/2:1/50, then we would require the following sequence of successive lines:

$$U_1, U_3^*, V_5, V_7^*, U_9, U_{11}^*, V_{13}, V_{15}^*, \dots$$

instead of the present HD-MAC sequence

$$U_1, U_5^*, V_5, V_9^*, U_9, U_{13}^*, V_{13}, V_{17}^*, \dots$$

Half of the colour information is therefore slightly displaced (by two lines) in the vertical direction. The degradation this introduces when an HD-MAC signal is displayed on a standard MAC receiver is very slight, however.

Quality comparisons

To permit a general comparison to be made between various television signal standards, we have collected the leading parameters of a number of existing and proposed television system in *Table I*. Besides the

Table I. A comparison of the leading parameters of a number of television systems.

TV system	Time-compression factor		Aspect ratio	Equivalent bandwidth (MHz)		Number of pixels			
	Y	U,V		Y	U,V	Horizontal		Vertical	
						Y	U,V	Y	U,V
PAL	1	1	4:3	3.7	1.0	370	100	290	230
MAC	3/2	3	4:3 and 16:9	5.6	2.8	560	280	290	144
	5/4	5	4:3 and 16:9	6.7	1.6	670	160	290	144
HD-MAC	3/2	3	4:3 and 16:9	16	8	800	400	520	260
	5/4	5	4:3 and 16:9	20	5 [*]	1000	250	520	260
MUSE	c ₁ [**]	c ₂ [**]	16:9	22	7	1000	330	520	260

[*] Because of the nonlinear operations carried out in the 'Colour-Transient Improvement' this value is in fact larger.

[**] $c_2/c_1 = 4$.

in the frequency range above 3.5 MHz (see fig. 15). Since this distortion is not compensated in the standard MAC receiver, it introduces — as close observation reveals — a barely noticeable, rapidly moving dot pattern comparable with the 'dot crawl' found in the existing NTSC television system, but it is less annoying, because its structure is almost twice as fine. Moreover, it is not present in unpatterned areas but only in 'textured' areas.

The processing of the colour-difference signals U and V in the HD-MAC system also introduces what is really a very small error in a standard MAC dis-

time-compression factors and the aspect ratio, we have shown the equivalent bandwidths of the signals Y, U and V (for stationary images). We have also shown the number of pixels for the same signals in the horizontal and vertical directions (N_H and N_V). These numbers relate to the active part of each line period and the active fraction of all the lines (see also the small print on p. 204 and p. 207). The Table mentions the MUSE signal [14], which we have not mentioned

[14] Y. Ninomiya, Y. Ohtsuka and Y. Izumi, A single channel HDTV broadcast system — the MUSE, NHK Laboratories Note No. 304, 1984 (12 pp.).

previously, of the type 1125/2:1/60. This has been zealously promulgated by the Japanese broadcasting organization NHK. At first sight the figures for MUSE and HD-MAC are comparable. However, the fact that the MUSE signal is compatible with no other signal (and so is *not* evolutionary in nature), swings the balance irrefutably in favour of HD-MAC.

The introduction of the HD-MAC system at the transmitting end makes it possible to change over gradually to better receivers at the receiving end, without existing equipment becoming unusable. Four gradations in quality can be distinguished here:

- reception with a standard MAC receiver;
- improvement of reception by the addition of (for example) a horizontal comb filter to eliminate the dot crawl resulting from subsampling;
- reception with a fairly simple HDTV receiver, in which use is made of movement detection and switching between inter-field signal processing and intra-field signal processing;
- reception with an advanced HDTV receiver, in which for example a movement vector is used and information about this is transmitted via the MAC/packet multiplex.

After all that we have said it should be abundantly clear that the HD-MAC system opens the way to a compatible and therefore smooth and attractive evolution towards HDTV.

Summary. As the dimensions and the light intensity of television receivers increase, the limitations of the existing television systems (NTSC, PAL and SECAM) become more obvious, and the need grows for a new system, known as high-definition television (HDTV). This must offer an increase in the resolution, a reduction in the interaction between colour and luminance information, an increase in the aspect ratio and stereophonic sound with 'hi-fi' quality. It is ultimately of vital importance here that a compatible system is chosen, so that existing equipment can still be used and the viewer can experience a gradual ('evolutionary') progress of television quality through the gradual acquisition of new equipment. The MAC system (MAC means Multiplexed Analog Components) recently standardized for use in satellites and cable systems will permit such an evolution. This article describes — by way of example — how various advanced signal-processing operations at the transmitting end will provide an 'HD-MAC' signal. This signal is suitable both for reception with standard MAC equipment, and for reception with a special HD-MAC set that can display a picture of HDTV quality. The signal processing required is described with the aid of one-, two- and three-dimensional spectra. Separate attention is paid to the transmission of moving images and colour information. The article concludes with the description of a complete HD-MAC television system, a discussion of the compatibility of HD-MAC and standard MAC and a quality comparison for various television systems.

1937

THEN AND NOW

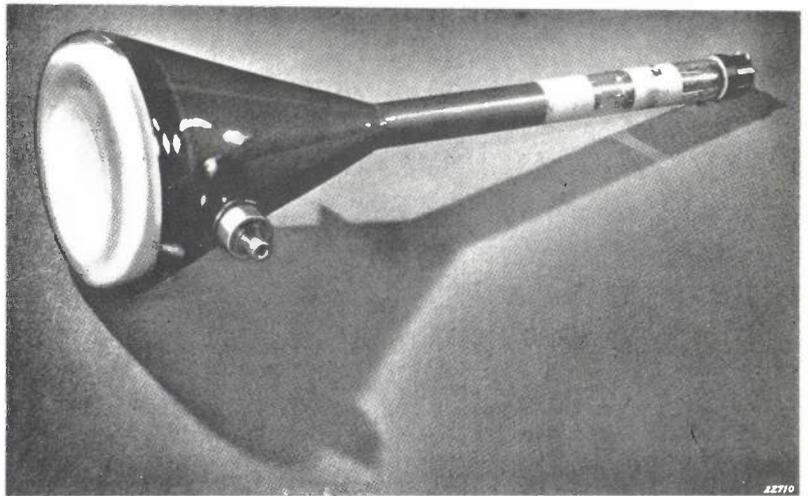
1987

Television projection tubes



Even before World War II, when television was still in its infancy, Philips Research Laboratories were working on projection television to give larger pictures (e.g. 40cm × 50cm) than the screen of the largest cathode-ray tube. To obtain the sharpest pictures the electron beam in the projection tube was focused magnetically instead of electrostatically; also, so that existing optical lenses could be used, the front of the projection tube was curved to match the image-field curvature of the optical lens (black-and-white photo) [*].

In 1987 the technical requirements demanded of television projection tubes are much more onerous. Colour pictures are produced by three separate tubes — one for each primary colour — which have to work in close cooperation; the diagonal dimension of the projected picture should be about 1 m or even larger; to keep the complete receiver compact the projection tube should have a rectangular front face and a large internal de-



flexion angle. For high-definition television (HDTV) — soon to be with us — the requirements will be even more onerous.

Recent research has led to the projection tube as shown on the colour photo. The electron beam is again focused electrostatically through the application of a new electron-optical principle. This includes the use of a special type of pre-focus lens in the electron gun, to give the very high resolution required. The screen is almost rectangular, about 8.5cm × 11cm, and the total length is 25 cm.

[*] from Philips Technical Review, August 1937.

Selecting quartz for resonators

J. C. Brice and W. Koelewijn

A chain is only as strong as its weakest link. In an electronic device or system that 'weakest link' is the worst component, of course, and we therefore have to give due attention to the manufacture of even the humblest component. Our care and concern must go right back to the selection of the raw material. This will often determine the performance of the component — so our choice must be soundly based. The article below discusses the criteria for selecting quartz for resonators.

Devices which require accurate timing or well-defined frequency features contain quartz resonators to meet these requirements. Without these resonators the use of colour television, mobile radio, telephone systems etc. would not be as widespread as it is now. In order to achieve satisfactory yields of devices meeting high technical specifications, it is necessary to select the best available raw material. The reasons for selecting a particular type of quartz and the method by which selection can take place are described in this article. Three simple tests appear to be sufficient for an accurate selection. First we shall briefly go into the physical background of the quartz resonators.

Materials with the crystal symmetry of quartz (threefold symmetry axis) or a lower degree of symmetry show piezoelectric behaviour. As a result of a mechanical distortion of a quartz crystal in a specific direction, electrons and nuclei are moved by different amounts. Thus a positive charge appears on one face of the crystal and a negative charge on the opposite face. The inverse effect — an applied voltage resulting in a mechanical distortion — also occurs. An alternating voltage produces cyclic deformation. At resonant frequencies this exchange of electrical and mechanical energy can be significant. A quartz resonator is a plate

of piezoelectric quartz cut in appropriate directions so that its natural resonance occurs at a particular well-defined frequency. Piezoelectricity in quartz only occurs if the mechanical distortion is in a direction that is not parallel to the threefold symmetry axis of the crystal. This implies that plates that are to be used as resonators must be cut according to this condition. Research into the possible orientations of the crystals revealed that some specific orientations resulted in resonators with frequencies that are independent of temperature^[1], a desirable feature in view of the widespread use of these devices. Specifications usually call for errors in the cutting direction of less than 0.01°.

Resonators are characterized by the electrical quality factor Q_e , the fundamental frequency of the mechanical distortion in response to an applied alternating voltage divided by the half-width of this response (see *fig. 1*). The great advantage of the quartz resonator is that it is nearly lossless so that it can have a very large electrical quality factor Q_e . This means that the resonance is very sharp. Typical coil-and-capacitor circuits have $Q_e \leq 200$, while quartz crystals can have a Q_e up to several million but typically in the range 10000 to 100000. The resonant frequencies f of a quartz plate are determined by its thickness: $f = nc/t$, where n is an odd integer, c is the velocity of sound in quartz and t is the thickness of the plate. Since maximum frequency deviations of no more than 1 ppm are

Dr J. C. Brice has recently retired from Philips Research Laboratories, Redhill, Surrey, England; Ing. W. Koelewijn is with the Elcoma Division, Philips NPB, Doetinchem, the Netherlands.

commonly requested, it must be possible to obtain crystals with well-defined thicknesses. The actual shape of the plate and the strain in the material also have an influence on the resonant frequencies.

Plates with the correct orientation and thickness are made in the following manner. Quartz is obtained from suppliers in batches consisting of a number of bars grown in one process. The crystal bar is oriented using an X-ray technique and then sawn in parallel-sided plates which are lapped to a thickness slightly greater than required. By way of an etching process damaged material is removed from the surface and the required thickness is attained. The electrodes necessary for applying the voltage are added by evaporating electrode material on the faces of the plate. These electrodes have an influence on the resonant frequency eventually reached. During the evaporation of the electrode material on the crystal the resonant frequency is monitored. Evaporation is stopped when the desired frequency is attained. Fig. 2 shows a diagram of a plate and the electrodes together with the shear motion as a result of the applied voltage.

The features of the crystal that, apart from the shape and the thickness, determine its frequency and its electrical quality factor Q_e are related to the crystal structure. It is therefore reasonable to expect that purer material gives better devices. Until 1955 only natural quartz was used, giving a yield of 10 to 30%. Since that date synthetic quartz has been available giving considerably higher yields. This is mainly due to the fact that synthetic quartz, as a result of its controlled growth, is easier to cut in the correct orientations, thus giving less waste material.

Our approach to the problem of establishing selection criteria has been to a large extent statistical. With the help of the results of experiments done by many colleagues we have assessed the performance of devices made from quartz with known characteristics, and so identified the material parameters (purity and perfection) which affect device performance and the interrelations between these parameters in the quartz available [2]. Since suppliers use different growth processes and different sources of raw material, we also considered the data for each supplier separately besides the overall trend. The necessity for this approach will become evident later (see fig. 9). Obviously data for a particular supplier show less spread than an overall distribution. However, overall trends are useful: they tell us what happens if we choose our supplier at random and change our choice from time to time.

Synthetic quartz is almost universally specified by an infrared quality factor Q_{IR} . This quality factor is inversely related to the extinction coefficient α^* due to hydrogen absorption. The latter is determined by

measuring the transmission in a quartz sample with parallel polished surfaces at wavelengths where hydrogen absorbs energy from the infrared beam. We really determine the total extinction coefficient in this way

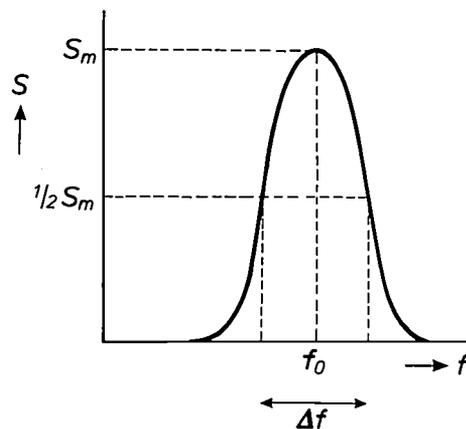


Fig. 1. Plot of the mechanical distortion S of a quartz plate in response to an applied alternating voltage of frequency f . The fundamental frequency of this plate is f_0 and Δf is the width of the response at half the maximum value of the response. These quantities determine the electrical quality factor $Q_e (= f_0/\Delta f)$.

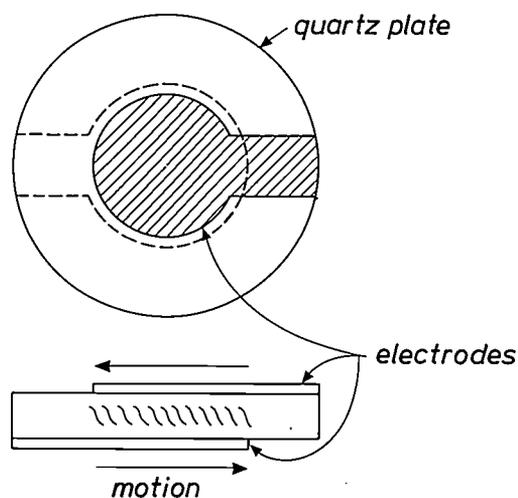


Fig. 2. A schematic representation of a simple device structure. The shaded area indicates the upper electrode; the lower electrode is indicated by a dashed line. In the lower part of the figure, the shear motion as a result of an applied voltage on the electrodes is indicated.

[1] J. C. Brice and W. S. Metcalf, Quartz-crystal resonators using an unconventional cut, Philips Tech. Rev. 40, 1-11, 1982.

[2] We received much help from our suppliers who provided us with many samples of material which is not normally commercially available. For some of our experiments, we needed particularly good and particularly imperfect samples. Data from these non-representative samples are, of course, excluded from statistics which show what is usually available. In total we investigated samples from 15 suppliers.

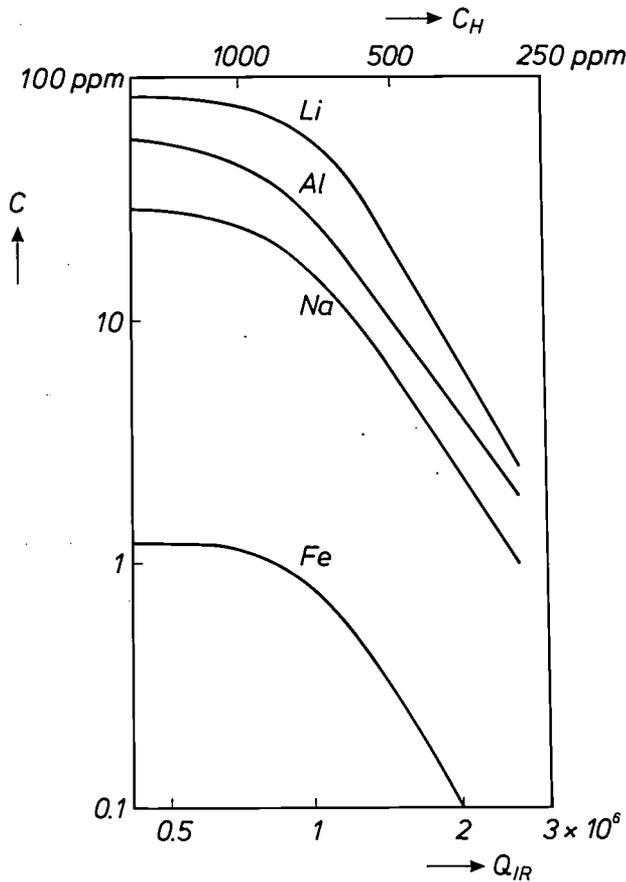


Fig. 3. Overall median impurity concentrations C as a function of Q_{IR} . The data are for the z -zone of crystals. Material grown on other faces is much less pure. For any supplier the distribution of $\log C$ at a given Q_{IR} -value is approximately Gaussian.

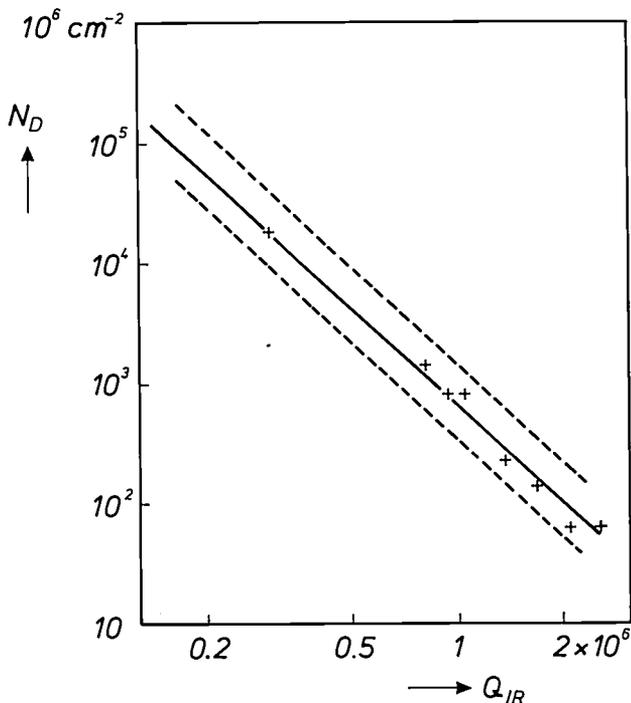


Fig. 4. Overall median dislocation density N_D as a function of the infrared quality factor Q_{IR} (continuous line and crosses). For any supplier the logarithm of the dislocation density has a Gaussian distribution. Results for the two most extreme suppliers are indicated by dashed lines.

but by doing reference measurements at nearby wavelengths the contribution due to hydrogen absorption can be established separately [9]. The parameters we have investigated (purity, dislocation density, strain, homogeneity) will be shown as a function of Q_{IR} .

Fig. 3 is a plot of the median [4] purity of the samples as a function of Q_{IR} (lower axis) or the equivalent hydrogen content (upper axis). Data from all the suppliers in the investigation have been used. High Q_{IR} corresponds to fairly pure material. The figure shows the major impurities (lithium, aluminium and sodium), which together with iron, are the only impurities known to be harmful. Other impurities show the same form of variation. Different suppliers have distributions of purity which differ from these data by factors of two or three, and different batches from one supplier at different Q_{IR} -values also show similar differences. However, the worst results differ by less than a factor of five from the median. By experience we learnt that impurities at levels less than 10 atomic ppm usually have negligible effects on devices. At higher levels, aluminium together with sodium can adversely affect yields of very-high-frequency devices: they change the etching characteristics so that very thin plates become porous. Lithium in high concentrations affects device stability [1] and iron decreases the radiation resistance of devices used in high-radiation environments.

Fig. 4 shows the variation of overall median dislocation densities as a function of Q_{IR} . Different suppliers' results can differ appreciably from this average behaviour. Data from all the suppliers in the investigation have been used. High dislocation concentrations are disastrous in three ways. First, high densities of dislocations reduce the yield of high-frequency devices [6]. Second, they make the material brittle, as discussed below, and third they reduce the performance of low-frequency devices. (Dislocations affect device performance by scattering the acoustic waves. Scattering is significant when the wavelength exceeds the distance between dislocations. Low-frequency-devices are therefore particularly affected.) This is shown in fig. 5, where the median equivalent series resistance of devices is given as a function of the dislocation density.

The correlation between dislocation density and Q_{IR} is not easily explained. What is found is that the density varies with the hydrogen content to the power of 2.5. The factor of proportionality varies with the supplier. Dislocations involve breaking bonds in the lattice and broken bonds can be terminated by hydrogen, so we might expect a linear relation ($N_D \propto C_H$). Some dislocations come from the seed crystals but most originate at inclusions [8]. Thus apparently

either the number of inclusions rises rapidly or, more probably, one inclusion results in a number of dislocations. Of course several mechanisms may operate simultaneously and the relation may be more complex than $N_D \propto (C_H)^{2.5}$.

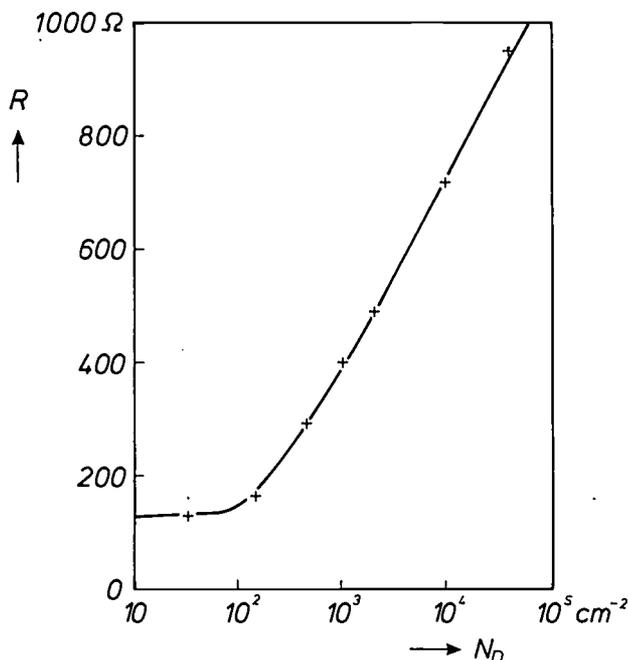


Fig. 5. A plot of the median equivalent series resistance R of devices as a function of dislocation density for 1.3 Mhz devices. An oscillator can be represented as an LC-circuit with a resistance in series. The quality factor Q_e of the device decreases with increasing equivalent resistance R . High dislocation densities imply a high R and thus a low Q_e .

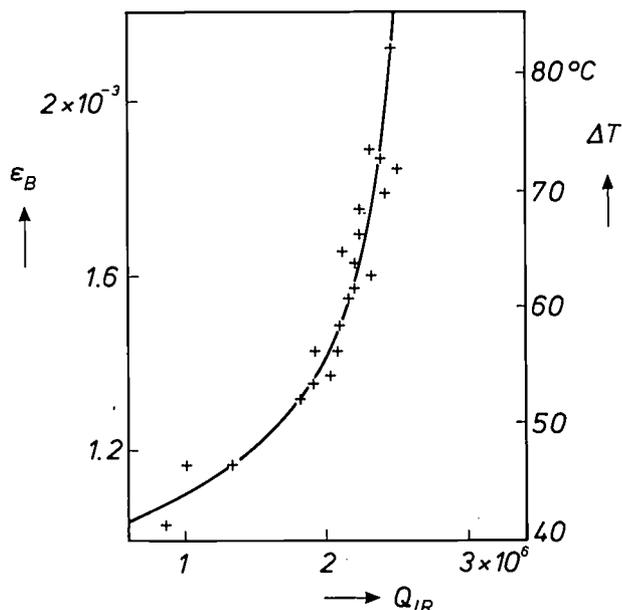


Fig. 6. Breaking strain ϵ_B (left-hand axis) and the temperature shock ΔT (right-hand axis) as a function of Q_{IR} . Samples were placed in a hot bath at a temperature T_B for 30 minutes and then transferred quickly to a cold bath at T_C . The samples of one batch were tested with increasing temperature difference $\Delta T = T_B - T_C$. The infrared quality factor Q_{IR} was measured for each sample.

In the course of fabrication devices are subjected to various stresses. It is therefore useful to know something about the strength of the material used. Strength in a brittle material like quartz is a fairly complex concept. The stress (force per unit area) required to break a plate depends on the surface finish. A plate with a perfect surface requires a fairly large stress to break it. In a perfect plate we have to nucleate a crack. In an imperfect plate we only have to extend an existing crack, which is much easier. Because we know the elastic constants of quartz, we can measure the breaking strain, rather than the breaking stress. We measured this strain by applying a thermal shock. We heated crystals in a water bath so that they were uniformly hot and plunged them into a cooler bath. If we cool them very rapidly by $\Delta T^{\circ}C$ the surface contracts an amount $\alpha \Delta T$, where α is the thermal expansion coefficient. The corners are under a rather larger tension (about three times the tension at the surface) so that cracks nucleate most easily there. The result obtained is shown in fig. 6. The correlation is clear but we were able to show that the dislocation density is also important [5]. The importance of the breaking strain becomes obvious when we consider figs 6 and 7. Fig 7 shows the yield of unbroken 0.5 mm thick plates in a production process.

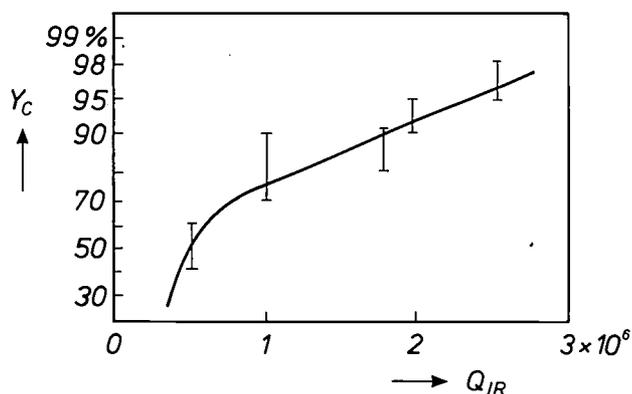


Fig. 7. Yield Y_C of unbroken 0.5 mm thick plates in a production process as a function of Q_{IR} . The data are for plates made with reciprocating slurry saws. When rotating wheels are used, yields are worse [6]. Note that the graph was drawn on probability paper so that the scale of Y_C varies in a nonlinear manner.

[3] For a detailed description of our procedure for measuring the infrared quality factor see: J. C. Brice, Crystals for quartz resonators, Rev. Mod. Phys. 57, 105-146, 1985.

[4] If the results found are written as a list in order of increasing numerical value, then if the list has an odd number of entries, the median is the one in the middle. If the list has an even number of entries, the median is the average of the middle two. The median value is exceeded by half the sample and is often more meaningful than the arithmetic mean. Consider a sample containing the values 0.5, 0.7, 0.9, 1.0, 1.3, 1.5, 10. The median is 1.0 but the mean is about 2.3.

[5] J. C. Brice, The specification of quartz for piezoelectric devices, Proc. 38th Ann. Frequency Control Symp., Philadelphia 1984, pp. 487-495.

We believe that during our processing procedures we expose the plates to stresses equivalent to the strain produced by a 50 °C shock. We therefore subjected every bar (46 000 in total) in a number of batches to a 50 °C thermal shock. Fig. 8 gives the results that we

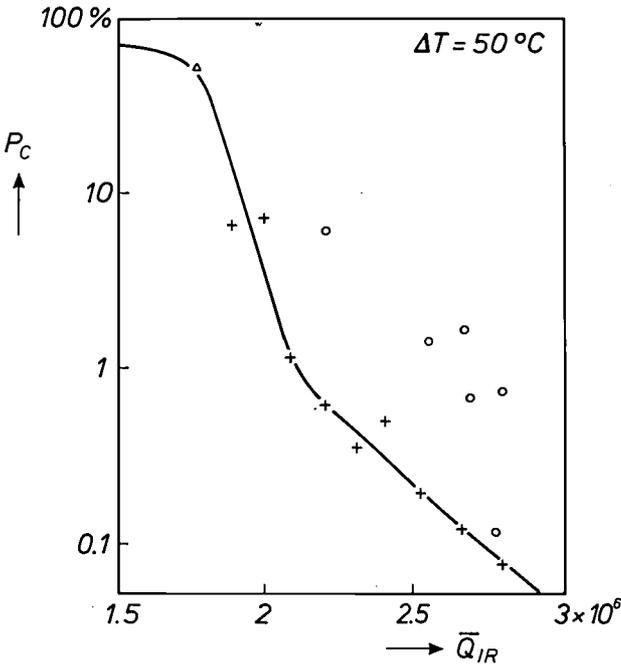


Fig. 8. The probability P_c that a bar will break with a 50 °C shock as a function of the mean Q_{IR} of the batch. Up to 30% of bars with $Q_{IR} < 1$ million survive a 50° shock. The crosses represent the results from normal batches and the circles represent the results for batches which contained noticeable numbers of deformed bars. The point marked Δ is taken from fig. 6. It is the value of Q_{IR} for which $\Delta T = 50^\circ$. A batch with a mean Q_{IR} of this value would have a failure rate of 50%.

found. In this figure the crosses represent the results from normal batches and the circles represent the results for batches which contained significant numbers of deformed bars. Clearly, these batches contain a significant excess of bars which break with a 50 °C shock (on average they contain 6 times as many bars which break).

Fig. 8 uses a batch mean Q_{IR} to describe the material. For this to be meaningful, we also need to know how the Q_{IR} of individual bars in a batch vary. These data are given for two suppliers in fig. 9.

The final piece of statistical evidence that we need to know is how homogeneous the crystals are. When we look at this we find that only in one direction is there an appreciable variation. Parallel to the z-axis of the crystals, the hydrogen content decreases as we move outward from the seed. Fig. 10 gives the data for two suppliers in the form of a graph of the gradient of α^* , the extinction coefficient, as a function of Q_{IR} . Data for most of the other suppliers lie between these two curves.

The results of our experiments show a qualitative relation between the infrared quality factor on the one hand and the other parameters of interest (impurity concentration, dislocation density, strain and homogeneity) on the other. The best yield of the production process will be reached if quartz with a high infrared quality factor ($Q_{IR} > 2$ million) is used. All the other requirements are then automatically satisfied.

In selecting our material we take the following line. First of all, every bar from each batch is inspected visually. We reject bars which contain clusters of inclu-

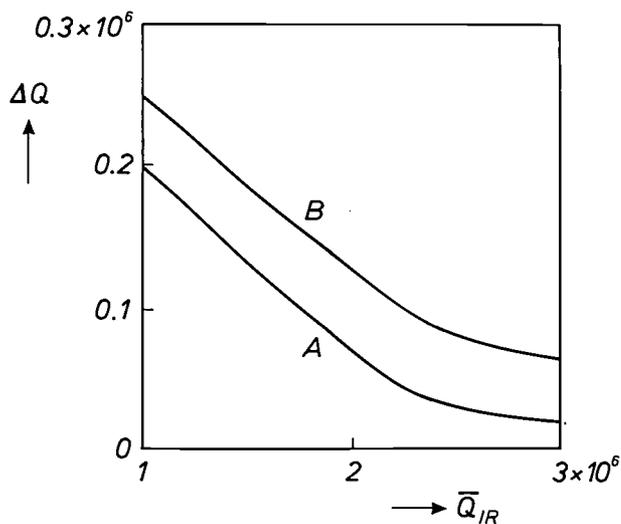


Fig. 9. The variation of Q_{IR} in batches as a function of batch mean Q_{IR} . In this figure ΔQ is the standard deviation of Q_{IR} in a batch. Curve A is for the supplier giving the greatest uniformity. Curve B is for the supplier referred to in fig. 8. Note that the data are actual variations of Q_{IR} . Because we measure Q_{IR} in three independent ways we know the experimental errors absolutely [6] and can therefore eliminate their effect before presenting the data.

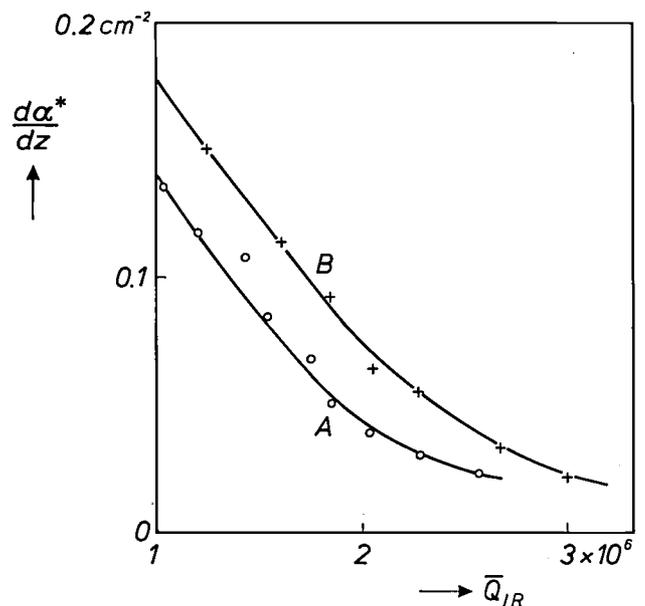


Fig. 10. The gradient of α^* (the extinction coefficient) in the z-direction as a function of the batch mean Q_{IR} . For A and B see caption of fig. 9.

sions large enough to be visible with the naked eye or which are grossly misshaped. Then we determine Q_{IR} for a small number of the remaining bars of each batch. The exact number depends on our knowledge of the suppliers' production process. If we believe that we know the standard deviation of Q_{IR} values in the batch (i.e. the appropriate curve on fig. 9) we can manage with a very small number (between 3 and 6 bars from a batch of 1000 or more). If we know nothing about the production process more bars should be examined, but only rarely do we need more than 12. Finally, every batch is subjected to a 50 °C thermal shock test on a 100% basis. For this test we have a tentative reject level of 10% but we have seen only one batch which had more than 10% breaking *and* a

Q_{IR} over 2 million. Batches passing the tests can be used to make devices to very tight technical specifications with high yields (95%).

Summary. To meet the high specifications usually required, quartz resonators must be made from quartz of a high quality. In this article the different aspects of this quality are discussed and the outcome of the investigations into the selection criteria are reported. Three simple tests appear to be sufficient for an accurate selection. Every quartz bar is inspected visually for gross defects and is exposed to a temperature shock of 50 °. For a small number of bars the infrared absorption at a wavelength where hydrogen absorbs infrared light is determined. Thus an infrared quality factor is assessed. Bars that do not show gross defects, do not break as a result of the temperature shock and come from a batch with a mean infrared quality factor of 2 million or more, can be used to make devices with the required technical specifications in production processes with high yields.

Looking at interfaces

F. Meijer

The address below was delivered by Dr F. Meijer at his inauguration on 7th March 1986 as Visiting Professor of Heterogeneous Catalysis at the University of Leiden. In his address Prof. Meijer does not confine himself to his own special field, but gives us, in a light-hearted style, a refreshingly relative view of the science of interfaces. The more objective sphere of the research he discusses includes measurement methods such as TEM and STM, which now approach or even reach the limit where atoms can be observed individually.

The concept of interface

Interface is a word that springs readily to the lips of technical people nowadays. The Oxford Dictionary gives the following two meanings of interface:

'surface forming common boundary between two regions',

'place or piece of equipment where interaction occurs between two systems'.

Both definitions are germane to my use of the word in the title of this inaugural address. The first definition brings out the essential aspect from which interface phenomena must be studied: the 'common boundary'. The second aspect, the interaction, the 'place where interaction occurs', follows immediately from this. The two aspects of the interface interact with one another.

A great deal goes on at interfaces.

In interface chemistry this includes all chemical reactions in heterogeneous systems. Heterogeneous means that there are two or more phases: a solid and a gas, or a solid and a liquid, or two different solids in contact with each other.

Corrosion, for example, such as the rusting of the metal of your car, is a chemical reaction at the interface of metal and humid air. Another phenomenon of practical importance is the adhesion between different materials, the adhesion of paint to a metal, or the adhesion of a metal film on a plastic substrate, as in the case of the Compact Disc. The word adhesion immediately brings out the concept of interaction.

There is also the very important field of heterogeneous catalysis, which concerns chemical reactions

at the interface of a gas or liquid with a solid, the catalyst. I shall return to this subject presently.

In physics and engineering the interface plays the leading role in the working or processing of a material; turning on a lathe, grinding, polishing, sputtering are all actions that take place at the interface. A sculptor is an artistic interface mechanic. He shapes the work at the interface between the tool and the stone. In a sense one might speak of the interaction between the sculptor and his material.

In the medical world the successful implantation of prostheses in the body, such as artificial hip joints and substitute arteries, depends to a very great extent on the interface between implant and body.

The word interface has also entered the language of electronics, defined as 'place or piece of equipment where interaction occurs between two systems', for example the interface between two computers or between a computer and a terminal. One of these two systems can be a human being, giving us the man/machine interface and the concept of 'user-friendliness'. The social sciences have also discovered the word interface and will undoubtedly add to its connotations.

There are many kinds of interface, and something usually happens at them — and that 'something' is worth examining more closely. 'Taking a closer look' is the subject of this address.

However, before I confine myself to a few examples of 'looking at interfaces' in chemistry and physics, I should say something about the interfaces between these disciplines. According to the secondary-school text-books there is a distinct difference between chemi-

cal and physical phenomena. In the first kind the molecules change in composition, in the second kind they do not. This is a characteristic example of a boundary between chemistry and physics — an artificial boundary. Nature did not itself create the separation between chemistry, physics and engineering. Man required it for classification. However, now that we are penetrating deeper and deeper into natural phenomena and understand more about them, the lines of demarcation between the old disciplines are rapidly becoming less distinct, and indeed they are often a nuisance. Between interface chemistry and interface physics there is no dividing line to be drawn. Nowadays we draw other dividing lines, for example between low-energy physics, the world in which atoms are the building blocks, and high-energy physics, the world contained within the atomic nucleus.

Heterogeneous catalysis

There is a direct connection between looking at interfaces and the subject of my own discipline 'heterogeneous catalysis'. What is heterogeneous catalysis?

A catalyst is a substance that speeds up a chemical reaction, without itself being consumed in the process. In heterogeneous catalysis the catalyst is a solid surface, while the reaction partners are in the gaseous or liquid phase. The chemical reaction does not take place in the gaseous or liquid phase, or only to a very slight extent, but the reaction is vastly accelerated if the reaction partners, the reacting molecules, meet at the surface.

The surface acts as a meeting place. You might compare it with a marriage bureau, where men and women enter alone and leave in pairs or as married couples, whereas the bureau, the catalyst, remains unaffected. The reaction does not take place in the outside world, but at the surface, the interface, the barrier is lowered and the reaction is able to take place there. To understand how this works, it must be possible to measure and look at the interface. The problem here is that one should really look at the system actually working during the reaction. This is a considerable scientific challenge.

Catalysis is of considerable economic importance, because it opens up reaction routes that give faster reactions, can work at lower temperatures and can be more specific, that is to say cause fewer side-reactions, less fuss.

In addition to the familiar fields of catalysis in the bulk-chemicals industry and in biochemistry (enzymes) a completely different and interesting field has emerged: the growth of thin films on a substrate that acts as a catalyst. A gaseous mixture is conducted

over a surface in a state (temperature and pressure) in which it does not react in the gaseous phase, but does react at the substrate surface, and in such a way that the reaction product remains on the substrate in the form of a closely defined thin layer, or film. The surface of each newly formed layer catalyses the subsequent reaction and the layer can continue to grow from the gaseous phase. This process is known as chemical vapour deposition, CVD, and has come into such wide use that in some quarters anyone using the process is said to be 'CVDing'. An example of the process is the growth of thin layers of silica by the reaction $\text{SiCl}_4 + \text{O}_2 \rightleftharpoons \text{SiO}_2 + 2\text{Cl}_2$.

Other applications are to be found in many areas, ranging from the application of extremely hard coatings to tools to the growth of single-crystal layers on semiconductors for integrated circuits and lasers. In general, the keywords in catalytic reactions are: subtle molecular reaction mechanisms, which take place at interfaces. This brings me back to my theme: 'looking at interfaces'. The catalyst here can take all manner of forms, from very small grains on a carrier material to crystalline solids with 'large' surfaces.

The examples I have chosen relate to interfaces where one side is a solid and the other can be solid, liquid or gaseous, with 'vacuum' as a special case.

Clean surfaces

The simplest interface can be produced in a hypothetical experiment if I draw a line on a sheet of white paper and write the word 'solid' on one side of it. The other side is empty and represents an ideal vacuum, absolute 'nothingness'. This does not imply, however, that absolutely nothing happens there. The atoms at the surface of the solid are not surrounded in the same way as their friends deep inside the solid. These outer atoms will therefore have different bonds, with different bond lengths, and they will therefore occupy different positions and have a different electron structure, causing changes in all kinds of chemical and physical properties: chemical properties such as chemical reactivity and physical properties such as effective optical constants. The interface does not react *with* the vacuum, but *on* the vacuum.

As soon as we take a step towards reality and break away from the theoretical line on the sheet of paper, we see that the interface becomes more complicated. There is no such thing as an ideal vacuum; it is really a very greatly reduced gas pressure, which we describe as ultra-high vacuum. The pressure is in fact 10^{-12} to 10^{-15} atmospheres. Figures like this do not mean much except to the vacuum specialist. A vacuum with a pressure of 10^{-13} atmospheres would therefore be

better described, I think, as something one can buy for about 5000 dollars per litre (a stainless-steel shell with pumps) and as a situation in which an atom at the surface of a solid inside this vacuum will only collide with a gas molecule once every three hours. Many of these collisions will cause a chemical reaction, and the surface of the solid will very slowly become 'dirty'. There is time enough, however, to look at it in the clean state.

To look at such surfaces we need measurement methods that provide specific information about the outer atomic layers. Investigations of this kind were first started in the late fifties. They made use of the strong interaction between low-energy electrons or ions and a solid. Scattering, reflection, emission and diffraction of such particles give specific information about the surface but not about the bulk material beneath it, because the particles cannot penetrate into the bulk and then emerge again to bring information back to the detector. Examples are Low-Energy Electron Diffraction (LEED), Secondary-Ion Mass Spec-

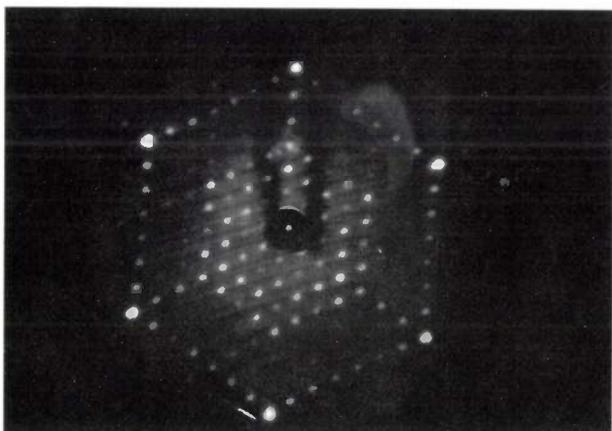


Fig. 1. The diffraction pattern of low-energy electrons (38 eV) obtained from a clean (111) silicon surface. The very bright reflections correspond to a structure that would be expected from the arrangement of the atoms inside the crystal. In between there are weak reflections at $1/7$ the spacing of the bright reflections. The fainter reflections show that the unit cell at the surface is 7×7 times larger than the unit cell corresponding to the atomic arrangement in the rest of the crystal. Diffraction micrographs such as this led to the proposal in 1959 of a model with the 7×7 structure for the surface of clean silicon, as shown in fig. 2.

trometry (SIMS), and Electron Energy-Loss Spectroscopy (EELS)^[1]. These methods give information about atomic arrangement, atomic species and molecular structure.

The use of these measurement methods relies on a relatively high vacuum, since the electrons and ions cannot travel freely in a gaseous atmosphere or a liquid. This gave a considerable impetus to research on solid/vacuum interfaces. Surfaces in vacuum could be de-

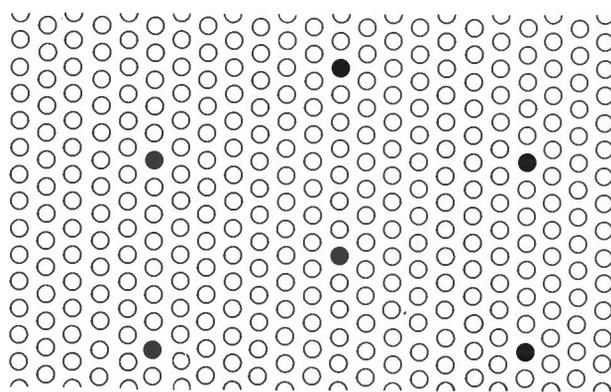


Fig. 2. Diagram of the atomic arrangement in the (111) surface of annealed silicon (the 7×7 structure). The diagram shows the arrangement for a minimum deviation of the positions of the atoms, corresponding to a diffraction pattern as in fig. 1, with open circles for the normal positions and filled circles for the deviant positions^[2].

finned more exactly, and vacuum-measurement methods became available, especially for those with generous budgets.

Many very interesting effects were found, and the phenomenon of 'surface' became much better understood. The snag was, however, that the research had a strong bias towards the use of model systems. The initial optimism that the model system could be translated into a real system was frequently followed by disappointment.

For instance, the adsorption of oxygen on an atomically clean single-crystal silicon surface provided little information about the growth of thin oxide layers in the processes for making transistors and integrated circuits. You might compare this to some extent with the situation in which a creature from outer space finds himself on receiving instructions to study the behaviour of human beings on the surface of the Earth. His attempts to do so are so hampered by clouds, rain, wind and mist that instead he writes a voluminous report on his study of two astronauts in vacuum on the surface of the moon. His conclusion is that the most important human activity is picking up stones and putting them into numbered sacks. The moral of this tale is that with good observations you can still be completely wide of the mark if you elevate your model system to the status of reality.

A familiar example of a solid/vacuum interface where the changes with respect to the bulk were directly observed is the surface of silicon. As long ago as 1959 it was demonstrated, from the diffraction of low-energy electrons, that the 'clean' silicon surface was a 'reconstructed' surface, and in fact reconstructed in a very complicated manner (fig. 1)^[2]. The unit cell on the silicon (111) surface was found to be 7×7 times larger than that of the bulk. This means that the atoms in the surface layer arrange themselves dif-

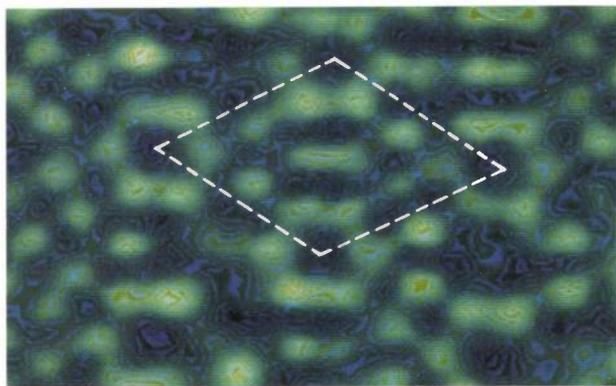


Fig. 3. A recent STM micrograph of the (111) surface of annealed silicon, in which the separately observable silicon atoms exhibit the 7×7 structure, where the dashed white line represents one unit cell. The dimensions of this unit cell were derived in 1959 from diffraction results, but this arrangement of the atoms could not be uniquely predicted from the experimental information. (Made available by courtesy of J. E. Griffith, J. A. Kubby and R. S. Becker of AT&T Bell Labs, Murray Hill, NJ, U.S.A.)

ferently, so that the regularity repeats only after seven atomic spacings (*fig. 2*). The lattice-work pattern has become a pattern of large diamonds. The frustrating aspect of the method was that, while it was very simple to observe the 7×7 structure, the theory was inadequate — and still is — for determining the new atomic positions, for understanding how exactly the 49 atoms in the large diamond-shaped unit cell are arranged, whether atoms are missing or additional ones are present, and so on.

This gave researchers a free hand to devise models, and through the years there were no lack of imaginative proposals. Fierce debates were heard at congresses and no one was able to prove that he or she was right while no method was available that could make the atoms at the surface directly visible. The breakthrough came in 1983, when the results of Scanning Tunneling Microscopy^[3] were published, a novel method, to which I shall return presently. The microscope scans the surface 'on an atomic scale', resulting in a relief map of the surface, on which the individual atoms are visible. It turned out that one of the many structures proposed was reasonably in agreement with this 'direct' observation (*fig. 3*). The matter now seems to be settled, although even with this method there are still a few reservations about the interpretation, particularly on the part of those who have a different model in mind.

As yet, however, there is not a single practical application to be found for this specific knowledge of the silicon surface, except of course for the manufacturers of surface analytical equipment, who display the Si(111)- 7×7 structure as a distinguishing feature on their publicity material. The immediate utility of the elucidation of such a structure is rather the better un-

derstanding of matters such as the behaviour of electrons in solids. Clean surfaces present theoreticians and experimenters with a problem they can really exercise their talents on.

The understanding achieved does not have to be specifically applicable. It provides the broad substructure that technical applications will be able to extend.

Measurement methods

The ideal measurement method

In my discourse so far I have shown that 'clean' surfaces in vacuum can readily be studied, but that they are not entirely satisfactory as models of the reality. Looking at 'true' interfaces calls for different methods of measurement. What would the ideal method look like?

The ideal measurement method would not disturb the interface, and would determine the chemical composition and molecular structure along it, in the *x*- and *y*-directions, and in the *z*-direction as a function of the distance from the interface. The ruler for measuring *x*, *y* and *z* has divisions in angstroms.

The main problem is the requirement that the method of measurement should be non-destructive. This means that the interface must not be disturbed either in the preparation of the specimen or during the measurement itself. The measurement must be made *in situ*.

If this requirement is left out, a great deal can be done at present. I shall touch briefly on two methods that provide information on an atomic scale, but are limited in this non-destructive, *in situ* aspect.

TEM

The first method that will 'observe' on an atomic scale is Transmission Electron Microscopy (TEM). High-energy electrons are fired through a thin specimen and form an image. Diffraction can also occur and the resultant diffraction image can be interpreted as an image of the atoms.

There are two methods of looking at an interface with TEM. For solid/gaseous or solid/liquid interfaces the solid has to be made thin enough for the electrons to pass right through it; this implies a thickness of a few dozen microns. The specimen is then introduced into the vacuum of the microscope and the solid side of the interface is observed. In the second method, which is used primarily for solid/solid interfaces, a cross-sectional specimen is made; this is again ex-

[1] H. H. Brongersma, F. Meijer and H. W. Werner, *Philips Tech. Rev.* **34**, 357-369, 1974.

[2] R. E. Schlier and H. E. Farnsworth, *J. Chem. Phys.* **30**, 917-926, 1959.

[3] G. Binnig, H. Rohrer, C. Gerber and E. Weibel, *Phys. Rev. Lett.* **50**, 120-123, 1983.

tremely thin, and the interface now runs like a line through the surface of the specimen.

An example of the first method is an examination of the nucleation of a surface with palladium, so that another metal, such as copper, can be catalytically grown on the surface. The palladium nuclei are observed on a titanium-oxide layer a few microns thick. The titanium oxide itself is supported by a silicon-nitride film about ten microns thick [4]. The nuclei consist of a few hundred palladium atoms (*fig. 4*). The crystal structure and the distances between the atoms can be determined

and the arsenic hydride. The crystal is built up layer upon layer.

The results of the synthetic activity can be analysed afterwards by TEM. The specimens, slices of material perpendicular to the interface, are etched to make them extremely thin, so that the electrons can pass through the specimen and form an image. The TEM micrographs show layers with a resolution of one atomic layer, and even the step in the interface of one atomic layer can be seen. TEM measurements of this kind provide useful information about the CVD process.

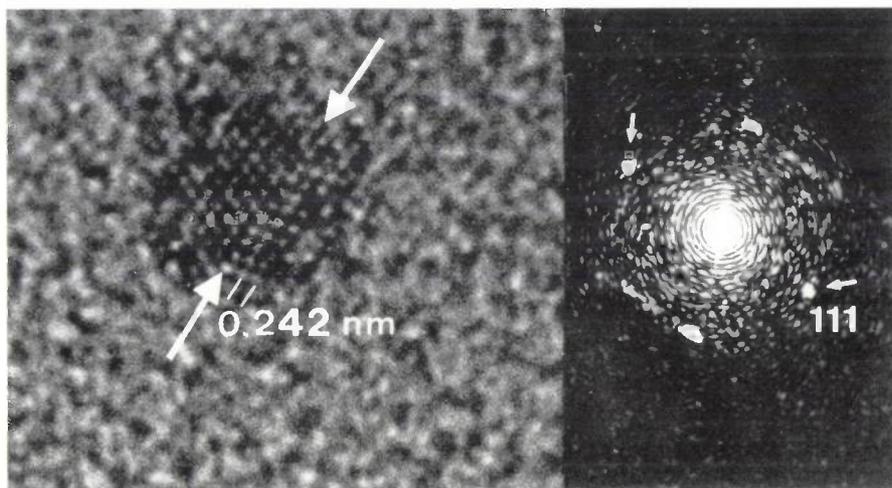


Fig. 4. *Left*: TEM micrograph of a metallic nucleus consisting of no more than a few hundred palladium atoms (obtained on a titanium-oxide layer a few microns thick, supported by a silicon-nitride film a few tens of microns thick). It can clearly be seen that the nucleus has a crystalline structure, in which the familiar effect of 'twinning' has taken place (the mirror plane corresponds to the direction of the arrows). *Right*: TEM diffraction micrograph of the nucleus [5].

from the diffraction of the electrons. In this example it was concluded that what was observed was indeed metallic palladium with the normal crystal structure. A palladium particle of a few hundred atoms behaves like a piece of metal. In some cases, especially after oxygen adsorption, an expansion of the lattice of 5 to 10% could be measured [6]. This information is immediately relevant to the understanding of the catalytic action.

Another example is to be found in semiconductor materials. Beautiful micrographs have been made of cross-sections of epitaxial layers of GaAs and AlGaAs (*fig. 5*) [6], where metal-organic compounds were used to grow successive monoatomic epitaxial layers on a substrate surface.

'Epitaxial' means that the crystal lattice continues from one composition, GaAs, to the next, AlGaAs. The metal-organic compounds, such as trimethyl gallium, dissociate on the surface and the gallium atom arrives at exactly the right place in the crystal lattice. The same thing happens with the triethyl aluminium

Transmission electron microscopy has aspects of the ideal measurement method. It does however require an elaborate preparation technique, followed by measurements in vacuum. Still pictures are obtained after the event, rather than a cinematographic record of events as they take place. The method tells us a good deal about the solid side of the interface, but hardly anything about the chemical structure of the molecules that react with the surface of the catalyst.

STM

Another method I shall touch on briefly, after TEM, is Scanning Tunnelling Microscopy, STM, a relatively new method with resolution on an atomic scale.

Electrons from the surface of a specimen tunnel towards a very sharply defined point of metal, located at a distance of a few angstroms. The tunnelling current depends very closely on the distance between surface atom and metal point. By scanning the surface with the point, at a constant tunnelling current, in the x -, y - and z -directions, an atomic topograph is obtained.

The method depends entirely on the fantastic mechanical precision with which the metal point can be manipulated in the sub-angstrom range. In this way, for instance, the 7×7 structure mentioned earlier of the silicon surface has been observed.

The STM method is used in vacuum, but in principle it can also be applied in a gaseous atmosphere or even in a liquid. The method has certain aspects of the 'ideal' measurement method, but it will never be possible to use it on solid/solid interfaces. Nor is it likely that it can be used to follow the reactions directly at

boundary faces. One of the principal advantages of this is that we can then study dynamic phenomena, reactions taking place at interfaces. In this respect there is still nothing that even approaches an 'ideal' method. Likely methods include those that use 'light', X-rays or acoustic vibrations.

The interesting thing about optical methods is that photons, unlike the electrons and ion beams I was discussing just now, interact very little with matter. In methods that use electrons or ions, however, the surface sensitivity is a direct result of the strong interac-

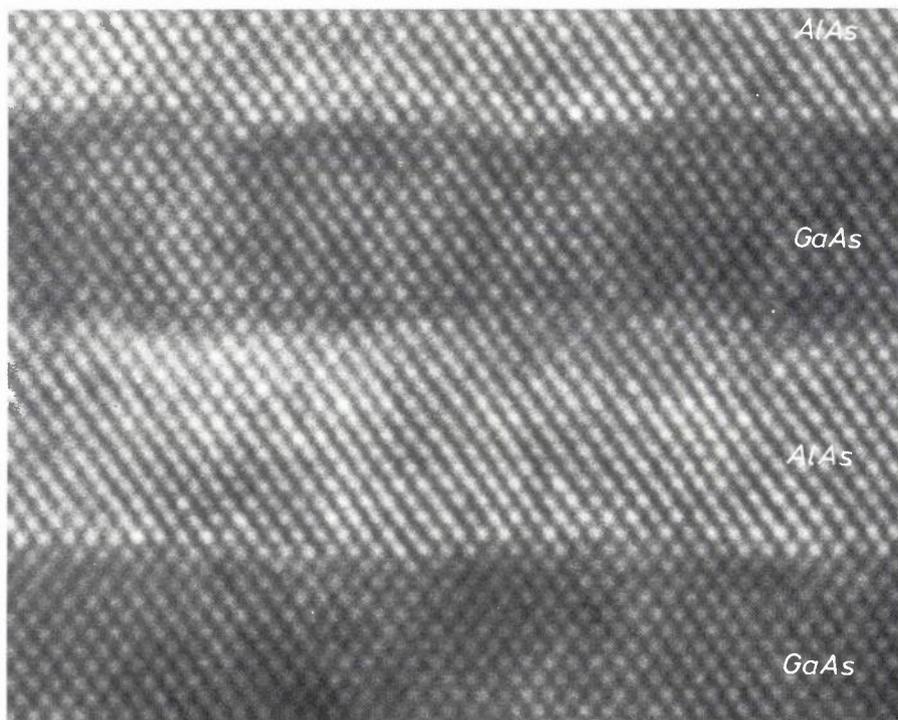


Fig. 5. TEM micrograph of a cross-section of epitaxial layers of GaAs (dark layers) and AlAs. The resolution is never less than two atoms: Ga and As or Al and As. The crystal lattice continues uninterrupted at the transition from one composition to the other, and here and there a step in the interface of one atomic layer can be seen^[6].

interfaces. TEM and STM are examples of highly advanced measurement methods that provide information about the atomic structure from surfaces. These methods and various other analytical techniques, each with their own strengths and weaknesses, have undoubtedly led to a better understanding of interface phenomena. More scientific challenges have still to be met, however, and much research of practical relevance remains to be done.

Non-destructive measurement methods

It seems to me that a field still lies fallow for non-destructive measurements on interfaces. With such measurements, interfaces can be observed without first separating them, without having to make two

tion with matter, and this has the disadvantage that the measurements have to be performed in a vacuum, since the electrons and ions are scattered in air or any other medium. With light rays (or, to be more accurate, electromagnetic radiation in certain wavelength ranges) there is no such limitation, yet specific information about the interface can nevertheless be obtained in addition to the signal that comes from the bulk. We can explain this in general terms with an example.

[4] J. W. M. Jacobs and J. F. C. M. Verhoeven, *J. Microsc.* **143**, 103-116, 1986.

[6] J. W. M. Jacobs and D. Schryvers, *J. Catalysis* **103**, 436-449, 1987.

[6] M. P. A. Vieggers, A. F. de Jong and M. R. Leys, *Spectrochim. Acta* **40B**, 835-845, 1985; M. R. Leys, M. P. A. Vieggers and G. W. 't Hooft, *Philips Tech. Rev.* **43**, 133-142, 1987.

A light ray is reflected when it encounters an optical change. This need not be associated with energy absorption; even the clearest non-absorbent glass reflects sun-rays. The reflected light ray contains information about the interface, because the interface is *different*, is different from the bulk. A nice example of this is the reflection of light at the surface of a silicon crystal. Silicon has a cubic crystal structure and is therefore optically isotropic, that is to say it behaves identically in all directions. The surface, on the other hand, is anisotropic because of the fact that it is a surface. In a recent experiment [7] measurable changes in the reflection coefficient were observed by rotating the surface around the axis perpendicular to the surface. These changes must be due to optical anisotropy and therefore contain direct information about the surface, because the bulk material below it is optically isotropic.

Optical reflection techniques — ellipsometry is a good example — can provide precisely that information about the interface in operation, *in situ*. The ellipsometric method measures changes in both the amplitude and the phase of the light upon reflection. This may be compared by an interference method, and ellipsometry is therefore extremely accurate in showing thin layers at an interface. The resolution in the *z*-direction easily approaches the range of monoatomic layers. In the reaction of a clean silicon surface with oxygen and many other molecular species it has been possible to demonstrate the presence of coatings of 0.01 of a monolayer, that is to say one oxygen atom for every hundred atoms on the silicon surface [8][9]. In the *x*- and *y*-directions the resolution is not in the atomic range but in the range of square millimetres to square microns, the limit being set by the wavelength of the light. The ellipsometer is therefore eminently suited for studying reactions of flat interfaces, as for example in CVD. There are also interesting examples of oxidation reactions on solid/solid interfaces in semiconductor technology that have been studied by ellipsometry because it was possible to 'look inside' [10].

The speed of the method is now about a hundred measurements per second, which is sufficient for measuring many reactions on interfaces. Ellipsometry makes use of spectroscopic capabilities. Measurements can be made in the near ultraviolet, in the visible and the near infrared. A logical but experimentally difficult extension into the infrared range of molecular vibrations has not yet been achieved.

I shall conclude this somewhat more detailed account of ellipsometry with the comment that here, too, there is strength in combination. A combination of the *in situ* measurement of dynamic effects by ellipsometry with an 'atomic micrograph' subsequently obtained by means of transmission electron microscopy

or with a 'molecular structure micrograph' obtained by means of EELS can bring us quite a bit closer to the ideal measurement method.

Without further explanation, so that only the initiated will be completely in the picture, I will just mention a few other methods that are of interest, or could be, for non-destructive investigations on interfaces.

- Surface-enhanced Raman scattering and infrared reflection spectroscopy, which can provide information about molecular structures.

- Fluorescence and phosphorescence of probe molecules [11], where the optical effect is sensitive to the immediate environment of the molecule, for example network rigidity in polymers. Many new developments are afoot in optical spectroscopy. It is very likely that the spin-off from these developments will be useful in interface studies, in both ordinary and nonlinear optics. There are prospects of *in situ* measurements in the X-ray wavelength range as well.

- Extended-X-ray absorption fine structure or EXAFS for short, which provides information about atomic distances in matter, including amorphous substances. The application of EXAFS in the study of rhodium catalysts in their reaction environment is yielding spectacular results [12]. In particular, measurements have been made of reversible reactions of the rhodium particle with carbon monoxide, in which the rhodium-rhodium-bonds were broken and then restored after removal of the carbon monoxide. To get some idea of such a catalyst particle you have to think in terms of a fragment consisting of 10 to 20 atoms.

Conclusion

In my address I have sought to emphasize that, while much is already known about interfaces, the step from model system to the reality still calls for a great deal of fundamental research. 'Looking at interfaces', using methods in which destructive preparation techniques and a high vacuum are not required, is a field that is open to interesting, advanced and applied research.

Teaching at the university first concentrates on the things that are known: the theory, the facts, the existing experimental methods. By doing research the student then learns to use this knowledge actively and creatively. The impending curtailment of the duration of undergraduate studies in the Netherlands must not lead to an impoverishment of that essential experience in actually doing research.

This brings me naturally to another very interesting interface, the interface between University or Polytechnic and Industry, or more generally the Employer. For the student this interface is the transition from the

past to the future. Nor must this interface be a mere common boundary; here again interaction is essential.

It is in everyone's interest that the man or woman who leaves the university should occupy the right position in working life, in the job that exactly matches his or her interests, capacity and education. These three parameters are not numbers, they are not scalar quantities, but vectors, directions in which development can take place at the proper rate.

In crossing the interface between university and industry these vectors must not suffer abrupt changes. This can best be guaranteed if the interface is both a 'common boundary' and the 'place where interaction occurs', as in the dictionary definitions.

Looking for a place in working life involves a choice of options, a choice by the job-seeker and the employer. The chance of making the right choice is greatest if the job-seeker has a realistic idea of what he or she wants. This idea is one that must be formed

during the years at university, and tried out and tested. It will also help the employer to assess a candidate's suitability for a particular kind of work.

The thread running through my address has been the concept of 'interface'. Looking at interfaces will greatly help us to understand the interactions we find there, and armed with this knowledge we can actively influence these interfaces. So we can prepare an efficient catalytic surface, improve the adhesion between materials, combat corrosion — and even provide chemistry students with the knowledge that will fit them for a job in industry.

'Looking at interfaces' may still be far from fully optimized, but methods both existing and new appear to offer many prospects of further progress.

[7] D. E. Aspnes, J. Vac. Sci. & Technol. B 3, 1138-1141, 1985.

[8] R. J. Archer and G. W. Gobeli, J. Phys. & Chem. Solids 26, 343-351, 1965.

[9] G. A. Bootsma and F. Meijer, Surf. Sci. 14, 52-76, 1969; F. Meijer and G. A. Bootsma, Surf. Sci. 16, 221-233, 1969; F. Meijer and G. A. Bootsma, Philips Tech. Rev. 32, 131-140, 1971.

[10] J. B. Theeten, D. E. Aspnes, F. Simonet, M. Erman and P. C. Müräu, J. Appl. Phys. 52, 6788-6797, 1981.

[11] This will be the subject of a forthcoming article in this journal.

[12] H. F. J. van 't Blik, J. B. A. D. van Zon, T. Huizinga, J. C. Vis, D. C. Koningsberger and R. Prins, J. Am. Chem. Soc. 107, 3139-3147, 1985.

Summary. The text is based on the address delivered by the author on 7th March 1986 at his inauguration as Visiting Professor of Heterogeneous Catalysis at the University of Leiden. After discussing the 'interface' concept, heterogeneous catalysis and clean surfaces, he mentions as the characteristics of the ideal measuring method a resolution of the order of atomic distances and the capability of performing non-destructive measurements *in situ*. Methods treated that show the first characteristics are TEM and STM. Methods exhibiting the second characteristic are those that depend on light, X-rays or acoustic vibrations, with special emphasis on ellipsometry. Combining both types of method should give even better results.

Scientific publications

These publications are contributed by staff of laboratories and plants that form part of or cooperate with enterprises of the Philips group of companies, particularly by staff of the research laboratories mentioned below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	A
Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	B
Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	E
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	R

T. Helzel & G. Martens	H	Optical slip ring for off-axis high-bit-rate data transmission	Appl. Opt. 25	775-779	1986
K. Kobs, H. Dimigen, H. Hübsch, H. J. Tolle, R. Leutenecker (<i>Fraunhofer Inst. für Festkörpertechnol., München</i>) & H. Ryssel (<i>FhG Arbeitsgruppe für integrierte Schaltungen, Erlangen</i>)	H	Improved tribological properties of sputtered MoS _x films by ion beam mixing	Appl. Phys. Lett. 49	496-498	1986
H. van Houten, B. J. van Wees (<i>Delft Centre for Submicron Technol., Delft</i>), M. G. J. Heijman & J. P. André	E,L	Submicron conducting channels defined by shallow mesa etch in GaAs-AlGaAs heterojunctions	Appl. Phys. Lett. 49	1781-1783	1986
K. Mohammed, D. A. Cammack, R. Dalby, P. Newbury, B. L. Greenberg, J. Petruzzello & R. N. Bhargava	N	Effect of lattice mismatch in ZnSe epilayers grown on GaAs by molecular beam epitaxy	Appl. Phys. Lett. 50	37-39	1987
B. Greenberg & J. Ladell	N	Modulation of Renninger scan intensity: A new x-ray technique to characterize epitaxial structures	Appl. Phys. Lett. 50	436-438	1987
P. Röschmann & C. Wetzel	H	Die Bedeutung der HF-Spulen für die Kernspintomographie	Biomed. Tech. 31	178-185	1986
G. Aissing*, R. Broer*, W. C. Nieuwpoort* (<i>*Univ. Groningen</i>) & L. F. Feiner	E	Interstitial 3d transition metal impurities in silicon: an AB initio cluster study	Defects in Semiconductors, H. J. von Bardeleben (ed.), Mater. Sci. Forum, Vol. 10-12, Trans Tech, Aedermannsdorf	711-716	1986
F. M. Klaassen	E	Device aspects of megabit RAMs	ESSDERC 1986, Cambridge 1986 (Inst. Phys. Conf. Ser. No. 82)	113-133	1987
H. Becher*, M. Schlüter*, D. G. Mathey*, W. Bleifeld* (<i>*Univ. Hamburg</i>), E. Klotz, P. Haaker, R. Linde & H. Weiss	H	Coronary angiography with flashing tomosynthesis	Eur. Heart J. 6	399-408	1985
R. Woltjer, R. Eppenga, J. Mooren, C. E. Timmering & J. P. André	E,L	A new approach to the quantum Hall effect	Europhys. Lett. 2	149-155	1986
W. J. A. Goossens	E	The smectic A-smectic C phase transition: a molecular statistical theory	Europhys. Lett. 3	341-346	1987
K. G. Freeman	R	The history of CRT colour television displays — a personal view	IEE Conf. Proc. No. 271	4pp.	1986
R. L. Maresca	N	A 400-Watt, tri-state switching controller for reciprocating linear motors	IEEE Power Electronics Specialist Conf., Vancouver 1986	587-592	1986

J. Khurgin, W. Seemungal, S. Colak & A. Hebling	N	Single longitudinal mode operation of the electron-beam-pumped semiconductor laser	IEEE QE-22	1158-1161	1986
F. Fonsalas & J. Y. Lejard	L	An edge enhancement algorithm using a polynomial transformation	IEEE Trans. CE-32	151-154	1986
J. P. Michel	L	8 bit converters for video applications	IEEE Trans. CE-32	630-635	1986
H. Sari, L. Desperben (<i>Alcatel Thomson Radioteleph., Gennevilliers</i>) & S. Moridi	L	Minimum mean-square error timing recovery schemes for digital equalizers	IEEE Trans. COM-34	694-702	1986
G. W. Bruning	N	A new high-voltage oscillator	IEEE Trans. IE-33	171-175	1986
V. Zieren, S. B. Luitjens, C. P. G. Schrauwen, J. P. C. Bernards, R. W. de Bie & M. Piena	E	Properties of one-sided probe heads on double-layer perpendicular recording media	IEEE Trans. MAG-22	370-372	1986
M. J. Vos*, S. B. Luitjens, R. W. de Bie & J. C. Lodder* (<i>Univ. of Technol., Enschede</i>)	E	Magnetization transitions obtained by deconvolution of measured replay pulses in perpendicular magnetic recording	IEEE Trans. MAG-22	373-375	1986
W. El-Kamali, J.-P. Grimm, R. Meierer* & C. Tsironis* (<i>SPAR Aerospace, Montreal</i>)	L	New design approach for wide-band FET voltage-controlled oscillators	IEEE Trans. MTT-34	1059-1063	1986
A. J. M. Houtsma (<i>Inst. for Perception Res., Eindhoven</i>), N I. Durlach* & D. M. Horowitz* (<i>Massachusetts Inst. of Technol., Cambridge, MA</i>)		Comparative learning of pitch and loudness identification	J. Acoust. Soc. Am. 81	129-132	1987
P. Hansen & M. Hartmann	H	Magnetic and magneto-optic properties of amorphous Gd-FeAu films	J. Appl. Phys. 59	859-863	1986
H. Mani*, A. Joullie*, F. Karouta* (<i>EM, Montpellier</i>) & C. Schiller	L	Low-temperature phase diagram of the Ga-As-Sb system and liquid-phase-epitaxial growth of lattice-matched GaAsSb on (100) InAs substrates	J. Appl. Phys. 59	2728-2734	1986
P. Hansen, C.-P. Klages & K. Witter	H	Magnetic and magneto-optic properties of praseodymium- and bismuth-substituted yttrium iron garnet films	J. Appl. Phys. 60	721-727	1986
M. Erman & J. B. Theeten	L	Spatially resolved ellipsometry	J. Appl. Phys. 60	859-873	1986
J. R. Morante*, J. Samitier*, A. Pérez*, H. Alterarrea* (<i>Univ. Barcelona</i>) & S. Gourrier	L	Analysis of the near-intrinsic and extrinsic photocapacitance due to the EL2 level in boron implanted GaAs	J. Appl. Phys. 60	1661-1669	1986
J.-P. Krumme, V. Doormann, B. Strocka & P. Willich	H	Selected-area sputter epitaxy of iron-garnet films	J. Appl. Phys. 60	2065-2068	1986
R. G. Pratt, J. Hewett, P. Capper, C. L. Jones* & N. Judd* (<i>Mullard, Southampton</i>)	R	Minority-carrier lifetime in doped and undoped n-type Cd _x Hg _{1-x} Te	J. Appl. Phys. 60	2377-2385	1986
B. Strocka	H	Selected-area liquid-phase epitaxy of iron garnet films applying local ion implantation of the substrate	J. Appl. Phys. 60	2977-2979	1986
S. D. Brotherton, J. P. Gowers, N. D. Young, J. B. Clegg & J. R. Ayres	R	Defects and leakage currents in BF ₂ implanted pre-amorphized silicon	J. Appl. Phys. 60	3567-3575	1986
J. W. M. Jacobs & D. Schryvers (<i>Univ. Antwerpen</i>)	E	A high-resolution electron microscopy study of photodeposited Pd particles on TiO ₂ and their oxidation in air	J. Catalysis 103	436-449	1987
P. F. Bordui, J. J. Zola, G. M. Loiacono & G. Kostecy	N	Turntable seed holder for crystal growth from low-temperature solution	J. Cryst. Growth 78	453-456	1986
P. Capper, B. C. Easton, P. A. C. Whiffin & C. D. Maxey	R	Electrical properties and annealing behaviour of Cd _x Hg _{1-x} Te grown by LPE and MOVPE	J. Cryst. Growth 79	508-514	1986
P. A. C. Whiffin, B. C. Easton, P. Capper & C. D. Maxey	R	Computer controlled deposition of CMT heterostructures by MOVPE	J. Cryst. Growth 79	935-939	1986
K. H. J. Buschow, D. B. de Mooij & H. M. van Noort	E	Properties of metastable ternary compounds and amorphous alloys in the Nd-Fe-B system	J. Less-Common Met. 125	135-146	1986
D. B. de Mooij & K. H. J. Buschow	E	Structure and magnetic properties of U ₃ Sb ₃ Cu ₂	J. Less-Common Met. 125	239-245	1986

- | | | | | |
|---|--|--|-----------|------|
| J. J. G. Willems (<i>Philips Lighting Div., Eindhoven</i>) & K. H. J. Buschow
<i>E</i> | From permanent magnets to rechargeable hydride electrodes | J. Less-Common Met. 129 | 13-30 | 1987 |
| R. B. Helmholtz (<i>ECN, Petten</i>) & K. H. J. Buschow
<i>E</i> | Crystallographic and magnetic structure of Ni_2MnSn and NiMn_2Sn | J. Less-Common Mater. 128 | 167-171 | 1987 |
| P. J. W. Severin
<i>E</i> | On the relation between passive optical fiber component properties and the network configuration | J. Lightwave Technol. LT-4 | 1425-1433 | 1986 |
| P. J. W. Severin & W. H. Bardeol
<i>E</i> | Differential mode loss and mode conversion in passive fiber components measured using the chromatic launching and the central spot far field techniques | J. Lightwave Technol. LT-4 | 1640-1646 | 1986 |
| J. N. Chapman*, I. R. McFadyen* & J. P. C. Bernardis* (<i>*Univ. Glasgow</i>)
<i>E</i> | Investigation of Cr segregation within rf-sputtered CoCr films | J. Magn. & Magn. Mater. 62 | 359-366 | 1986 |
| P. P. J. van Engelen & K. H. J. Buschow
<i>E</i> | Kerr effect in $\text{R}_2\text{Fe}_{14}\text{B}$ and $\text{R}_2\text{Co}_{14}\text{B}$ compounds | J. Magn. & Magn. Mater. 66 | 291-293 | 1987 |
| J. Khurgin
<i>N</i> | Theory of bistability in the face-pumped laser with bimolecular recombination | J. Opt. Soc. Am. B 4 | 86-102 | 1987 |
| E. W. Meijer & R. J. M. Zwiers
<i>E</i> | Molecular mobility of polyepoxide networks as revealed by emission spectroscopy and mechanical analysis | Macromolecules 20 | 332-338 | 1987 |
| G. McKinnon, P. Röschmann & R. Tischler
<i>H</i> | A contribution to the art of magnetic resonance imaging | Magn. Resonance Imaging 3 | 417-418 | 1985 |
| D. Kunz
<i>H</i> | Use of frequency-modulated radiofrequency pulses in MR imaging experiments | Magn. Resonance Medicine 3 | 377-384 | 1986 |
| D. Kunz
<i>H</i> | Double pulse echoes — a novel approach for fat-water separation in magnetic resonance imaging | Magn. Resonance Medicine 3 | 639-643 | 1986 |
| J. C. Jacco
<i>N</i> | The infrared spectra of KTiOPO_4 and a $\text{K}_2\text{O-P}_2\text{O}_5\text{-TiO}_2$ glass | Mat. Res. Bull. 21 | 1189-1194 | 1986 |
| F. Delahaye*, D. Dominguez* (<i>*LCIE, Fontenay-aux-Roses</i>), F. Alexandre (<i>CNET, Bagneux</i>), J. P. André, J. P. Hirtz** & M. Razeghi** (<i>**LCR, Thomson CSF, Orsay</i>)
<i>L</i> | Precise quantized Hall resistance measurements in $\text{GaAs/Al}_x\text{Ga}_{1-x}\text{As}$ and $\text{In}_x\text{Ga}_{1-x}\text{As/InP}$ heterostructures | Metrologia 22 | 103-110 | 1986 |
| A. W. Kofschoten
<i>E</i> | The role of noble gas ion bombardment in etching reactions | Nucl. Instrum. & Methods Phys. Res. B19/20 | 1001-1008 | 1987 |
| G. N. A. van Veen, F. H. M. Sanders, J. Dieleman, P. C. Zalm, E. D. J. Oostra* & A. E. de Vries* (<i>*FOM, Amsterdam</i>)
<i>E</i> | Extension of the model for Ar^+ ion induced etching of Si by SF_6 | Nucl. Instrum. & Methods Phys. Res. B19/20 | 1022-1025 | 1987 |
| M. A. Misdaq (<i>Univ. Rabat</i>), G. Blondiaux*, J. P. André, M. Hage Ali (<i>CNRS, Strasbourg</i>), M. Valadon*, G. J. Maggiore (<i>Los Alamos National Lab., New Mexico</i>) & J. L. Debrun* (<i>*CNRS, Orleans</i>)
<i>L</i> | Use of channeling in association with charged particle activation to study the position of light elements at trace level in crystals: the case of carbon in GaAlAs prepared by MO-VPE | Nucl. Instrum. & Methods Phys. Res. B15 | 328-332 | 1986 |
| J. W. Orton
<i>R</i> | Physics in an industrial research laboratory | Phys. Educ. 22 | 79-84 | 1987 |
| P. Dawson, K. J. Moore, G. Duggan, H. I. Ralph & C. T. B. Foxon
<i>R</i> | Unambiguous observation of the 2s state of the light- and heavy-hole excitons in GaAs-(AlGa)As multiple-quantum-well structures | Phys. Rev. B 34 | 6007-6010 | 1986 |
| K. J. Moore, P. Dawson & C. T. Foxon
<i>R</i> | Observation of luminescence from the 2s heavy-hole exciton in GaAs-(AlGa)As quantum-well structures at low temperature | Phys. Rev. B 34 | 6022-6025 | 1986 |
| P. K. Larsen, L. F. Feiner & P. Friedel
<i>E,L</i> | Electronic structure of CF_3 radicals on GaAs (001) | Phys. Rev. B 35 | 757-764 | 1987 |
| R. Woltjer, J. Mooren, J. Wolter (<i>Univ. of Technol., Eindhoven</i>), J.-P. André & G. Weimann (<i>Forschungsinst. der Deutschen Bundespost, Darmstadt</i>)
<i>E,L</i> | Four-terminal quantum Hall and Shubnikov-de Haas measurements with pulsed electron fields | Physica 134B | 352-355 | 1985 |

- | | | | | | |
|--|-----|---|--|-----------|------|
| W. A. P. Claassen (<i>Philips Elcoma Div., Nijmegen</i>) | | Ion bombardment-induced mechanical stress in plasma-enhanced deposited silicon nitride and silicon oxynitride films | Plasma Chem. & Plasma Processing 7 | 109-124 | 1987 |
| E. W. Meijer | E | The synthesis and structure of bis(acetylacetonato) (2-hydroxyphenolato)-aluminium: a stable dimer | Polyhedron 6 | 525-529 | 1987 |
| J. G. Kloosterboer & G. F. C. M. Lijten | E | Deuterium isotope and oxygen effects on the rate of photopolymerization of a diacrylate | Polym. Commun. 28 | 2-5 | 1987 |
| J. L. Gentner, J. N. Patillon, J. P. André, B. G. Martin, C. Mallet-Mouko, J. P. Chané & G. M. Martin | L | A new semi-planar InGaAs/InP pin photodiode prepared by metalorganic vapour phase epitaxy | Proc. 12th ECOC, Barcelona 1986 | 281-284 | 1986 |
| E. H. L. Aarts & J. H. M. Krost | E | Simulation of learning in parallel networks based on the Boltzmann machine | Proc. 2nd Eur. Simulation Congr., Antwerp 1986 | 391-398 | 1986 |
| R. Woltjer, J. Mooren, J.-P. André & G. Weimann (<i>Forschungsinst. der Deutschen Bundespost, Darmstadt</i>) | E,L | Four-terminal quantum Hall and Shubnikov-de Haas measurements with pulsed electric fields | Proc. Conf. on Condensed matter, Berlin 1985 | 419-427 | 1985 |
| M. Fouassier, C. Piaget & E. Roaux | L | Experimental and theoretical evaluations of 2nd and 3rd generation intensifier viewing ranges | Proc. IEE Conf. on Photoelectronic imaging, London 1985 | 9-12 | 1985 |
| M. Lemonier, J. C. Richard, C. Piaget, M. Petit & M. Vittot | L | Photon-in and electron-in CCD arrays for image read-out in tubes | Proc. IEE Conf. on Photoelectronic imaging, London 1985 | 74-77 | 1985 |
| H. Sari, L. Desperben & S. Moridi | L | A new class of frequency detectors for carrier recovery in QAM systems | Proc. IEEE Int. Conf. on Commun., Toronto 1986 | 482-486 | 1986 |
| H. Sari, S. Morida, L. Desperben & P. Vandamme (<i>CNET, Lannion</i>) | L | Baseband equalization and carrier recovery in digital radio systems | Proc. IEEE Int. Conf. on Commun., Toronto 1986 | 1460-1465 | 1986 |
| A. A. Shaulov, M. E. Rosar, W. A. Smith & B. M. Singer | N | Composite piezoelectrics for ultrasonic transducers | Proc. IEEE Int. Symp. on Applications of ferroelectrics, Bethlehem, PA, 1986 | 231-234 | 1986 |
| W. A. Smith | N | Composite piezoelectric materials for medical ultrasonic imaging transducers — a review | Proc. IEEE Int. Symp. on Applications of ferroelectrics, Bethlehem, PA, 1986 | 249-256 | 1986 |
| P. Gamand | L | A complete small size 2 to 30 GHz hybrid distributed amplifier using a novel design technique | Proc. IEEE MTT-S Int. Microwave Symp. Digest, Baltimore, MD, 1986 | 343-346 | 1986 |
| B. Fitzpatrick, J. Khurgin, P. Harnack & D. de Leeuw | E,N | Low threshold visible electron-beam-pumped-lasers | Proc. Int. Electron Devices Meeting, Los Angeles, CA, 1986 | 630-632 | 1986 |
| M. Wolny, P. Chambery, A. Briere & J. P. André | L | Low noise high electron mobility transistors grown by MOVPE | Proc. Int. Conf. on Highspeed electronics, Stockholm 1986 | 148-150 | 1986 |
| G. Martens, T. Helzel & J. Kordts | H | Optical detection of respiration motion and heart beats for magnetic resonance imaging (MRI) applications | Proc. SPIE 586 | 250-259 | 1985 |
| J. Kosanetzkey, G. Harding & U. Neitzel | H | Energy resolved X-ray diffraction CT | Proc. SPIE 626 | 137-142 | 1986 |
| P. M. Frijlink | L | An introduction to OMVPE of III-V compounds | Proc. Winter School, Heterojunctions and Semiconductor Superlattices, Les Houches 1985 | 226-236 | 1986 |
| J. C. Brice | R | A numerical description of the Cd-Hg-Te phase diagram | Prog. Cryst. Growth & Charact. 13 | 39-61 | 1986 |
| J. C. Brice, P. Capper, C. L. Jones* & J. J. G. Gosney* (<i>*Mullard, Southampton</i>) | R | ACRT: a review of models | Prog. Cryst. Growth & Charact. 13 | 197-229 | 1986 |

- L. M. G. Feijs & H. B. M. Jonkers *E* Transformational design: an annotated example Program specification and transformation, L. G. L. T. Meertens (ed.), Elsevier Science, Amsterdam 89-112 1987
- P. Röschmann & R. Tischler *H* Surface coil proton MR imaging at 2T¹ Radiology **161** 251-255 1986
- E. P. Honig *E* Theory of oscillatory experiments with a coaxial cylinder viscometer Rheol. Acta **26** 2-6 1987
- M. Duseaux, S. Martin, J. P. Chevalier (C.E.C.M.-CNRS, Vitry-sur-Seine) *L* Microprecipitates in bulk GaAs: EL2 distribution in bulk and annealed GaAs Semi-Insul. III-V Mater. 221-226 1986
- J. C. M. Henning & J. P. M. Ansems *E* A new model of deep donor centres in Al_xGa_{1-x}As Semicond. Sci. Technol. **2** 1-13 1987
- D. E. Lacklison & P. Capper (Mullard, London) *R* Minority carrier lifetime in doped and undoped p-type Cd_xHg_{1-x}Te Semicond. Sci. Technol. **2** 33-43 1987
- M. G. Collet *E* Solid-state image sensors Sensors & Actuators **10** 287-302 1986
- R. J. Nicholas*, R. G. Clark*, A. Usher* (*Clarendon Lab., Oxford), C. T. Foxon & J. J. Harris *R* High order fractional quantisation in a two dimensional system Solid State Commun. **60** 183-187 1986
- F. Berz *R* Transport of electrons in monolithic hot electron Si transistors Solid-State Electron. **29** 1213-1222 1986
- P. W. J. M. Bouwmans & J. J. A. M. Vrakking *E* The widths and shapes of about 350 prominent lines of 65 elements emitted by an inductively coupled plasma Spectrochim. Acta **41B** 1235-1275 1986
- J. Biesterbos, A. Bouwer, G. van Engelen, G. van de Looij & J. van der Werf *E* A new lens for submicron lithography and its consequences for wafer stepper design Proc. SPIE **633** 34-43 1986
- D. Visser *E* Design and measurement of replicated aspheric compact disc objective lenses Proc. SPIE **645** 49-52 1986
- J. J. Harris, C. T. Foxon, D. E. Lacklison & K. W. J. Barnham (Imp. College, London) *R* Scattering mechanisms in (Al, Ga)As/GaAs 2deg structures Superlattices & Microstruct. **2** 563-568 1986
- B. A. Joyce, P. J. Dobson, J. H. Neave & J. Zhang (Imp. College, London) *R* Surface effects and growth dynamics in MBE of III-V compounds Surf. Sci. **178** 110-123 1986
- M. H. Verhaegen (NASA Ames Res., Moffet Field, CA) & P. van Dooren *B* A reduced order observer for descriptor systems Syst. & Control Lett. **8** 29-37 1986
- P.-J. Courtois & P. Semal *B* Bounds on conditional steady-state distributions in large Markovian and queueing models Teletraffic analysis and computer performance evaluation, O. J. Boxma, J. W. Cohen, H. C. Tijms (eds), Elsevier Science, Amsterdam 499-520 1986
- P. Houdy, E. Ziegler & L. Névot (Inst. d'Optique, Orsay) *L* Sputtering techniques applied to realization of ultrathin layered stacks (of the order of nanometres): *in situ* ellipsometry control system Thin Solid Films **141** 99-109 1986
- R. Memming, H. J. Tolle & P. E. Wierenga *E,H* Properties of polymeric layers of hydrogenated amorphous carbon produced by a plasma-activated chemical vapour deposition process II: Tribological and mechanical properties Thin Solid Films **143** 31-41 1986
- R. Memming *H* Properties of polymeric layers of amorphous hydrogenated carbon produced by a plasma-activated chemical vapour deposition process I: Spectroscopic investigations Thin Solid Films **143** 279-289 1986
- O. Bonnefous & P. Pesqué *L* Time domain formulation of pulse-doppler ultrasound and blood velocity estimation by cross correlation Ultrason. Imaging **8** 73-85 1986
- J. A. J. Roufs, M. A. M. Leermakers, M. C. Boschman (Inst. for Perception Res., Eindhoven) *E* Criteria for the subjective quality of visual display units Work with display units 86, K. Knave & P.-G. Widebäck (eds), Elsevier Science, Amsterdam 412-418 1987

Automatic segmentation of speech into diphones

J. P. van Hemert

With the advance of automation the number of people who transact business directly with the public is steadily declining. It is becoming increasingly common, for example, to draw money from automatic cash dispensers. In such situations the customer has to carry out instructions given by the machine. Usually the instructions appear on a screen, but sometimes they are spoken. In such a case the number of instructions that can be given is limited. A method of producing spoken instructions just as freely and flexibly as text on a screen consists in joining small speech segments together and making them audible. This article describes the preparation of the 'library' of segments (diphones).

Introduction

Speech has been with us since time immemorial. And today, even in situations where it is possible to communicate in some other direct way, e.g. with a keyboard and display screen, many people will still prefer the spoken word. One possible reason for this is that people can then make other 'observations' — visual ones, perhaps — while they talk or listen. On the other hand, communication between man and machine usually requires a keyboard and a display. At the Institute for Perception Research ways and means are being studied of bringing one direction of this two-way communication, the output from machine to man, into the form of speech. Speaking machines will have many useful applications in the future. Much time can be saved, for example, by using a spoken set of instructions instead of a printed one when setting up complex equipment. In situations, too, where eyes and hands are already occupied and actions have to be performed on the basis of new information, that information can best be presented in the form of spoken messages (as for instance in the CARIN system^[1]).

Ir J. P. van Hemert is with the Institute for Perception Research, Eindhoven.

Many systems with speech output depend on the use of complete messages spoken-in beforehand. The system would be more flexible and offer wider prospects if the spoken messages could be built up from smaller units, in the same way as written language. The units constitute the basic elements of what might be called a 'speech alphabet', by analogy with the alphabet of written language. A system with such a collection of basic elements would bring an unlimited vocabulary within reach.

Before going into our choice of the basic elements for the speech alphabet, we shall first take a look at natural speech and indicate possible divisions of the speech signal. Speech is a succession of audible air-

^[1] CARIN is an automobile navigation and information system comprising a Compact Disc player, a position-determining system and a computer. It uses geographical information stored in digital form on a Compact Disc to the CD-I (Compact Disc Interactive) standard. The system takes this geographical information and the destination entered by the driver, and maps out a route. It then uses the position-determining system to guide the driver to his destination. The directions are given as spoken instructions (at present they are complete sentences spoken-in beforehand, but they will probably be derived from synthetic speech in future). CARIN will be the subject of a forthcoming article in this journal.

pressure differences that are generated by the human speech-production system (vocal cords, mouth cavity, nasal cavity and pharynx). Fig. 1 shows an example of a speech sound recording. The mouth and nasal cavities and pharynx (the vocal tract) can be regarded as a resonant cavity that is excited by the periodic vibration of the vocal cords into the voiced elements of speech (e.g. the vowel [a:]) or by turbulences that occur at constrictions in the vocal tract (e.g. the sibilant [s]). This excitation can be characterized by three parameters: the amplitude, which is connected with the loudness of the signal, a parameter that indicates whether the sound is voiced or unvoiced [2] and the source frequency at which the vocal cords vibrate. The resonant cavity can be characterized by the centre frequencies and the bandwidths of the formants, i.e. the resonance peaks in the transfer function. The source frequency is found in the fine structure of the spectrum, the formants in the envelope. The spectrum of speech up to a frequency of 5 kHz usually contains five formants, so that the resonant cavity can be described by ten parameters [2], the five centre frequencies and the corresponding bandwidths.

In principle, synthetic speech can be generated by specifying the behaviour of the thirteen parameters as a function of time. But this requires a large number of complicated rules, some of them as yet unknown. In practice, this approach is difficult and the speech quality is poor.

Another possible way of generating synthetic speech consists in joining together ('concatenating') fragments of natural speech. In its simplest form this amounts to the auditory presentation of *complete messages* spoken-in beforehand, as mentioned before.

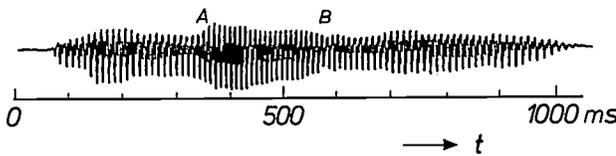


Fig. 1. The speech signal, showing the pressure changes as a function of time, for the nonsense word 'nenaanene' spoken by the test subject. At A and B, for example, the speech signal changes, where there is a transition between two sounds.

More scope is then offered by concatenating *separate words* to form complete sentences. But if all conceivable speech utterances are to be possible with such a system, an enormous list of words will be necessary. We can limit the volume to some extent by using *syllables* instead of words as our basic units. But even then, we still require a quite substantial 'library' of basic units. The library can be reduced to the mini-

mum if we take as our units the smallest segments of speech that still carry differences in meaning, the *phonemes* [3]. To maintain the fluent character of speech, however, it is then necessary to generate the transition from each phoneme to the next by sets of rules that are difficult to establish.

We have no need of such rules if we take *diphones* as our basic units. A diphone is a brief segment of speech that contains the last part of one phoneme, the transition, and the first part of the next phoneme. (Fig. 2 shows the phonemes and diphones contained in the word 'spoon'.) The number of basic units is now

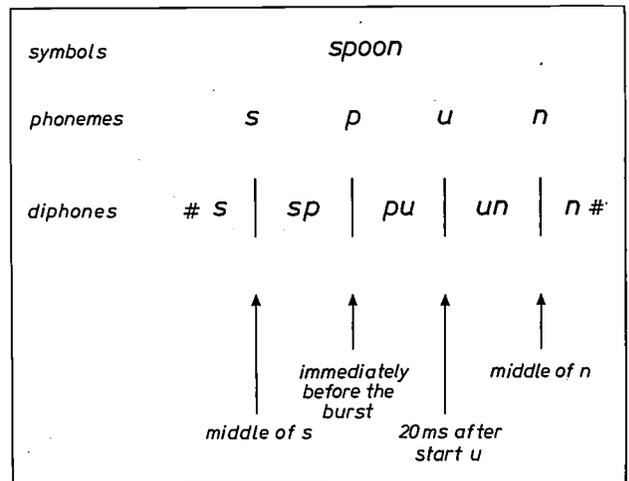


Fig. 2. An example of the location of the diphone boundaries in the word 'spoon'. After the phoneme boundaries have been determined, the diphone boundaries are fixed by rules. The first diphone is the transition from silence [#] to [s], the next one from [s] to [p], and so on.

larger, but we do not need to add any rules for the transitions. The Dutch language for example contains some forty phonemes. A library of diphones thus consists of some 1600 basic units. For each of these diphones a parameter representation must be available in the speech alphabet. The individual diphones of this library can be joined together to produce words or sentences.

A diphone library can be obtained in the following way. Tape recordings are made of a large number of words, which are spoken by a selected speaker. The list of words is compiled in such a way that it contains all the possible combinations of speech sounds. Four-syllable nonsense words may often be used, such as 'nenaanene' [nəna:nənə] and 'pəpaapepe' [pəpa:pəpə], in which the second syllable has to provide the diphones. The procedure is as follows. The words are digitized at a sampling rate of 10 kHz, and then the parameters mentioned above are determined and stored as a continuous series in a memory. This is used for generating (resynthesizing) a new speech signal. The phoneme boundaries are determined from the

parameter values and by listening to the resynthesis. These boundaries are thus established non-automatically. The diphone boundaries are then calculated from the phoneme boundaries, using a set of rules. The rules for positioning the diphone boundaries are different for each phoneme. In vowels the diphone boundary is located at a fixed time after vowel onset. The reason for this is that a vowel in natural speech is not affected, or only slightly, by the consonant immediately preceding it, but it is affected by the consonant immediately following it. The duration of a vowel is thus determined by the phoneme following the vowel; see for example fig. 2, where the boundary between the [pu] and the [un] diphone is located 20 ms after the start of the [u]. In plosives [b], [d], [p], [t] and [k] the diphone boundary is located immediately before the burst and in all other consonants it is positioned in the middle of the phoneme. These rules are used for determining the diphone boundaries. The parameters between the boundaries are stored in the diphone library.

This non-automatic method of building up a diphone library is tedious and time-consuming. A Dutch diphone library has been prepared in this way for one speaker (B.A.G. Elsendoorn)^[4]. The quality of the synthetic speech produced by concatenating diphones from this library is fairly good. There is a need for diphone libraries for more languages and more speakers for each language. (Although there are considerable differences between the diphone libraries for different languages, the principle for the concatenation of diphones can be used for various languages). The procedure described above has to be repeated, however, for each language and for each speaker. This is why efforts have been made to find an automatic method of preparing diphone libraries from recordings of spoken words. Besides the time argument there is a second reason for researching automatic segmentation. An automatic method may be expected to yield a final result that is more easily reproducible and more consistent than that of a non-automatic method.

In the rest of this article we shall consider a number of possible methods of automatic segmentation into phonemes, from which the diphones are then calculated from a set of rules, and it will be shown that a combination of two such methods gives a satisfactory result.

Automatic segmentation into phonemes

Segmentation methods described in the literature can be classified into two groups. On the one hand there are the methods that look for recognizable features in the speech sound, such as the boundaries be-

tween more or less steady-state parts of the speech signal. In the methods of this kind no use is made of the fact that the word to be segmented is known. There is therefore no phonetic identification, that is to say no attempt to link a phoneme with the associated segment of the signal. On the other hand there are methods that make the division on the basis of the correspondence between the speech utterance and a known reference pattern of each sound in the speech utterance. Both methods have their advantages and disadvantages. When a method of the first kind is used the boundaries between the segments are relatively accurately defined, but the segments delimited by two boundaries are not provided with any phonetic identification. A method of the second kind does not in practice appear to meet the requirements imposed by a speech-synthesis system for the accuracy of the boundaries^[5]. We shall first consider the method of the first kind that we are using. It will then be shown how satisfactory results can be achieved by combining this method with a method in which use is made of reference patterns.

Segmentation into steady-state segments

In general the behaviour of the acoustic properties of speech sound as a function of time is rather unpredictable. Nevertheless, it is possible to identify 'steady-state'^[6] segments of the speech signal, which is usually portrayed in terms of the pressure changes as a function of time. Fig. 1 shows the speech signal of the nonsense word 'nenaanene'. At the points *A* and *B*, for instance, transitions can be seen between two steady-state segments. The transitions are also visible in the spectral composition of the signal at the different times.

[2] The sounds are either 'voiced' or 'unvoiced', depending on the vibration of the vocal cords during articulation. All the vowels and the consonants [m], [n], [l], [b], [g] and [d] are voiced, and the plosives [p], [t] and [k] and the fricatives [s] and [f] are unvoiced. See also: J. 't Hart, S. G. Nootboom, L. L. M. Vogten and L. F. Willems, Manipulation of speech sounds, Philips Tech. Rev. 40, 134-145, 1982.

[3] In phonetic script phonemes are written as follows: short vowels with the associated letter (e.g. [æ] as in 'mat', [ɔ] as in 'rot' etc.), long vowels with the associated letter plus a colon (e.g. [a:] as in 'father', [ɔ:] as in lawn, etc.), consonants are also represented by the associated letter (e.g. [p] as in 'put', [f] as in 'find' etc.). The *a* as in 'ago' is indicated by [ə], the 'ur' as in 'turn' by [ɜ:], [#] indicates silence and [ʔ] the glottal stop. More information about phonetic representation will be found in most bilingual dictionaries.

[4] B. A. G. Elsendoorn and R. J. H. Deliege, A speech synthesis system using diphones: a portable communication aid for the speech impaired, Proc. Eur. Conf. on Technology & Communication Impairment, London 1985.

[5] L. R. Rabiner, A. E. Rosenberg, J. G. Wilpon and T. M. Zampini, A bootstrapping training technique for obtaining demisyllable reference patterns, J. Acoust. Soc. Am. 71, 1588-1595, 1982.

[6] These parts of the signal are not truly steady-state. The shape of the curve changes, but not so much as the relative change between the segments.

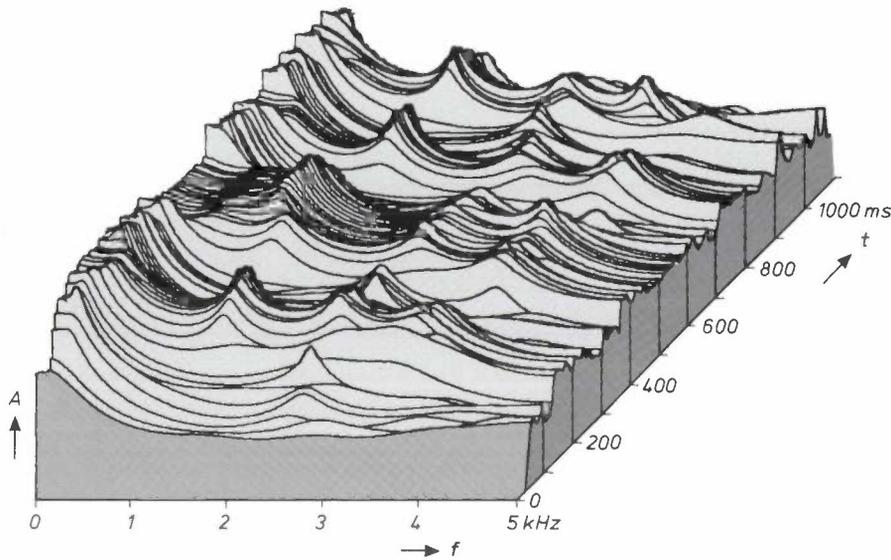


Fig. 3. The spectra (the amplitude A as a function of frequency f) as determined every 10 ms from the speech signal of the word 'nenaanene'. In most spectra five maxima can be recognized. These are the formants, and their heights and positions vary with time. These signals give the basis for dividing the utterance into steady-state parts.

Fig. 3 shows a number of 'pictures' of the amplitude spectrum at intervals of 10 ms. The formants, the peaks in the spectra, are subject to changes. The centre frequency and the bandwidth vary with time. With the aid of such a representation, amplitude spectra at intervals of 10 ms, the utterance can be divided into steady-state segments. Two spectra belong to the same steady-state segment of the utterance when they resemble each other to some extent. We must therefore indicate a measure of the similarity between the spectra. To do this we proceed as follows. We determine the extent of the similarity between two spectra by calculating the correlation between them. This correlation $c_{i,j}$ is defined as:

$$c_{i,j} = \frac{(s_i \cdot s_j)}{\sqrt{(s_i \cdot s_i)(s_j \cdot s_j)}} \quad (1)$$

where $(s_i \cdot s_j) = \int_0^{5 \text{ kHz}} W(f) s_i(f) s_j(f) df$, the scalar product of the spectra s_i and s_j , where f is the frequency and $W(f)$ the spectral weighting function that approximates to the sensitivity of the human ear. This spectral weighting function is given by:

$$W(f) = \begin{cases} 0 & \text{if } 0 < f \leq 200 \text{ Hz} \\ 1 & \text{if } 200 < f \leq 1000 \text{ Hz} \\ \frac{1000}{f} & \text{if } f > 1000 \text{ Hz.} \end{cases}$$

If two spectra are identical, the correlation $c_{i,j}$ is equal to 1. But as the similarity of the spectra diminishes, the more $c_{i,j}$ will differ from 1. The correlation $c_{i,j}$ between the i th spectrum and the ten neighbours on both sides ($i - 10 \leq j \leq i + 10$) is shown schemati-

cally in fig. 4. The curve has a maximum 1 at $j = i$. After choosing a threshold value c_t we can identify a segment belonging to spectrum i as the series of successive spectra (from i_b to i_e) whose correlation with spectrum i is greater than the threshold value c_t . The boundary between two segments is established in the following way. For all segments we determine the 'centre of mass' m_i , which is defined as:

$$m_i = \frac{\sum_{j=i_b}^{i_e} j(c_{i,j} - c_t)}{\sum_{j=i_b}^{i_e} (c_{i,j} - c_t)} \quad (2)$$

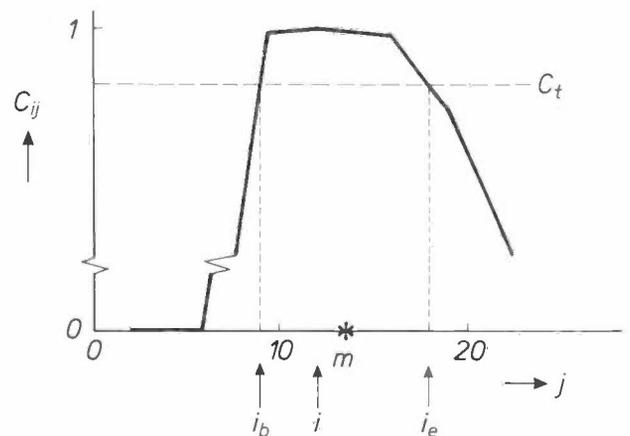


Fig. 4. Diagram representing the correlation $c_{i,j}$, the correlation between the spectra of the i th ($= 12$) spectrum and the j th spectrum (j goes from 2 to 22). The i th and the j th spectra belong to the same steady-state segment of the speech utterance if the correlation between the corresponding spectra is greater than c_t . In this schematic example the spectra $j = 9 (= i_b)$ to $j = 18 (= i_e)$ belong to the same steady-state segment. The 'centre of mass' m_i (see text) is also indicated in the figure.

The distance d_i from spectrum i to the centre of mass m_i of a segment is used for establishing the boundaries: $d_i = m_i - i$. Exactly how this is done can be seen by referring to the example given in fig. 5. For our nonsense word 'nenaanene' we have a representation of the speech signal, the correlations $c_{i,j}$ between the spectra, and the distance function d_i . For all i , the curve $c_{i,j}$ gives the correlation between spectrum i and the other spectra j in this speech utterance. From $i = 14$ to $i = 24$ the correlation curve changes only slightly. Near $i = 25$ we find a transition, and then the maximum of $c_{i,j}$ moves to larger j -values. The centre of mass of the curve is indicated by a short vertical line. The short horizontal lines indicate the distance d_i from spectrum i to the centre of mass m_i of the correlation curve. This distance is also plotted at the bottom of the figure. The zero crossings of the distance function from plus to minus are all close to the centre of a steady-state segment. The zero crossings from minus to plus are close to a transition, and are taken as the boundaries between two steady-state segments.

The n th boundary determined by this method is called $g_{1,n}$. We shall return to this later.

Phoneme boundaries have been determined in this way for a number of nonsense words. A comparison of these results with results obtained non-automatically shows that the phoneme boundaries have been fairly accurately determined. It can happen, however, that some boundaries are not found or that a phoneme is divided incorrectly into more than one segment. The number of incorrectly determined boundaries is so small, however, that this method can be used in an automatic segmentation program. For our real purpose, the segmentation into diphones, however, we do have to know which segments and phonemes correspond. In the following subsections we shall see how segments are given their phonetic identification. We shall also consider which boundaries have to be ignored or included to obtain the correct number of phonemes in an utterance.

Segmentation with the aid of reference patterns

In the division into steady-state segments as discussed in the previous subsection, we did not establish which steady-state segments and phonemes correspond. We do however know how many phonemes there are in the word to be segmented, and which ones they are. We can use this information to identify the steady-state segments. This can be done, for instance, by connecting the successive steady-state segments with the phonemes in the sequence occurring in the word to be segmented. This does not always work, however. There are some phonemes, the diphthongs (e.g. [ei], [ou], [ai]) and the plosives (e.g. [b], [d], [t]), that consist of two more or less steady-state sections. Furthermore, errors can be introduced by determining too many boundaries or too few. It is therefore necessary to ascertain whether the spectra of a steady-state part of the utterance do indeed correspond to the spectra of the intended phoneme. The procedure we follow for this is now described.

The word to be segmented is split into 'spectral states'. Most phonemes correspond to one spectral state; only the diphthongs and plosives correspond to two spectral states. We characterize each spectral state by a spectrum that we obtain from a series of spectra corresponding to the phoneme. We shall call these spectra the reference spectra. Next we compare the spectra of the speech utterance (we shall call these the test spectra) with the reference spectra of all the states that occur in the speech utterance. We do this in the same way as described in the previous subsection, by determining the correlation between the reference spectra and the test spectra. Let us illustrate this by an example. The word 'nenaanene' mentioned earlier

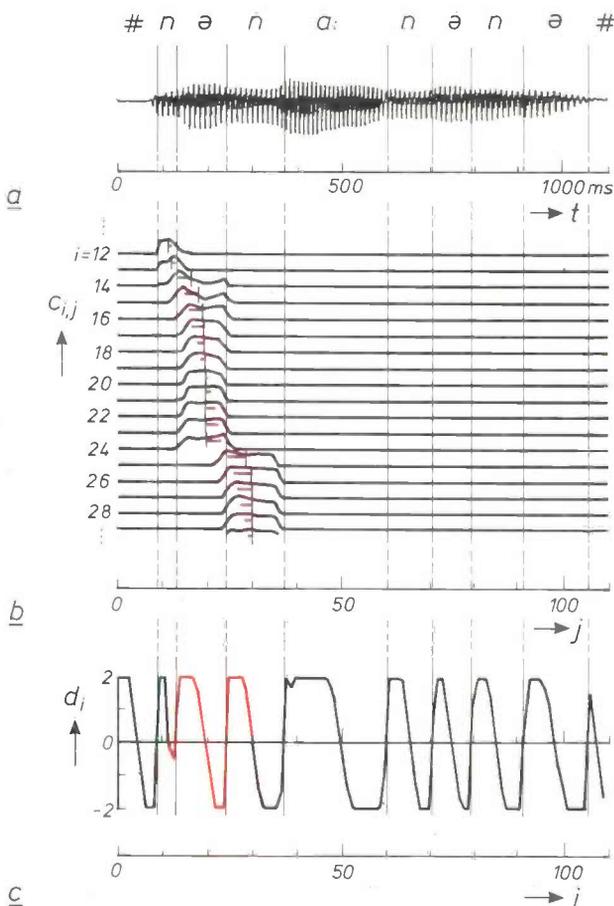


Fig. 5. The speech signal (a), the correlation $c_{i,j}$ for a part of the speech signal (b) and the distance function d_i (c) derived from the correlation. The lines indicating the phoneme boundaries derived from the zero crossings (from - to +) of the distance function are extended to the speech signal.

contains ten spectral states: the silence before the word [#], the phonemes [n], [ə], [n], [a:], [n], [ə], [n], [ə] and the silence after the word [#].

Fig. 6 shows the reference spectrum n for each of these spectral states and the correlation $c_{i,n}$ of all the test spectra i with the different reference spectra. The correlation is high if the test spectra closely resemble the reference spectra. The correlation of all the test

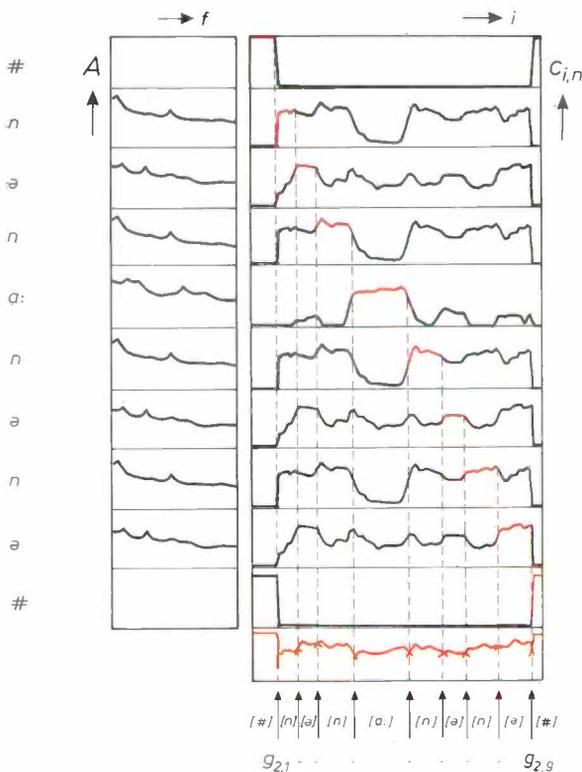


Fig. 6. The spectral states, the associated reference spectra n and the correlation $c_{i,j}$ between the reference spectra and the test spectra in the word 'nanaanene'. The arrows below indicate the boundaries $g_{2,n}$, as determined subject to the condition that the total sum of the correlations between the boundaries must be as large as possible. A is the amplitude, i the test spectrum and f the frequency.

spectra i with a reference spectrum n shows — depending on the number of times the spectral state occurs — a number of regions where the correlation is high. The correlation curve of the reference spectrum of the [n], for example, has four such regions, one for each occurrence of the [n] in the word. In principle, this spectral information can be used as the basis for segmenting an utterance. The spectral variability of speech is so great, however, that the correlation between the test spectra and a reference spectrum can always be lower than the correlation with other reference spectra, which would give the corresponding spectral state a zero duration. To prevent this, restrictions are imposed on the duration of each spectral state. The minimum and maximum durations of each spectral state are estimated by measuring the duration of a number of 'realizations' of that state.

We can now formulate a rule for the segmentation. The word to be segmented contains N spectral states and is available in the form of I test spectra. The last test spectrum associated with the n th spectral state we call g_n . The boundaries g_n for $n = 1$ to $(N - 1)$ are the free parameters that have to be determined. The boundary g_0 is set at the start of the speech utterance and the boundary g_N is set at the end: $g_0 = 0$ and $g_N = I$. The free parameters have to be determined in such a way that the resemblance between reference and test spectra is as close as possible, subject to the condition that the duration of each spectral state must lie between the associated minimum and maximum values (\min_n and \max_n). This requirement is formulated as:

$$\max_{g_n} \sum_{n=1}^N \sum_{i=g_{n-1}+1}^{g_n} c_{i,n}, \quad (4)$$

with the condition $\min_n \leq g_n - g_{n-1} \leq \max_n$. This maximization can be performed by 'dynamic programming' [7]. The results of these calculations can be found in the last graph of fig. 6. This shows parts of the correlation curves between the boundaries calculated as described above. The spectral state is also shown. In general, the inaccuracy in the determination of the boundaries by this method is too great for our purpose (see also fig. 10). In the next section we shall see that this method can however be used for generating the phonetic identification for the segments with steady-state properties, as obtained by the method described in the previous section.

Combination of the two methods of segmentation

Table I shows the advantages and disadvantages of the two automatic segmentation methods discussed earlier. It can be seen that the methods are more or less complementary. It is therefore logical to combine the two methods in such a way as to retain the advantages of each and avoid the disadvantages. This can be done by using the determination of the segment boundaries of the first method (which is accurate) and the phonetic identification of the second method. For

Table I. The advantages and disadvantages of the two automatic segmentation programs.

	Advantages	Disadvantages
First method	accurate determination of boundaries	too few boundaries or too many phonetic identification not known
Second method	correct number of boundaries phonetic identification known	boundaries not accurately determined

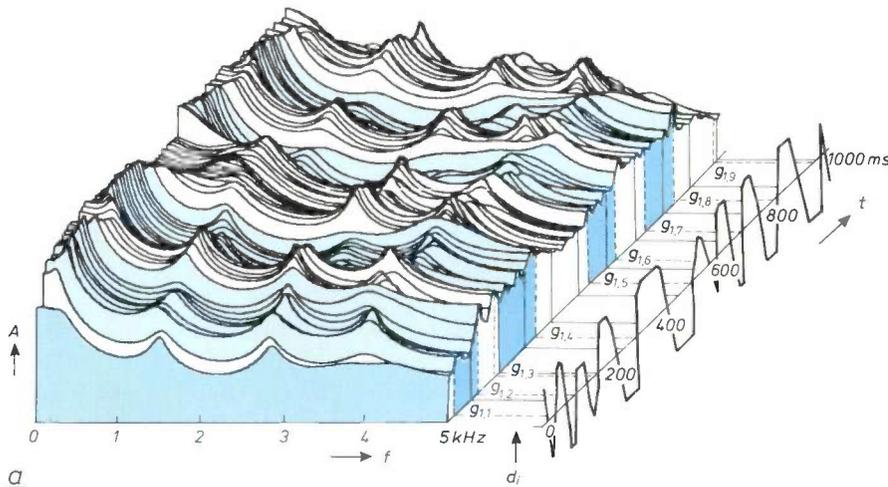
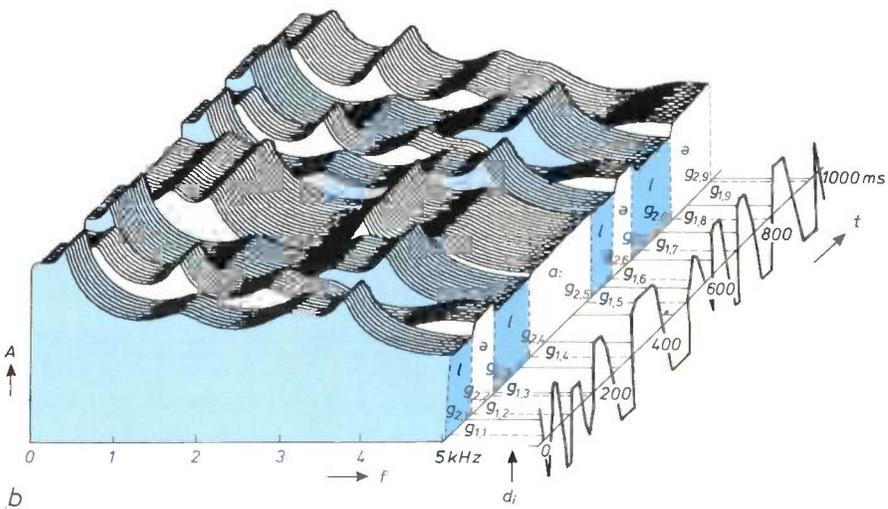


Fig. 7 a) The spectra of the word 'lelaalele' (the amplitude A as a function of frequency f at successive times). Segmentation into steady-state parts is done with the aid of the distance function d_i (equations 1, 2 and 3). The distance function d_i and the boundaries $g_{1,1}$ to $g_{1,9}$ derived from it are shown along the time axis. **b)** The boundaries $g_{2,1}$ to $g_{2,9}$ are found with the second automatic segmentation method. The diagram shows, between two boundaries, the reference spectra with which the spectra give the maximum correlation at the associated times (equation 4). Simply placing these reference spectra in sequence gives a picture (and also a sound) that is different from that of (a). The distance function d_i and the boundaries $g_{1,1}$ to $g_{1,9}$ derived from it are also indicated. The boundaries determined by the two methods do not correspond with one another everywhere.



this we need an algorithm that combines the boundaries determined by both methods. This implies that any extra boundaries found by the first method are omitted, and that any boundaries not found are inserted.

Fig. 7 shows the results of both methods once again. Fig. 7a gives the spectra of the nonsense word 'lelaalele' [lələ:lə]. The distance function d_i is shown along the time axis, with the boundaries $g_{1,n}$ derived from it (first method). The boundaries determined by the second method ($g_{2,n}$) can be seen in fig. 7b. This figure only gives the reference spectra that we used for determining these boundaries. Note that there is a considerable difference between the actual spectra (fig. 7a) and the steady-state 'caricature' of them obtained by the second method (fig. 7b). The distance function d_i and the boundaries $g_{1,n}$ derived from it are again shown along the time axis. At some places the

boundaries $g_{1,n}$ and $g_{2,n}$ are different. The boundaries closest to one another, $g_{1,n}$ and $g_{2,n}$, are linked. The final location of the boundaries is derived from the first method and the phonetic identification from the second method. A diagrammatic representation of this procedure is given in fig. 8. The link between two boundaries is represented as a dashed line. The linking of the boundaries can be continued until a number of boundaries remain that cannot be linked without the dashed lines crossing. The dashed lines must not

[7] 'Dynamic programming' is a mathematical procedure for solving an optimization problem. The solution is reached in different stages, in each of which a number of decisions can be made. The method can be efficiently programmed. More information can be found in books such as:
 O. I. Elgerd, Control systems theory, McGraw-Hill, New York 1967;
 R. E. Bellman, Dynamic programming, Univ. Press, Princeton, NJ, 1957.

cross because the sequence in which the segments have to be found is fixed. The boundaries $g_{1,n}$ that cannot be linked are omitted. The boundaries $g_{2,n}$ that cannot be linked are retained. The final location of these boundaries is derived from the second method. In the other cases, in which the boundaries can be linked, the final location of the boundaries is derived from the first method.

Fig. 9 shows an example of the segmentation obtained with the two methods combined, showing the boundaries ultimately established in the speech signal.

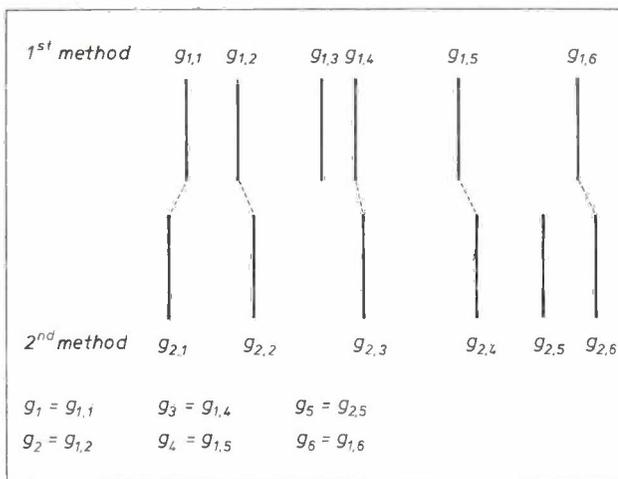


Fig. 8. A stylized example of the algorithm that combines the boundaries found by the two automatic segmentation methods. The upper vertical lines indicate the boundaries found by the segmentation into steady-state segments ($g_{1,1}$ to $g_{1,6}$) and the lower vertical lines give the results of the reference method ($g_{2,1}$ to $g_{2,6}$). The boundaries are linked by dashed lines in such a way that the total length of all the dashed lines has the minimum value. The dashed lines must not cross. The final results for the phoneme boundaries are also shown.

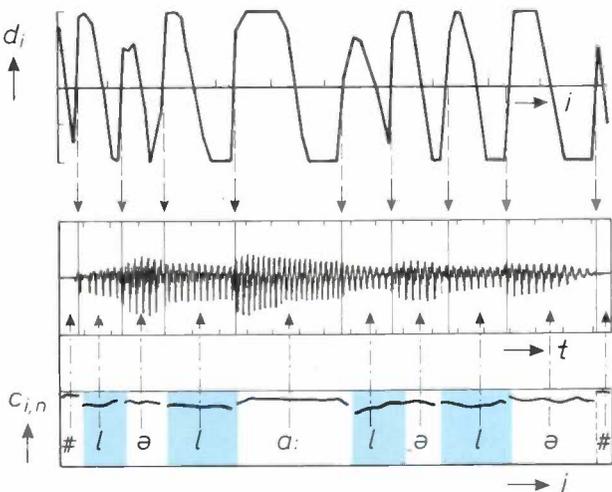


Fig. 9. The distance function d_i , the speech signal and the correlation $c_{i,n}$ between the spectral states $\{ \#, [l], [ə], [l], [a:], [l], [ə], [l], [ə], [\#] \}$ and the spectra for the word 'lelaalele'. The phoneme boundaries finally determined are shown in the speech signal.

Synthetic speech built up from diphones

The methods discussed above have been programmed on a VAX 11/780 computer. These programs have been used to make diphone libraries for Dutch, German and British English. The diphones are stored in parametrized form in these libraries. All possible speech utterances can be made audible by retrieving the appropriate diphones from the library, concatenating and synthesizing them. An example will be given to illustrate the final result. Fig. 10 shows the behaviour of the parameters as a function of time for the simple sentence 'The sun rises.' The parameter values of the diphones are shown in sequence. To make the synthetic speech more natural the source frequency is added by means of a set of rules [8]. In this example the stress is on the word 'sun'. At the boundaries between the diphones (vertical dashed lines) there are occasional discontinuities. These are not usually audible in the synthesis made by means of this parametrization. Before considering the result of an evaluation made by subjects who listened to a number of such synthesized utterances, we shall first look at a comparison of phoneme boundaries obtained by automatic and non-automatic means.

The phoneme boundaries of a number of words were determined by the usual non-automatic method, the second method and the two automatic methods combined. The results of these determinations are compared in fig. 11. The figure shows the percentage of phoneme boundaries for which the difference between the automatically and non-automatically determined boundaries is less than or equal to a particular multiple of 10 ms (the time interval between two spectra). Assuming that the non-automatically determined boundaries are somewhere near what we need for our purpose, we find that the results obtained by combining the two automatic methods give a definite improvement. If the required accuracy is set at 30 ms, 96% of all phoneme boundaries are correctly determined.

The agreement between automatically and non-automatically determined phoneme boundaries does not in itself show that the library of diphones obtained automatically is useful in practice for generating synthetic speech. Information about this can be obtained by asking subjects to assess the intelligibility and quality of the synthetic speech. A Dutch-language version of a diphone library has been used for this perceptual evaluation. Constant-pitch syntheses were therefore made of 320 words of one or two syllables from both the 'automatic' and the 'non-automatic' diphone libraries. Syntheses of both types were played in random order to seven subjects, who were asked

[8] J. 't Hart and R. Collier, Integrating different levels of intonation analysis, J. Phonet. 3, 235-255, 1975.

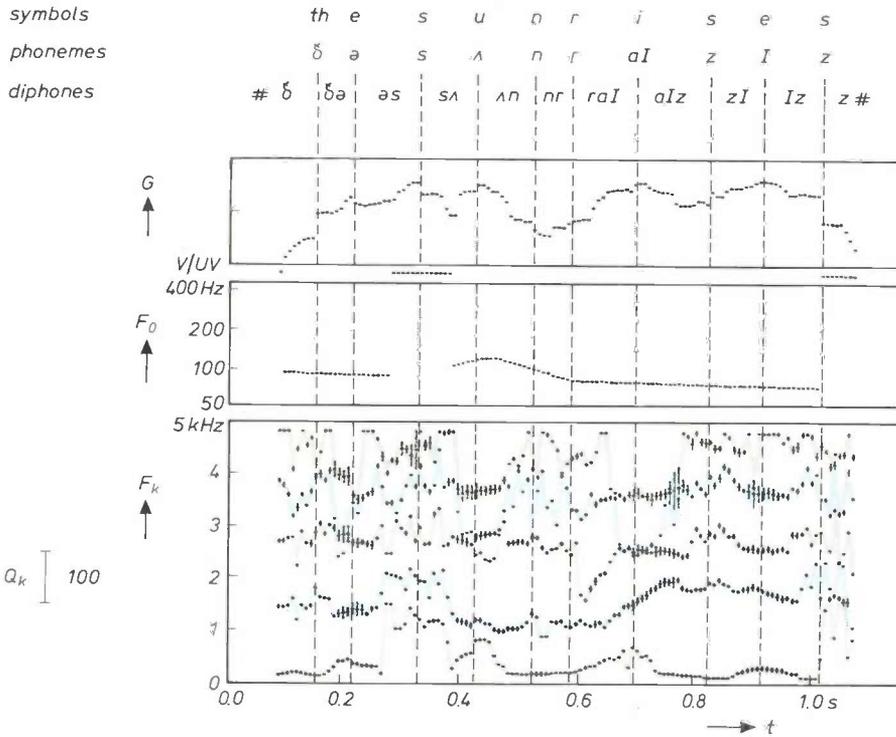


Fig. 10. Example of a concatenation of diphone parametrizations from the 'automatic' diphone library. From top to bottom: the sentence 'The sun rises.', the corresponding phonemes, the corresponding diphones, the mean amplitude G of the speech signal, the value of V/UV — unvoiced is shown black, voiced white — the source frequency F_0 in the voiced parts of the sentence, the centre frequency F_k and the Q (quality factor) Q_k ($k = 1$ to 5) of the five formants that characterize the resonant cavity. The centre frequency is indicated by a horizontal dash. The Q , which is inversely proportional to the bandwidth, is indicated by a vertical dash through the corresponding value of F_k . The value of the source frequency does not depend on the concatenation of the diphone parametrizations but has been added. The vertical dashed lines indicate the boundaries between the successive diphones.

which word they had heard. The subjects could also evaluate the quality of the synthetic speech by awarding points. The scores are presented in Table II. It can be seen that the 'automatic' diphone library is certain-

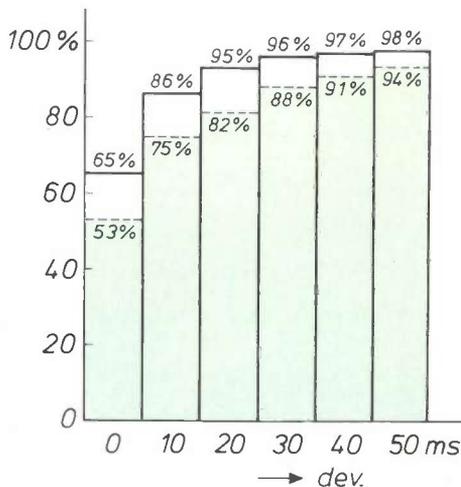


Fig. 11. The percentage of 'automatic' phoneme boundaries that deviate by less than an indicated multiple of 10 ms from the 'non-automatic' phoneme boundaries. The dashed lines with the corresponding percentages indicate the results of the second method, the continuous lines give the results of the combination of the two methods. In total 979 phonemes boundaries were determined in 90 words.

ly no worse than the 'non-automatic' one. Although this does not prove that the automatic segmentation always gives better results, it does show that 'automatic' diphone libraries can compete well with 'non-automatic' ones. In view of the other advantages, automatic segmentation is preferable to non-automatic segmentation. The fact that words are misunderstood in an intelligibility test is not solely attributable to erroneous segmentation of speech into diphones, of

Table II. Results of the perceptual evaluation. Altogether 320 words were synthesized from both the 'automatic' and the 'non-automatic' libraries. Seven test subjects were asked to identify the spoken word. The subjects were also able to give their evaluation of speech quality.

	Diphone set	
	'automatic'	'non-automatic'
Percentage of words wrongly understood by everyone	2%	3%
Percentage of words wrongly understood by at least one person	29%	35%
Total percentage of wrong scores (100% corresponds to 7 × 320 words)	9.5%	12.5%

course. Errors of analysis, the way in which the syntheses are made, the nature of the equipment, discontinuities and other factors all come into the picture.

The evaluation of the quality of the speech by different subjects is difficult to bring to a common denominator because their judgement is subjective. It can however be concluded that six of the seven subjects considered the quality of the syntheses made with the 'automatic' diphone library to be better than that of the syntheses made with the 'non-automatic' diphone library.

One aim of the investigation, to shorten the time it takes to prepare a diphone library from which synthetic speech of good quality can be produced, has been achieved. The segmentation of the speech utterances into diphones, which is a part of the process, can now be done in a fraction of the time that a 'non-automatic' method would have required. Synthetic speech built up from diphones obtained in this way will in future be subjected to further perceptual evaluation. An effort will also be made to improve various

details of the method. The strategy described can be used to build up a diphone library in less time than before. The availability of high-quality diphone libraries for more speakers and more languages could in future lead to promising applications, such as multilingual talking dictionaries and dialogue systems.

Summary. Systems in which synthetic speech is generated often use a library of sound transitions (diphones) as units for forming words in the same way as the letters of the alphabet are used for writing. The libraries are built up by extracting these units from spoken text. This can be done automatically by a combination of two automatic methods. In the first method, boundaries can be placed in words at positions where one steady state changes into another one. In the second method the sounds in a word are compared with reference patterns of particular sounds, and the phonetic identification is assigned to the segments determined by the first method. The second method is also used for correction if too many boundaries or too few are found. The boundaries obtained in this way are not very different from the results obtained by non-automatic methods. With these methods diphone libraries can be built up in a fraction of the time that would be required with a non-automatic method. Evaluation by test subjects of synthetic speech built up from 'automatic' diphone libraries is certainly no worse than that of synthetic speech built up from 'non-automatic' diphone libraries.

1937

THEN AND NOW

1987

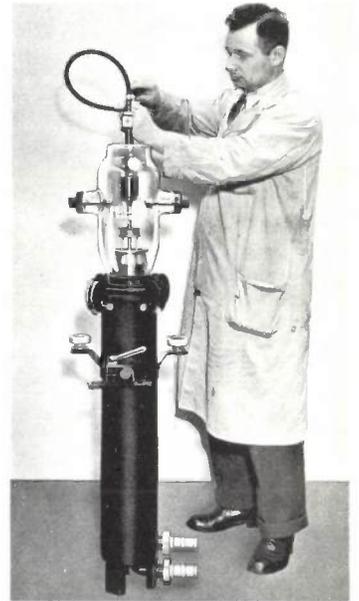
Transmitting valves

Fifty years ago the TA 20/250 transmitting triode was the largest transmitting valve in the Philips production range; see the photograph at the upper right [1]. This valve was 1.40 m high, and could supply a continuous-wave output power of 65 kW or a peak pulse power of 260 kW. The electrode system in this type of transmitting valve took the form of concentric cylinders, with the cathode at the centre and the anode at the outside. The heat dissipation was greatest at the anode, which was therefore force-cooled by a current of water.

The TA 20/250 transmitting triode was used in long-wave broadcast transmitters, since the valve could only be used at relatively low frequencies (1 MHz and below) because there was glass insulation between the electrode connections and the electrode system was relatively long. (The length of the electrodes should be small compared with the wavelength of the transmitted signal.) The water-cooled anode could dissipate a maximum of 100 W/cm².

Triodes that will provide very high powers are still produced today, and their construction remains very much the same. The largest type in the present range has a continuous-wave output power of 500 kW and a peak pulse power of 2 MW. The height of this valve is 0.80 m; see the valve at the back in the lower photograph. The heat transfer has been greatly improved, so that the length of the electrode system need only be less than half that of the TA 20/250 transmitting triode of 1937. This and the ceramic insulation with very low dielectric losses make these valves useful at considerably higher frequencies, so that they can be used — in tetrode form — in short-wave broadcast transmitters at frequencies up to 30 MHz. As triodes they are also used for heating plasmas in controlled nuclear fusion. The frequency of about 100 MHz then used is equal to the cyclotron frequency of the ions in the plasma.

The two valves at the centre of the lower photograph are air-cooled, the other three are water-cooled. Because the water cooling is very efficient the anode can dissipate a maximum of about 1000 W/cm². The high heat transfer has been achieved by the use of special grooves on the outside of the anode.



[1] From Philips Technical Review, April 1937.

Black-cobalt coating for solar collectors

B. Vitt

The efficiency of a solar collector can be improved considerably by using an absorber plate with a black-cobalt coating which is spectrally selective. Until recently, making suitable coatings was mainly an empirical process. The article below describes an investigation into the fundamental properties of such coatings. The results can be used successfully to optimize coatings for high-temperature applications.

Introduction

In a solar collector, an absorber plate with a black coating is used to convert the solar energy into thermal energy. The heat collected is transferred to a circulating medium (e.g. water or oil) and can be used for a variety of purposes, e.g. for space heating or in industrial processes. Various measures can be taken to compensate for the increase in thermal radiation losses at the plate when it is operating above ambient temperature. These include a concentration of the solar energy by reflectors and lenses, and a reduction of convection and conduction losses by evacuation of the collector tube. The efficiency of the energy conversion can be improved further by using a spectrally selective black coating which has a high absorptance for solar radiation, i.e. at wavelengths shorter than about 2.5 μm , and a small emissivity for thermal radiation at longer wavelengths.

Several materials have been investigated for this application^[1]. As a part of the Philips activities on evacuated solar collectors^[2] an efficient black-cobalt coating was developed in our laboratories. An electroplating process was used to successfully apply this to the inside walls of silver-coated glass tubes in an experimental collector^{[2][3]}. Later on, a quantity-production process was developed at the Philips Plastics and Metalware Factories in Eindhoven for deposition on copper-coated steel plates. These were used in the

commercial Philips evacuated tubular collector, in which the heat gained by the absorber is extracted efficiently via a heat pipe^[4].

The efficiency and life of a solar collector depend strongly on the selective properties and stability of the black coating. In the vacuum environment of an evacuated tube there is no stability problem as long as the collector is designed for 'load' temperatures no higher than 150 °C. However, in view of the good selective properties of cobalt coatings a modification of the collector design seems to be possible in a way that allows efficient production of process heat in the high-temperature range of 150 to 250 °C^[5]. In this case stagnation temperatures of more than 350 °C can arise, requiring extended knowledge of the long-term stability of the coating under such conditions.

During our investigations cobalt coatings were deposited on various steel substrates. The samples were characterized chemically and optically, with emphasis on the stability properties at elevated temperatures in various atmospheres^[6]. In addition, the effect of the substrate on the optical properties and stability was investigated. The results have been related to the collector performance. It was found that coatings well-suited for high-temperature applications can be obtained on smooth substrates, after sufficient annealing and outgassing in an inert atmosphere.

In this article we first deal with the structure and chemical composition of the coatings. Next we shall

Dr B. Vitt is with Philips GmbH Forschungslaboratorium Aachen, Aachen, West Germany.

discuss their optical properties and their behaviour in non-inert atmospheres. Finally we shall show how the coating properties affect the collector performance.

Structure and composition

For the investigations on the structure and composition black-cobalt coatings were prepared on different steel substrates previously coated with a thin infrared reflecting copper film. The substrates are designated as smooth, medium or rough, corresponding to surface roughnesses of less than 0.05, 0.1 or $1\ \mu\text{m}$, respectively. The electroplating was carried out at room temperature with a solution containing 25 g/l of cobalt chloride (CoCl_2) and 25 g/l of potassium thiocyanate (KSCN), at a pH-value of 4.5. The current density was $1\ \text{A}/\text{dm}^2$ and the electrolytic thickness of the deposit, expressed as the charge transferred per dm^2 , was varied up to $500\ \text{C}/\text{dm}^2$.

The structures of the deposits were characterized by scanning electron microscopy (SEM) as well as by transmission electron microscopy (TEM) for samples which had been made sufficiently transparent to electrons. Like similar electrolytic coatings, the deposits exhibit a rather inhomogeneous and porous structure, depending on the deposition conditions and the roughness of the substrate. SEM micrographs of two coatings on substrates with different roughness are shown in *fig. 1*. The nucleation conditions on the smooth substrate obviously lead to a denser and more compact deposit.

Annealing at $350\ ^\circ\text{C}$ and $450\ ^\circ\text{C}$ introduces a significant change in the surface structure, and this is accompanied by a decrease in the coating thickness. This decrease amounts to 25% after the $450\ ^\circ\text{C}$ step and is related to a partial decomposition and to the observed emission of considerable amounts of H_2O and CO_2 . In addition to showing that the structure is inhomogeneous, TEM micrographs indicate that the material itself is built up inhomogeneously from a granular distribution of particles in a solid matrix of filling material. The particles seem to vary in size from 2 to 10 nm.

Coating material removed from the substrate was found to be ferromagnetic. This indicates the presence of a cobalt-rich phase, since common cobalt compounds such as oxides and sulphides are paramagnetic. This interpretation is supported by energy-dispersive X-ray fluorescence analysis during SEM and TEM, giving an average atomic cobalt/sulphur ratio of 10:1. Infrared reflection spectra reveal the presence of OH^- and CO_3^{2-} ions, pointing to a significant amount of $\text{Co}(\text{OH})_2$ and CoCO_3 . X-ray diffraction does not show any detectable lines, so the deposits are either amorphous or polycrystalline with very small



Fig. 1. SEM micrographs showing that a cobalt coating of $30\ \text{C}/\text{dm}^2$ on a smooth substrate (a) is much denser than that on a rough substrate (b).

particle sizes. However, after inert annealing at $250\ ^\circ\text{C}$ diffraction lines of crystalline CoO are obtained. These lines are intensified on annealing in air. For some samples that were annealed at $350\ ^\circ\text{C}$ the formation of Co_6S_6 was detected.

These observations were correlated with a quantitative chemical analysis of coatings of maximum thickness of about $10\ \mu\text{m}$. The percentages by weight obtained for the elements Co, S, O, C and H were con-

- [1] See for example O. P. Agnihotri and B. K. Gupta, *Solar selective surfaces*, Wiley, New York 1981.
- [2] H. Hörster (ed.), *Wege zum energiesparenden Wohnhaus*, Philips Fachbuchverlag, Hamburg 1980; H. O. Jungk and H. Scholz, *Deutsche Offenlegungsschrift P 25567162*; H. Bloem, J. C. de Grijs and R. L. C. de Vaan, *An evacuated tubular solar collector incorporating a heat pipe*, *Philips Tech. Rev.* **40**, 181-191, 1982.
- [3] K. R. Schreitmüller and K. Vanoli, *Proc. 8th Biennial Congr. Int. Sol. Energy Soc.*, Perth 1983, pp. 816-824.
- [4] J. C. de Grijs and R. J. H. Jetten, *Proc. 8th Biennial Congr. Int. Sol. Energy Soc.*, Perth 1983, pp. 1071-1076.
- [5] B. Vitt, *Solar systems with highly efficient collectors*, *Proc. Summer School on Solar Energy, Igls (Austria) 1985*, pp. 55-61.
- [6] B. Vitt, *Characterization of a solar selective black coating*, *Sol. Energy Mater.* **13**, 323-350, 1986.

Table I. Composition of a black-cobalt coating before and after inert annealing at 450 °C. For simplicity the cobalt-sulphur compound is shown as CoS.

Element or compound	Mol. %	
	Before annealing	After annealing
Co	62.7	63
CoS	9.6	10
Co(OH) ₂	18.8	—
CoCO ₃	8.3	—
CoO	—	27
Residue	0.6	—

verted into molecular ratios, as given in *Table I*. In view of the TEM results and the ferromagnetic properties it is unlikely that the composition given corresponds to a homogeneous alloy. From energy-dispersive X-ray fluorescence analysis it can be shown that the spatial distributions of Co and S are highly correlated. An inhomogeneous model may therefore be more reasonable, which accounts for a granular distribution of small Co₁₀S particles in a matrix mainly consisting of the 'filling materials' Co(OH)₂ and CoCO₃.

As would be expected, Co(OH)₂ and CoCO₃ decompose on annealing in an inert atmosphere, leading to the formation of CoO and the emission of H₂O or CO₂. The emission of some CO, H₂, H₂S and SO₂ may also become significant at higher temperatures. However, the decay of the sulphur content in the samples is very limited below 450 °C. These results show that the coatings should be carefully outgassed at 450 °C before using them in vacuum collector tubes. *Table I* gives the approximate composition after such a heat treatment. Vacuum annealing at even higher temperatures should then result in a porous layer consisting of cobalt metal plus some mol. % of CoO.

Optical properties

An important quantity in characterizing an absorber plate for solar collectors is the 'solar absorptance' α , which is defined as the fraction of the incident solar power absorbed by the plate. Another important quantity is the thermal emissivity ϵ , which is the radiated power density divided by the power density radiated by an ideal black surface at the same temperature. To obtain solar collectors with a high efficiency, α should be as large as possible and ϵ as small as possible.

These properties were evaluated for coatings of different thickness by measuring the reflection spectra in the infrared wavelength range (2.5 - 50 μm) and in the solar wavelength range (0.3 - 2.5 μm). *Fig. 2* shows the spectra for coatings on smooth steel substrates on

which 1 μm of copper had been electroplated from a cyanide solution. The spectra exhibit pronounced interference structures, which shift towards the infrared with increasing thickness. Analysis of the maxima and minima reveals that the geometrical thickness is proportional to the plating time or the electrolytic thickness. This holds at least as far as the 500 C/dm² sample, for which a geometrical thickness of $10 \pm 1 \mu\text{m}$ was determined by microscopy. The spectra also demonstrate that the requirements of a high absorptance in the solar wavelength range and of a high infrared reflectance (low thermal emissivity) are only met by coatings with an electrolytic thickness not exceeding about 30 C/dm².

The spectra of *fig. 2* can be described by using unique wavelength-dependent functions for both the refractive index n and the absorption index k . Computer simulations were found to fit the measured spectra very well, except for short wavelengths ($\lambda \leq 0.8 \mu\text{m}$) and for the vibrational OH⁻ and CO₃²⁻ peaks [7]. The appropriate optical functions are shown in *fig. 3*. Whereas n increases monotonically with λ , k decreases from the high value at short wavelengths. However, in the infrared there is an increase in k , which is almost the same as the increase in n , indicating some metallic behaviour. This increase in k implies that thick coatings are not transparent in the infrared (see also *fig. 2*). In spite of the combination with a reflecting metal film, a thick black-cobalt film is therefore more like a black-body radiator than a selective radiator.

The n and k functions also provide a fairly satisfactory description of the directional thermal emissivity ϵ_ϕ , obtained by radiometric measurements of the heat emitted at different values of the angle ϕ with respect to the normal to the surface. If ϵ_ϕ is expressed as a function of $\sin^2 \phi$, an expression for the hemispherical emissivity ϵ_h can easily be derived:

$$\epsilon_h = \int_0^1 \epsilon_\phi d(\sin^2 \phi).$$

Results for coatings at 127 °C are shown in *fig. 4*. It can be seen that the coating thickness is the most important parameter for the directional emissivity. For thicker coatings the shape of the curve for ϵ_ϕ differs increasingly from the shape for the uncoated copper surface. The ratio $\epsilon_h/\epsilon_\perp$ of the hemispherical emissivity to the emissivity along the normal first increases from the value of 1.3 for uncoated copper to a maximum of about 1.9 and then decreases to values close to 1. The increase in the high-angle emission with very thin coatings is due to the increased optical pathlength at higher angles. For coatings of 60 C/dm² or more the emission characteristics of a weakly absorbing dielectric [8] are obtained.

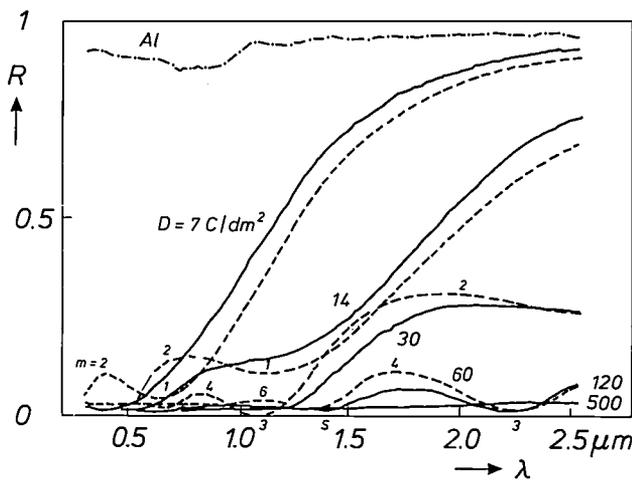
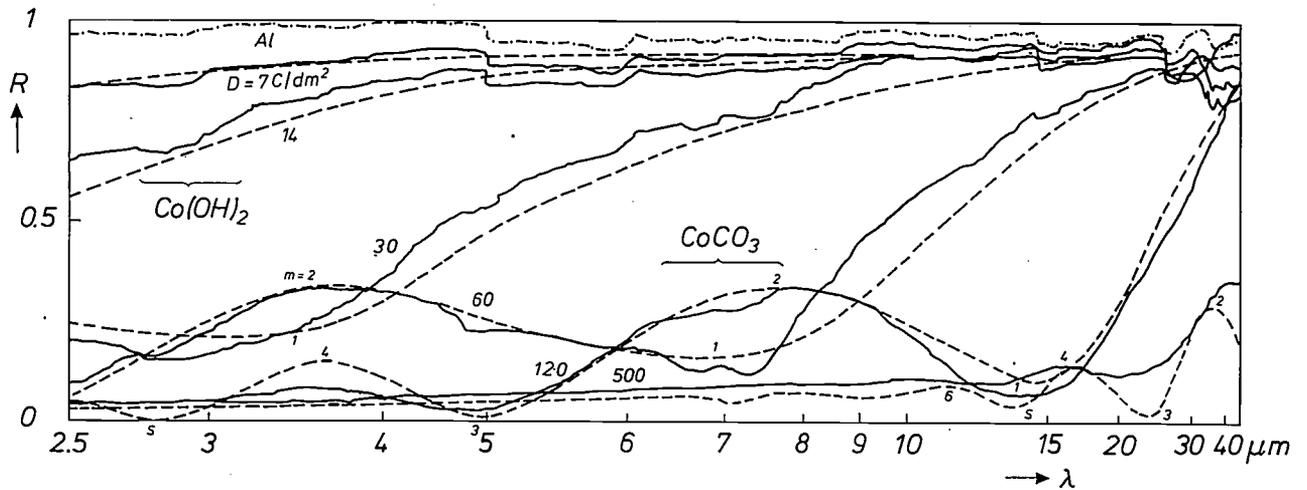


Fig. 2. Measured reflectance R plotted against wavelength λ in the infrared (a) and the solar wavelength range (b), for black-cobalt coatings with different electrolytic thickness D . For comparison the reflectance of an aluminium film (Al) is also given. The dashed curves refer to calculations made using wavelength-dependent functions for the refractive index and the absorption index, and assuming a proportional relation between electrolytic and geometrical thickness (500 C/dm^2 corresponds to 10 μm). The interference maxima and minima of order m are also indicated. On increasing thickness, the reflectance in the near infrared decreases strongly and the interference maxima and minima shift to longer wavelengths.

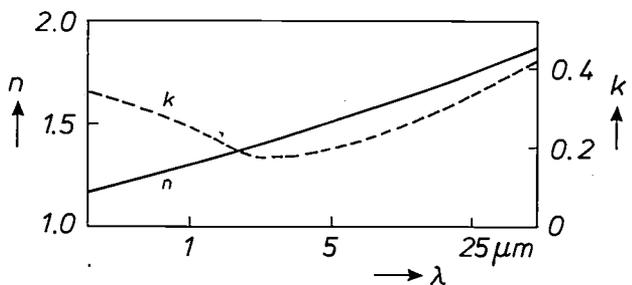


Fig. 3. Refractive index n and absorption index k of black-cobalt coatings as a function of the wavelength λ . Whereas n increases monotonically with λ , there is a minimum for k at about 2.5 μm .

Annealing in high vacuum results in two significant changes in the optical properties. After annealing at 350 °C the reflection spectra reveal a significant increase in k for wavelengths between 1 and 15 μm , to an almost constant value of about 0.35. This effect is accompanied by a reduction in thickness of several per cent. A subsequent treatment at 450 °C has no measurable effect on k , but a thickness reduction of 25% is now observed which remains constant even after annealing for more than 24 hours.

The effect of annealing on the reflectance in the solar wavelength range is shown in fig. 5 for a coating of 14 C/dm^2 . The mean value for the solar absorptance still exceeds 90% for a sample annealed at 450 °C. It is important to note that the infrared emissivity of this sample is essentially the same as for a non-annealed sample, since the increase in infrared absorptance is almost exactly compensated by the reduction of 25% in the coating thickness. If the coatings are assumed to be homogeneous with compositions as given in Table I, a decrease in thickness of 33% would be expected. However, the annealing is associated with changes in the porous structure, as has been demonstrated by SEM micrographs.

The considerable influence of the substrate roughness on the morphology of the coatings (fig. 1) has significant consequences for their optical properties. This is illustrated in fig. 6, which shows the directional emissivity for coatings of different thicknesses on various substrates. The emission characteristics of the uncoated substrates are found to be very similar, but

[7] See for example R. A. Nyquist and R. O. Kagel, Infrared spectra of inorganic compounds, Academic Press, New York 1971.
 [8] H. C. Hottel and A. F. Sarofim, Radiative transfer, McGraw-Hill, New York 1967.

in the case of the coated samples the substrate roughness has a pronounced effect, especially for coatings of 30 C/dm². In this case the hemispherical emissivity at 133 °C may vary by more than a factor of two. This can be mainly attributed to an enhanced porosity of the coating on rough substrates, which is associated with an increased thickness.

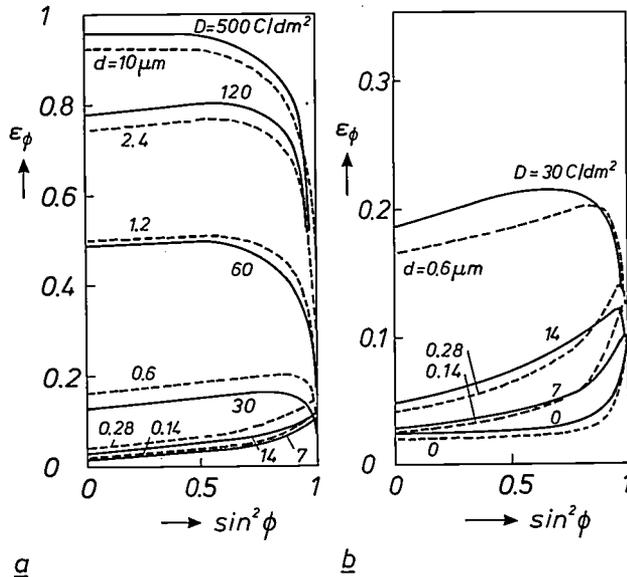


Fig. 4. Directional thermal emissivity ϵ_ϕ at 127 °C as a function of $\sin^2\phi$ for two series of black-cobalt coatings. The continuous curves were obtained from measurements on coatings with different electrolytic thickness D . The dashed curves refer to calculations for different values of the geometrical thickness d . For coatings of thickness up to 60 C/dm² the directional emissivity differs increasingly from that of an uncoated copper surface ($D = 0$ C/dm²). The thicker coatings have an emissivity that is characteristic of a weak absorber.

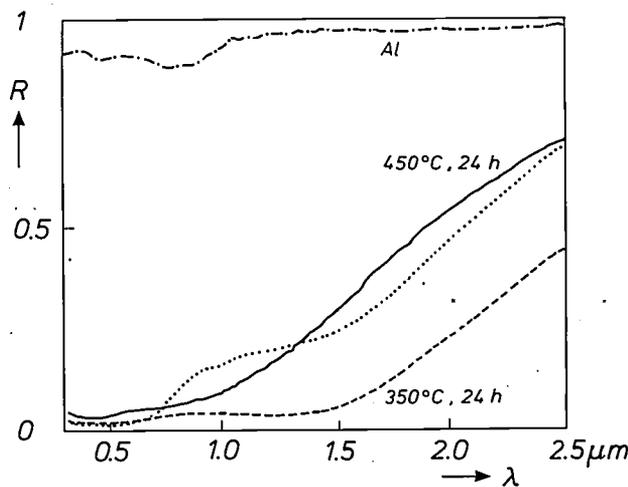


Fig. 5. Effect of vacuum annealing on the reflectance in the solar wavelength range, for a cobalt coating of 14 C/dm² on smooth steel coated with 5 μm of bright copper. Annealing for 24 h at 350 °C (dashed curve) results in an appreciable decrease in reflectance as compared with the non-annealed coating (dotted curve). After annealing for a further 24 h at 450 °C (continuous curve) the original near-infrared reflectance is almost completely restored.

A number of results are summarized in fig. 7, where the normal and hemispherical infrared emissivities are given as a function of the electrolytic thickness. Because of the broadband infrared absorption of the black cobalt, the coating thickness is the most important parameter for the emissivity. The same is true for the solar absorptance, which is also shown in fig. 7.

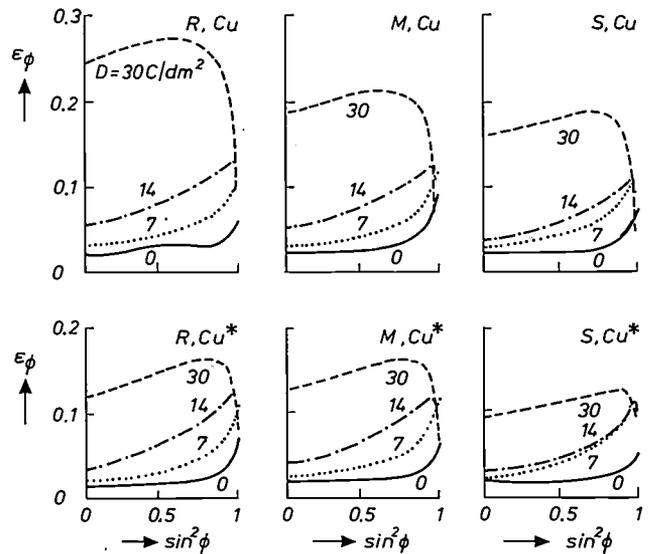


Fig. 6. Directional thermal emissivity ϵ_ϕ at 133 °C plotted against $\sin^2\phi$ for black-cobalt coatings of different electrolytic thickness D . The substrates used are rough (R), medium (M) or smooth (S) and they are coated with 1 μm of copper from a cyanide solution (Cu) or 5 μm of bright copper (Cu*). It can be seen that the coating thickness is the most important parameter for the directional emissivity.

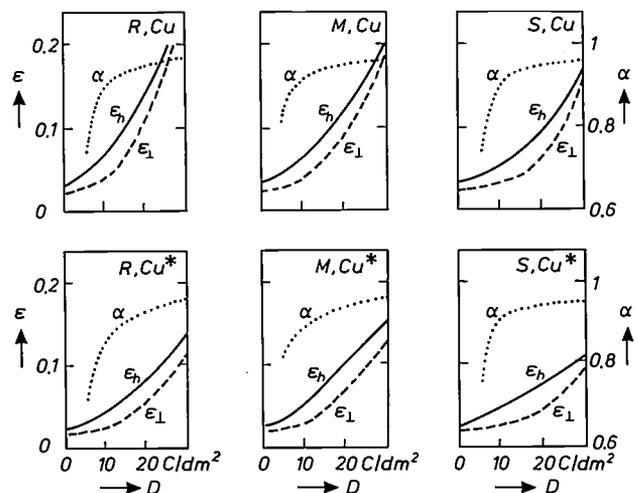


Fig. 7. Solar absorptance α (at room temperature) and the thermal emissivities ϵ_l and ϵ_h (both at 133 °C) as a function of the electrolytic thickness D for the samples of fig. 6. The properties of the substrate have a much smaller effect on α than on ϵ_l and ϵ_h . The best results are obtained on a smooth substrate with 5 μm of bright copper, giving $\alpha = 93\%$ and $\epsilon_h = 5\%$ for a cobalt coating of 14 C/dm².

However, whereas the emission properties depend significantly on the substrate properties, the solar absorptance is much less affected. The requirements of a high solar absorptance and a low infrared emissivity are best met with a coating thickness of 14 C/dm^2 on a smooth steel substrate coated with bright copper. In this case the solar absorptance is as high as 93% and the hemispherical emissivity at 133°C is as low as 5%.

Performance in a non-inert atmosphere

A thin black-cobalt coating changes considerably when exposed to air at elevated temperatures. For example, during a one-hour anneal in air at about 250°C its colour changes from black to a golden or brownish hue, with a marked reduction in the solar absorptance. SEM micrographs show that the coating is still present, and that its morphology is unchanged. The formation of a more transparent coating is due to the oxidation of the cobalt particles to CoO , as has been demonstrated by X-ray diffraction. Further exposure to air leads to the formation of the grey-black oxide Co_3O_4 ^[9] and an oxidation of the reflecting copper film, giving an increased infrared emissivity.

Coatings with reduced solar absorptance were also obtained after annealing under poor vacuum conditions ($\approx 1 \text{ Pa}$) at higher temperatures (350°C) and longer exposure times. The resulting reflectance in the solar wavelength range is shown in *fig. 8* for a coating of 14 C/dm^2 . For comparison, the spectra before annealing and before coating are also shown. Absorption data for CoO ^[10] were also used to calculate the spectrum for an oxidized coating. The resulting curve is close to the one measured after annealing, which confirms the above assumptions.

Annealing at 450°C in a poor vacuum ($\approx 1 \text{ Pa}$) may lead to severe degradation of the surface, related to the quality of the copper film. *Fig. 9* shows SEM micrographs of two surfaces with such a degradation. Both contain outgrowths which can be observed at different stages of the growth. X-ray diffraction and X-ray fluorescence analysis indicate that these outgrowths consist of crystalline $\gamma\text{-Fe}_2\text{O}_3$. This oxide has a broadband infrared absorption^[7] and the arrangement of the outgrowths also forms a resonance structure for the infrared wavelengths. This degradation can therefore be responsible for values of the emissivity which are between the value for black cobalt and a

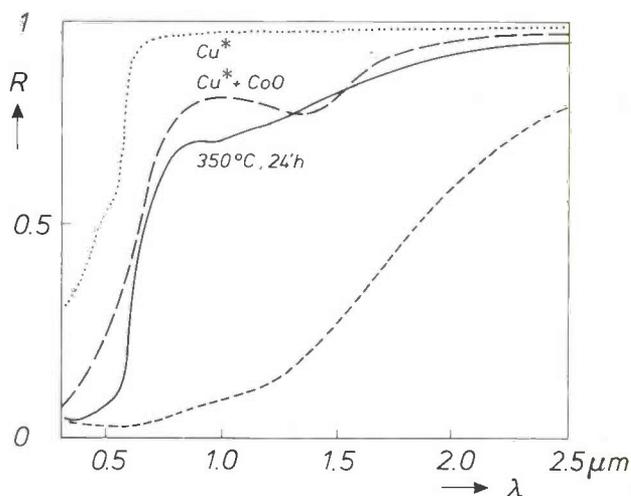


Fig. 8. Effect of annealing for 24 h at 350°C in poor vacuum ($\approx 1 \text{ Pa}$) on the reflectance of a black-cobalt coating of 14 C/dm^2 in the solar wavelength range. The dotted curve refers to the bright-copper reflector, the lower dashed curve and the continuous curve to the cobalt coating before and after annealing. The strong increase in reflectance can be simulated fairly well. The upper dashed curve was calculated for $0.26 \mu\text{m}$ of CoO on bright copper.

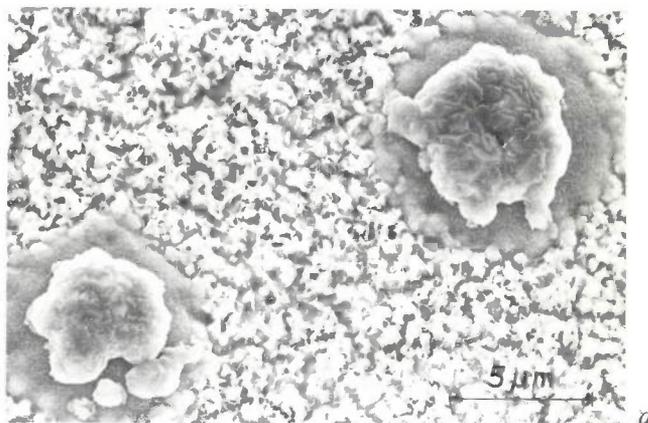


Fig. 9. SEM micrographs of two coatings, showing the effect of annealing for 24 h at 450°C in poor vacuum, with an insufficient coating of copper on the steel substrate. For a cobalt coating of 60 C/dm^2 on $1 \mu\text{m}$ of cyanidic copper (a) large outgrowths are formed. For a coating of 14 C/dm^2 on $0.5 \mu\text{m}$ of cyanidic copper (b) many outgrowths of different sizes are observed.

^[9] The formation and properties of this oxide have been described by several authors, including: T. Tanaka, *Jap. J. Appl. Phys.* **18**, 1043-1047, 1979; and C. Choudhury and H. K. Sehgal, *Sol. Energy* **28**, 25-31, 1982.

^[10] I. G. Austin, B. D. Clay and C. E. Turner, Optical absorption of small polarons in semiconducting NiO and CoO in the near and far infra-red, *J. Phys. C* **1**, 1418-1434, 1968.

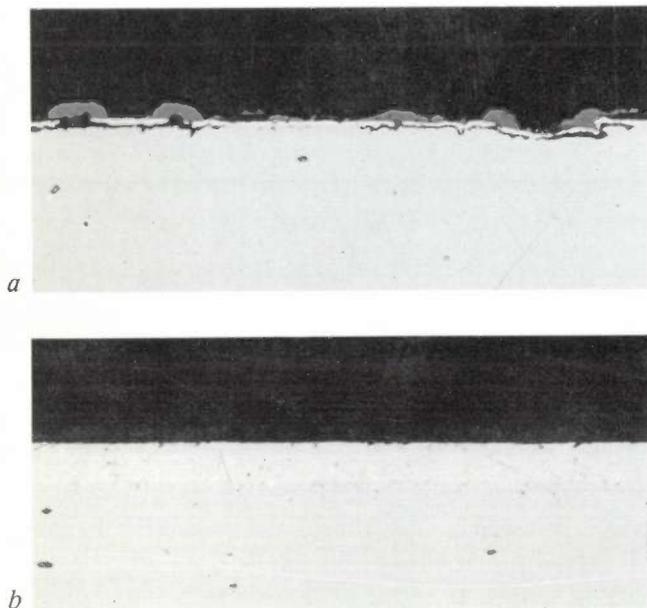


Fig. 10. Micrographs of polished cuts of two coated steel substrates with different copper films, after annealing under the same conditions. The degradation (a) found with the sample of fig. 9b does not occur when the substrate is coated with 5 μm of bright copper (b).

The infrared reflection spectra of some annealed copper films on steel substrates are shown in fig. 11. The thinner films exhibit a low reflection in the near infrared and characteristic Fe_2O_3 peaks at long wavelengths. The 5-μm film, on the other hand, gives a high infrared reflection which is not very different from that of a reflecting aluminium film.

Effect on collector performance

Some ten years of experience of the use of black-cobalt coatings in the Philips evacuated collector tubes^{[2]-[4]} has now been accumulated. In recent collector designs the maximum absorber temperature is about 250 °C. Further optimization may however allow increased stagnation temperatures, especially if an optimized coating is used. In such a case, the coating will have to withstand temperatures of up to 350 °C. Our present results indicate that an optimized coating, carefully annealed at 450 °C under inert conditions, has sufficient stability to maintain a high tube efficiency for more than fifteen years.

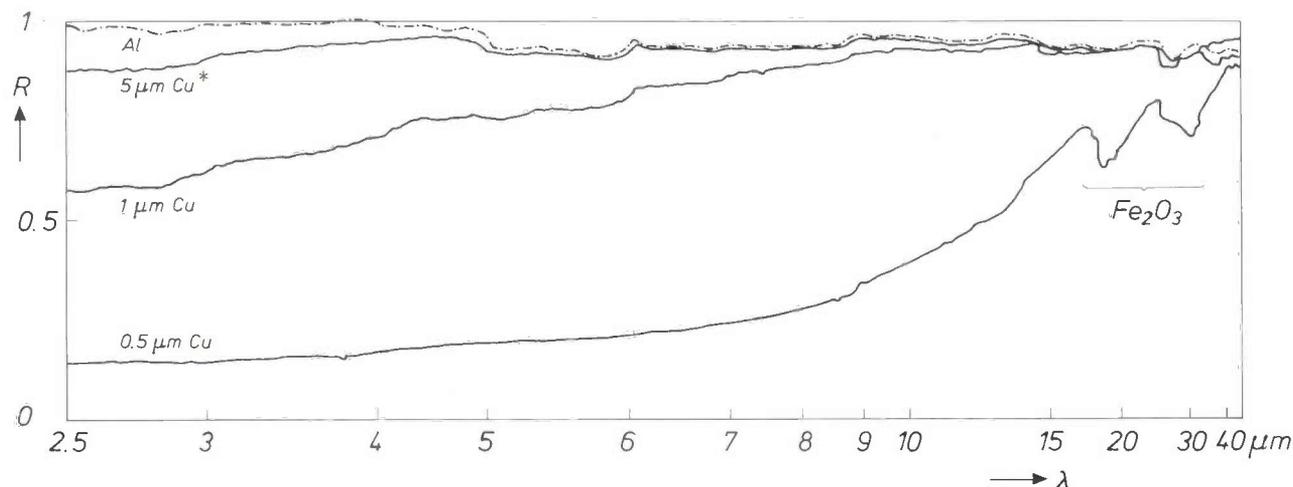


Fig. 11. Infrared reflection spectra of steel substrates coated with various copper films, after annealing for 24 h at 450 °C in poor vacuum. An insufficient coating results in too low an infrared reflectance and vibrational Fe_2O_3 peaks due to a partial oxidation of the steel. With 5 μm of bright copper the reflectance is very like that of a reflective aluminium film.

value close to 1, the emissivity of an ideal black-body radiator.

To explain this degradation, it is assumed that the steel of the substrate is attacked by oxygen or OH^- groups via the semipermeable copper film, which is 0.5 or 1 μm thick in the examples of fig. 9. The Fe_2O_3 formed by the oxidation is transported to the surface via tunnels in the copper film, as can be seen in the micrograph of fig. 10a. Chemical attack by oxygen or water and the resulting degradation can however be avoided by using a 5-μm film of bright copper; see fig. 10b.

The most important characteristics for the collector performance — the solar absorptance α and the thermal emissivity ε — both increase with coating thickness^[11]. The thickness dependence of the solar absorptance α and the hemispherical emissivity ε_h at 127 °C (400 K) for optimized coatings can be expressed as

$$\alpha = 0.975 - 0.815 \exp(-0.22 D),$$

$$\varepsilon_h = 0.022 + 0.0016 D + 4.5 \times 10^{-5} D^2,$$

where D is the electrolytic thickness in C/dm^2 . The increase in the thermal emissivity with temperature for

various thicknesses is illustrated in *fig. 12*. The calculated curves are also shown, made with the assumption that the dielectric properties of the coating are temperature-independent in this range. The emissivity increases linearly with temperature, and the slope γ (i.e. the temperature coefficient) increases linearly with the coating thickness:

$$\begin{aligned} \epsilon_h(D, T) &= \epsilon_h(D, 400 \text{ K}) + \gamma(D) \times (T - 400 \text{ K}), \\ \gamma(D) &= (1.52 D - 6.65) \times 10^{-5} \text{ K}^{-1}, \end{aligned}$$

for $D \geq 7 \text{ C/dm}^2$. These results can be explained by the shift of the Planckian radiation distribution to shorter wavelengths and the reduction in the near infrared reflectance with increasing thickness (*fig. 2*).

The results obtained were used for the evaluation of the collector performance. This was done for an improved version of the Philips VTR 361 solar collector tube [4], which has an aluminium reflector. The calculations were based on a validated model for evacuated tubular collectors [5].

The collector efficiency η is defined as the power delivered to the heat-transfer fluid (e.g. water or oil) in the collector system divided by the incident radiant power. The useful power delivered to the heat-transfer fluid is equal to the power Q_a absorbed by the plate less the thermal loss Q_1 , whereas the incident radiation power is given by the solar irradiance G_T (in W/m^2) multiplied by the aperture area A (in m^2) of the collector, so that:

$$\eta = (Q_a - Q_1)/AG_T.$$

The ratio Q_a/AG_T is referred to as the conversion efficiency η_0 , and Q_1/A is proportional to the overall heat-transfer coefficient U_1 (in $\text{W/m}^2\text{K}$) and to the average difference ΔT between the fluid temperature in the upper heat pipe of the collector [4] and the ambient temperature. The efficiency can therefore be expressed as

$$\eta = \eta_0 - F'U_1 \Delta T/G_T,$$

where $F' (\leq 1)$ is a dimensionless factor to take account of the heat resistance between absorber and heat-transfer fluid.

In *fig. 13* the calculated values of the conversion efficiency η_0 and the heat-loss coefficient $F'U_1$ are shown as a function of the coating thickness for various values of ΔT . The collector efficiency for any operating conditions can be obtained from these curves and the above equation.

The heat-loss coefficient $F'U_1$ is found to increase considerably with increasing coating thickness. Its temperature dependence is determined both by the T^4 radiation law and by the temperature dependence of the emissivity itself. For an adequate conversion

efficiency η_0 , a minimum electrolytic thickness of about 8 C/dm^2 is required. It should be noted that a performance with $\eta_0 \geq 0.70$ and $F'U_1 < 1 \text{ W/m}^2\text{K}$ is possible for low-temperature applications ($\Delta T \approx 50 \text{ K}$). In high-temperature applications ($\Delta T \approx 200 \text{ K}$) η_0 can

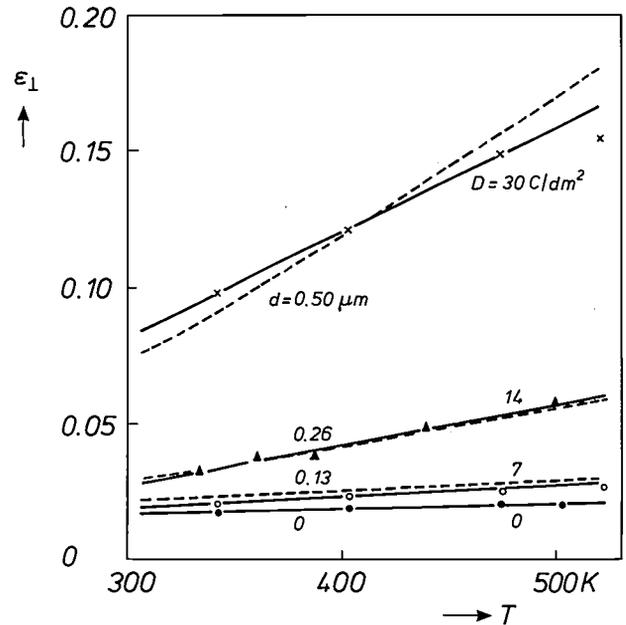


Fig. 12. Temperature dependence of the thermal emissivity ϵ_{\perp} for black-cobalt coatings with different electrolytic thicknesses D (continuous curves) and calculated for different geometrical thicknesses d (dashed curves). It is seen that ϵ_{\perp} depends linearly on the absolute temperature T , and that the slope increases for thicker coatings.

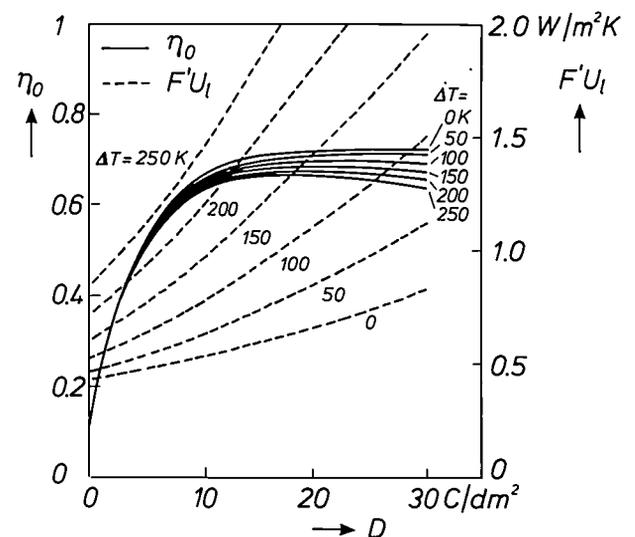


Fig. 13. Calculated performance of a solar collector (VTR 361 tubes with heat pipe and an aluminium reflector). The conversion efficiency η_0 and the heat-loss coefficient $F'U_1$ (see text) are plotted against the electrolytic thickness D of optimized coatings for different values of the average difference ΔT between the heat-transfer fluid temperature in the upper heat pipe and the ambient temperature.

[11] General information on the use of solar energy can be found in: J. A. Duffie and W. A. Beckman, *Solar engineering of thermal processes*, Wiley, New York 1980.

still be as high as 0.60 or more, with $F'U_1 \approx 1.20$ $\text{W/m}^2\text{K}$.

These performance results can be converted into annual efficiencies by a simple method^[5]. The values were calculated for climatic conditions characteristic of

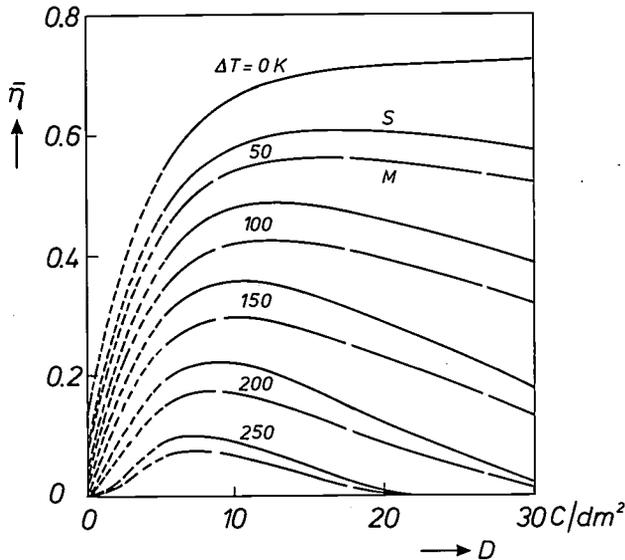


Fig. 14. Annual efficiency $\bar{\eta}$ of a Philips evacuated tubular collector as a function of the electrolytic coating thickness D , calculated for two different European climatic zones (S Southern Europe, M Central Europe) and for different values of the mean operating temperature difference ΔT . The optimum coating thickness can be derived for a given application with a known value of ΔT .

two European climatic zones; see fig. 14. Obviously, different coating thicknesses should be used: 15 C/dm^2 for low-temperature applications and 10 C/dm^2 for high-temperature applications. In the latter case annual efficiencies of more than 20% are possible in a sunny climate with permanent operation at $\Delta T = 200 \text{ K}$.

Valuable contributions to the work described here were made by our colleagues at the laboratories in Aachen and the Philips Plastics and Metalware Factories in Eindhoven.

Summary. In the evacuated receiver tubes of a solar collector an absorber plate with a spectrally selective black-cobalt coating is used for reducing the thermal radiation losses. The influence of the substrate roughness and annealing conditions on the structure and composition of the coatings greatly affects their high-temperature stability and thermo-optical properties. The most suitable coatings for high-temperature application are obtained on smooth substrates after vacuum-annealing at about 450°C . The optimum electrolytic coating thickness for the required optical selectivity is generally less than 30 C/dm^2 , depending on the intended temperature range in the collector. The stability of optimized coatings seems to be sufficient for use in solar collectors designed for high-temperature applications ($150\text{--}250^\circ\text{C}$).

Research on television glass

D. M. Krol and R. K. Janssen

The mass production of glass has long been a stronghold of empiricism, but science is now an increasing presence. Much of the initiative for this comes from the Central Glass Laboratory, set up in 1957, and from Philips Research Laboratories. Following on from an earlier study at Philips Research Laboratories on the 'fining' of glass, described previously in this journal, attempts have been made to obtain a scientifically based understanding of a wider field of the glass-formation processes, including both melting and fining. The glass compositions used in these studies were made as close as possible to those used in practice.

Introduction

The average home in the western world contains a television set. The glass for the picture tube is as likely as not made in a Philips pressed-glass works. Against this background and the fact that glass also forms an indispensable component for many other Philips products, it is not surprising that studies on glass are regularly reported in this journal^[1].

The glass for picture tubes is generally formed in a continuous process by subjecting a mix of raw materials to a treatment in which the temperature varies with time in a special predetermined pattern^[2]. The mix for this application consists of alkali metals and alkaline-earth metals combined as oxides and carbonates, SiO_2 , Al_2O_3 as feldspar and a number of special additives, and cullet. The temperature-time treatment takes place in furnaces constructed from separate, but closely fitting, blocks of refractory material. A diagram of a widely used version of such a furnace, which can contain 300 tonnes of glass, is shown in *fig. 1*; *fig. 2* is a view of the interior of such a furnace.

The mix is delivered via the 'batch hopper' *B* to the 'melting end' *M*. This is where the glass-formation reactions take place. The molten glass produced flows through the 'working end' *W* via feeders *F* to the glass-processing machines. After the glass has been formed in the 'melt' phase at about 1350 °C and has become completely fluid, the melt initially contains large quantities of gas bubbles; at this stage the composition of the melt is still highly inhomogeneous in other

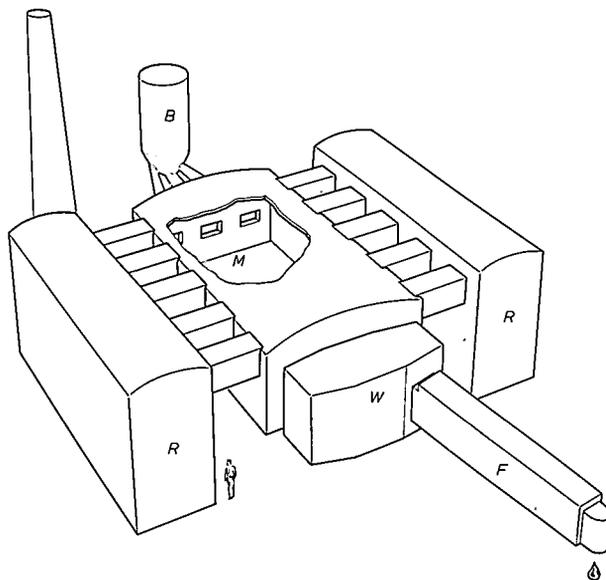


Fig. 1. Diagram of furnace for the production of glass for television picture-tube screens. *B* batch hopper. *M* 'melting end'. *W* 'working end'. *F* feeder. The air above the glass melt is heated by burners (not shown). The temperature is held at about 1350 °C in the melting end, and allowed to fall to about 1000 °C for shaping. The way in which the temperature is allowed to vary between these two final values has a considerable effect on the quality of the glass produced. The regenerators *R* preheat the air with heat derived in part from the flue gases. A furnace such as this can contain 300 tonnes of glass, with a production of 90 tonnes a day.

Dr D. M. Krol, formerly with Philips Research Laboratories, Eindhoven, is now with AT&T Bell Labs, Murray Hill, New Jersey, U.S.A.; Ing. R. K. Janssen is with Philips Research Laboratories, Eindhoven.

^[1] See for example the issue devoted entirely to glass: Philips Tech. Rev. 22, 281-341, 1960/61, and the recent article: A. Kats, Glass — outline of a development, Philips Tech. Rev. 42, 316-324, 1986.

^[2] See for example: G. E. Rindone, Glass Ind. 38, 489-528, 1957; J. Stanek, J. Non-Cryst. Solids 26, 158-178, 1977.

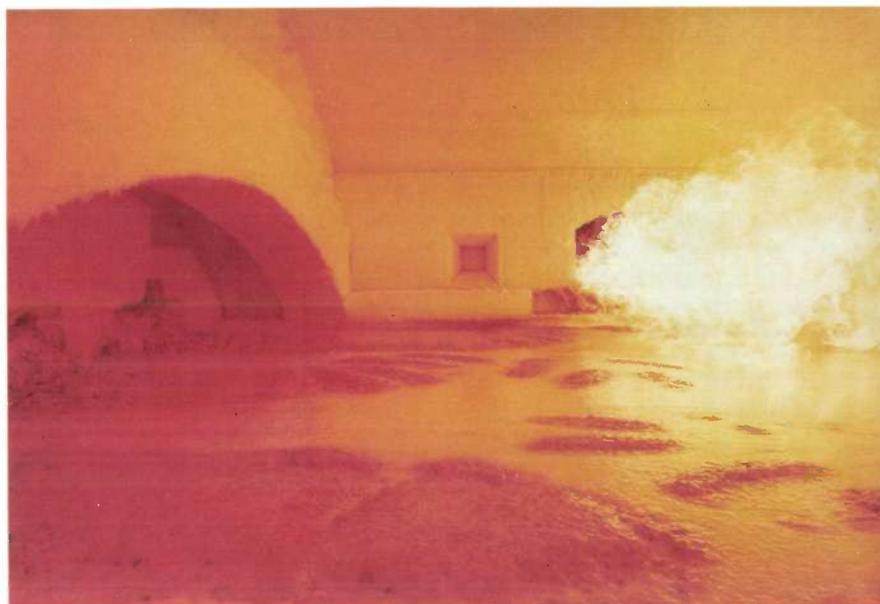


Fig. 2. Looking into a glass furnace at the Philips glassworks in Aachen, West Germany. The mix enters the furnace through the three openings on the left of the furnace, to form the 'sandbanks' further along in the 'river' of molten glass. The furnace is 7 m wide and 15 m long, and contains a glass melt about 1 m deep. Above is one of the burners that melt the mix at a temperature of about 1350 °C.

respects. As the melt continues to flow through the furnace, the bubbles disappear in a 'fining' process, and the melt becomes more homogeneous. Finally, the temperature decreases slowly to about 1000 °C, and the shaping processes can be started. The glass remains in the furnace for about 50 hours on average.

Among the special additives included with the raw materials, the fining agents are particularly important; they are often antimony compounds in combination with NaNO_3 . The cullet, which always accounts for a fairly large proportion of the ingredients, is added to reduce the process temperature. In an earlier study at Philips Research Laboratories it was shown that Raman spectrometry can be extremely useful for studying the fining of glass^[3], largely because the presence of both crystalline and liquid compounds can be established with this method. This study was performed on a glass system obtained by melting 70 mol% of SiO_2 and 30 mol% of K_2CO_3 ; only two or three different compounds are produced in the glass formation in this system.

In this article we describe new investigations at Philips Research Laboratories on the formation of glass, with particular relevance to glass for colour television screens. Since we wanted to start with a more realistic glass composition than in the previous study, which implied that the experiments would be more complex, we split the investigation into two parts: we investigated the melting stage and the fining stage separately. For the investigations on the melting stage we used a glass composition that is indeed closer to the

actual composition of the television glass: $15\text{Na}_2\text{CO}_3$ - 10BaCO_3 - 75SiO_2 , but is certainly not the same as the composition used in practice (see page 257). In the composition we used the amount of Na corresponds to that of the sum of all the alkali metals, the amount of Ba to that of the sum of all the alkaline-earth metals and the amount of SiO_2 to that of all the $\text{SiO}_2 + \text{Al}_2\text{O}_3$ in the actual composition.

The only other constituent of the mix that we included in the study was the NaNO_3 . We did not include the fining agent antimony, assuming that its action would only be observable in the fining stage. With these restrictions, we endeavoured to follow all the conversions by using laser-Raman spectrometry.

In the investigations on the fining stage we chose a type of glass corresponding so closely to the glass used in practice that it was no longer practical to use Raman spectrometry, because of the complexity of the spectra. In this study we have confined ourselves to the effect of the fining agent, antimony, with special attention to the effect of the $\text{Sb}^{3+}/\text{Sb}^{5+}$ ratio, which we considered relevant. We determined this ratio by potentiometric titration and emission spectrometry.

Investigation of the melting stage

Measurements on samples consisting entirely of raw material

In our first series of investigations we wanted to obtain a picture of the variety of phases, both crystalline and amorphous, that can arise in a relatively 'realistic'

glass composition during the melting as a function of temperature and time. To start with we therefore subjected samples consisting of the raw materials Na_2CO_3 , BaCO_3 and α -quartz (i.e. without special additives) to temperatures between 700 and 1400 °C for periods of 5 minutes to 16 hours, and then quenched them to room temperature. Fig. 3 gives a number of Raman spectra for the samples treated in this way^[4]. The sharp peaks correspond to crystalline phases, the broad bands to amorphous phases or — at higher temperatures — to liquid phases. The identification of

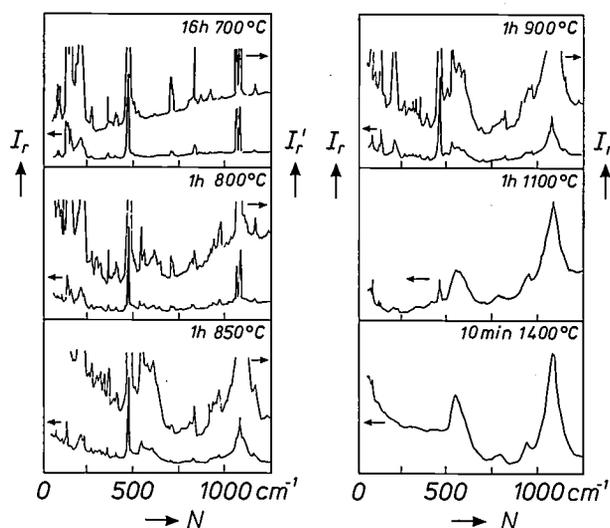


Fig. 3. Raman spectra of mixes consisting of 15 mol% of Na_2CO_3 , 10 mol% of BaCO_3 and 75 mol% of SiO_2 that have been heated for different times and at different temperatures (shown top right) and then quenched to room temperature. The intensity I_r is shown as a function of the wave number N in arbitrary units; at the four lowest temperatures it is also shown with a unit five times as large (I_r').

these peaks and bands was done largely with the aid of compounds we made specially for the purpose (Table I). Table II indicates the temperature range in which the various compounds can occur.

These Tables lead us to the following conclusions:

- Up to 800 to 900 °C the spectra chiefly originate from crystalline phases — produced in solid-state reactions. A significant amount of the raw material SiO_2 (about half) has already been used up in solid-state reactions at 800 °C.
- Above 1000 °C all the spectra are due to liquid phases, with the exception of α -quartz, which has been dispersed in the melt in crystalline form.
- As the temperature increases there is a 'condensation', in which silicates with an increasing content of silicon are formed. At 700-800 °C the silicates are mainly ortho- and pyrosilicates, at 850 and 900 °C they are mainly meta- and disilicates.

The influence of NaNO_3

From the results of a Raman-spectrometric investigation in which NaNO_3 was also involved^[5] (fig. 4a and b) it was found that the addition of NaNO_3 to the raw materials mentioned above led to a *more rapid* formation of the *same* reaction products in the temperature range around 700 °C, where the solid-state reactions are dominant; when the samples are heated at 1000 °C, with and without NaNO_3 , the spectra are identical.

We have been able to confirm the first result by a study in which the effect of heating for an hour at 700 °C with NaNO_3 was compared with the effect of heating for 16 hours at 700 °C without NaNO_3 (fig. 5). In spite of the much longer heating time fewer reaction products were formed in the sample without NaNO_3 .

Since the NaNO_3 is only found to have an effect in the temperature range where solid-state reactions are dominant, its action could possibly come about because it assists the formation of a liquid phase below 700 °C. In this liquid phase Na_2CO_3 and BaCO_3 in the dissolved state would be able to react more quickly with α -quartz than they would in the solid state.

Since no data are available about the phase diagram of the system Na_2CO_3 - NaNO_3 - BaCO_3 , we have tested this assumption about the action of NaNO_3 by heating a mix of these substances, in the correct proportions, to 700 °C. This mix did indeed go into the liquid phase at this temperature.

Investigation of the fining stage

In studying the very complex fining action the following three points should be borne in mind:

a) The glass melt contains impurities in the form of bubbles; in addition to nitrogen and water from the atmosphere the melt always contains large quantities of the CO_2 generated in the solid-state reactions. Since the purpose of the fining is to remove all the bubbles from the glass, it is important that the solid-state reactions should be fully completed before the fining mechanism comes into action.

b) The fining mechanism is brought into action by the conversion of Sb^{5+} from the added fining agent into Sb^{3+} with the formation of oxygen, which rises to the surface in the form of bubbles. This oxygen formation at the same time aids the growth of the other bubbles present in the glass, which also rise to the surface.

[3] See for example H. Verweij, Philips Tech. Rev. 40, 310-315, 1982.

[4] A more detailed account of this investigation is given in: D. M. Krol and R. K. Janssen, J. Physique 43 (Colloque C9), C9/347-C9/350, 1982.

[5] A more detailed account of this investigation is given in: D. M. Krol and R. K. Janssen, Glastech. Ber. 56K, 1-6, 1983.

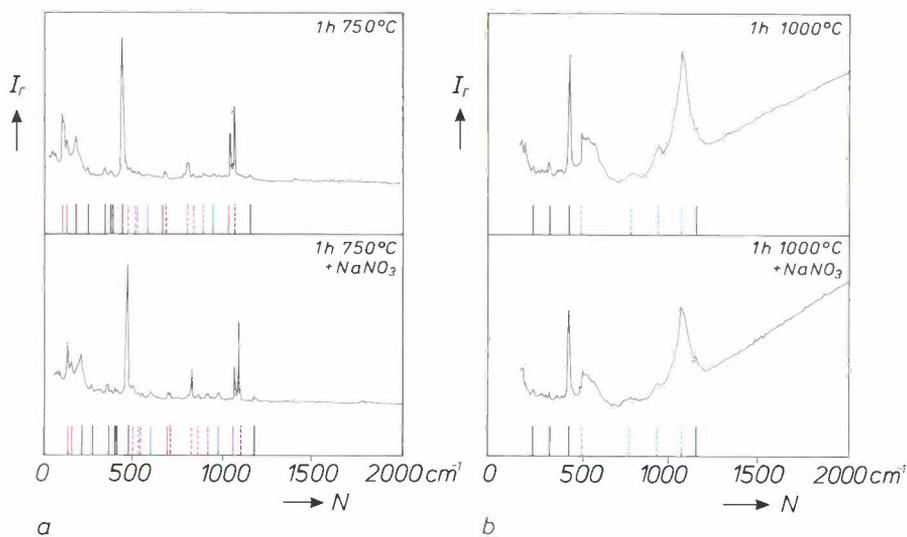


Fig. 4. As fig. 3, but here the spectra shown below originate from samples with added NaNO_3 . The peaks here have been identified as originating from α -quartz (short black vertical line), Na_2CO_3 (black dashed line); BaCO_3 (red line), $\text{Ba}_2\text{SiO}_4 + \text{Na}_2\text{Si}_2\text{O}_7$ (red dotted line), BaSiO_3 (blue line) and liquid disilicate (blue dashed line).

c) At the end of the fining stage the glass melt only contains small bubbles, which mainly consist of oxygen. If the temperature decreases, the oxygen in these bubbles goes into solution again in the melt, so that the bubbles are 'reabsorbed'.

In our study of the fining stage^[6] we have worked with a glass composition that was very close to the composition used in practice: the raw-materials mix

that we prepared consisted of SiO_2 , Al_2O_3 , K_2CO_3 , Na_2CO_3 , MgO , CaCO_3 , SrCO_3 , BaCO_3 in mole percentages of 72, 2.2, 4.9, 10.15, 2.15, 2.1, 0.15 and 5.6 respectively. Although earlier investigations included Sb_2O_3 , we did not use this now little-used material. Instead we used 0.1 mol% of $\text{NaSbO}_3 \cdot \text{H}_2\text{O}$, which is also used in practice. We also added 0.65 mol% of NaNO_3 .

These samples were heated for 30 minutes, 10, 20 or 200 hours at various temperatures and then quenched to room temperature. With this realistic glass composition, as we mentioned earlier, Raman spectrometry is no longer the preferable method for studying the reactions. We confined ourselves to following the change in the equilibrium of the redox couple $\text{Sb}^{3+}/\text{Sb}^{5+}$, by determining the quantity of Sb^{3+} ions by potentiometric titration and the total amount of Sb (Sb_T) by emission spectrometry.

Some results are given in Table III. On increasing the temperature, the oxidation-reduction equilibrium evidently shifts towards Sb^{3+} , which leads to the desired formation of oxygen. This seemed to indicate that at temperatures above 1000°C equilibrium would be reached after only 10 hours. However, further investigation showed that this conclusion was not justified. It was found, for example, that on heating at 1200°C for a much longer period, 200 hours, the equilibrium continued to shift, with the $\text{Sb}^{3+}/\text{Sb}^{5+}$ ratio changing from 0.58 (after heating for 20 hours) to 0.67 (see also Table IV).

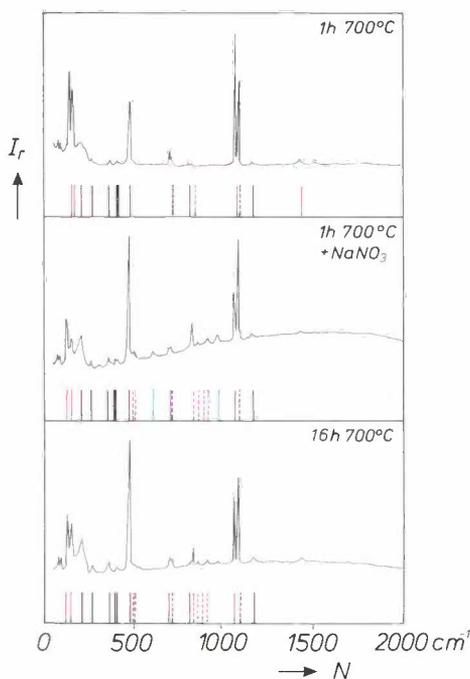


Fig. 5. As in fig. 4, but here the sample without NaNO_3 has been heated again for a further 16 hours (below). Heating for one hour with NaNO_3 (centre) is always found to produce more reaction products than heating for 16 hours without NaNO_3 .

[6] A more detailed account of this investigation is given in: D. M. Krol and P. J. Rommers, *Glass Technol.* 25, 115-118, 1984.

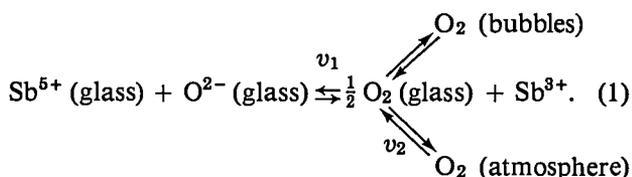
Table III. The ratio Sb^{3+}/Sb_T for 'raw-materials samples' as a function of the temperature T to which they have been subjected for 0.5, 10 or 20 hours.

T (°C)	Time		
	0.5 h	10 h	20 h
900	0.011 ± 0.006	0.018 ± 0.009	0.030 ± 0.004
1000	0.029 ± 0.015	0.100 ± 0.008	0.107 ± 0.009
1100	0.14 ± 0.02	0.28 ± 0.02	0.32 ± 0.02
1200	0.45 ± 0.03	0.60 ± 0.04	0.58 ± 0.04
1300	0.62 ± 0.05	0.79 ± 0.05	0.76 ± 0.05
1400	0.75 ± 0.05	0.89 ± 0.06	0.87 ± 0.06
1500	0.84 ± 0.05	1.00 ± 0.07	0.93 ± 0.06

Table IV. The ratio Sb^{3+}/Sb_T for three types of sample that we used; temperature and time for the treatment are given in the Table.

Types of sample	Temperature		
	1000 °C	1200 °C	1400 °C
Raw materials		0.67 (200 h) 0.58 (20 h)	
Cullet		0.95 (200 h) 0.94 (20 h)	
Powdered cullet 4 × 20 h	0.76	0.69	0.95
Powdered cullet 5 × 20 h	0.67	0.69	0.95

Our explanation for this is that the oxidation/reduction ratio Sb^{3+}/Sb^{5+} is determined by two processes:



First of all an equilibrium with the oxygen in the bubbles becomes established at a relatively high rate v_1 and an equilibrium with the oxygen in the furnace atmosphere only becomes established later, at a much slower rate v_2 .

Indications that this hypothesis is correct have also been found in another part of the investigation, where we were particularly interested in studying the effect of the addition of cullet.

We made 'cullet samples' from glass that had been obtained by melting the raw materials described earlier at 1500 °C. The standard procedure was then followed with these 'ordinary cullet samples': heating to temperatures between 900 and 1500 °C and quenching. With the 'powdered cullet samples' this procedure was modified by repeating the heating and cooling in periods of 20 hours, pulverizing the sample again after each cooling.

For the powdered cullet samples we were able to establish that at temperatures of 1200 °C and above there was no further change in the Sb^{3+}/Sb_T content, so that it may be assumed that an equilibrium with the oxygen in the atmosphere is attained at these temperatures after this treatment (Table IV). The value of Sb^{3+}/Sb_T found at 1200 °C is then indeed found to be just as high as the value measured for the 'raw-materials sample' at the same temperature after the lengthy treatment for 200 hours.

The values of Sb^{3+}/Sb_T that are found as a function of temperature and time for an 'ordinary cullet sample' are however very different (fig. 6 and Table IV). In this case the Sb^{3+}/Sb_T content at each temperature has a value in the region of 0.95, i.e. the sample maintains the equilibrium value that it acquired in the initial melt at 1500 °C.

In our view the explanation for this difference is that the oxygen that is necessary for the reoxidation of Sb^{3+} to Sb^{5+} in the non-powdered — massive — cullet sample has to travel a much longer diffusion path than it does in the powdered sample. This seems reasonable because even after 200 hours the diffusion length of O_2 at 1200 °C is only 4.4×10^{-2} cm, and thus negligibly small compared with the dimensions of the sample. In the powdered cullet sample the oxygen that will diffuse into the interior is present around each separate grain.

This result explains the well-known practical result that the fining agent does not work properly if too much cullet is added to the mix.

Similarly, in the raw-materials sample, if it was only the oxygen from the surrounding atmosphere that was

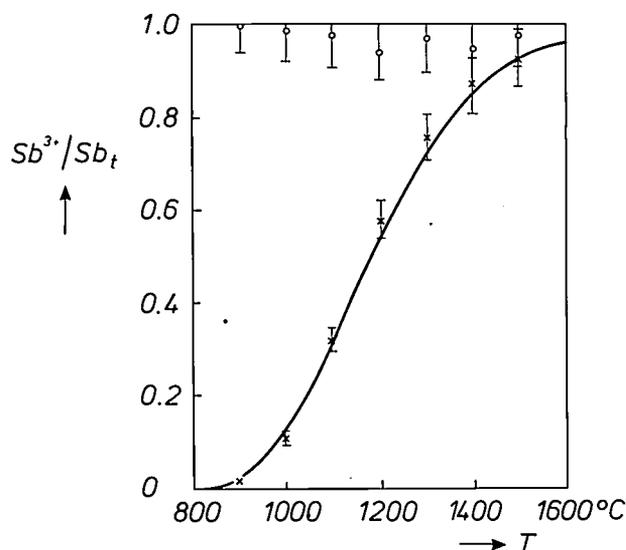


Fig. 6. The ratio Sb^{3+}/Sb_T for 'raw-materials samples' (crosses) and 'powdered cullet samples' as a function of the temperature to which the samples have been subjected for 20 hours.

significant, then there would be no reason for expecting that the effect of temperature would result in a noticeable shift of the $\text{Sb}^{3+}/\text{Sb}^{5+}$ couple, since this shift would be determined by the very slow oxygen-diffusion process. In the raw-materials sample, however, as we have seen, there is the possibility of oxygen occurring in the form of bubbles in the melt, with a partial oxygen pressure of 1 bar. This is a higher oxygen pressure than is reached in the equilibrium state of the powdered cullet sample (0.2 bar). This means that in the raw-materials sample the $\text{Sb}^{3+}/\text{Sb}_T$ ratio will be lower than in the powdered cullet (corresponding to a high p_{O_2} (glass)) for relatively short melting times (20 hours). From Table IV and fig. 6 we have been able to show that the ratio $\text{Sb}^{3+}/\text{Sb}^{5+}$ at 1200 °C (after 20 hours) is equal to 1.4 ± 0.2 for a raw-materials sample. For a powdered cullet sample this value is 2.2 ± 0.6 .

If in the equilibrium expression

$$K' = \frac{[\text{Sb}^{3+}](p_{\text{O}_2})^{\frac{1}{2}}}{[\text{Sb}^{5+}]},$$

we use the value $p_{\text{O}_2} = 0.2$ for the powdered cullet sample, then with the $\text{Sb}^{3+}/\text{Sb}^{5+}$ -values given above for the raw-materials sample we find a p_{O_2} -value of 0.5 bar. In view of the experimental uncertainties in the starting points for the calculation, this value agrees reasonably well with the value of 1 bar assumed earlier.

The influence of NaNO_3

Finally, we shall say a few words about the influence of NaNO_3 , which, as our results also indicate, becomes evident at the fining stage.

At a temperature of about 550 °C the NaNO_3 decomposes accompanied by the formation of oxygen. This means that the oxygen concentration (or more accurately, the oxygen activity) in the glass melt will be increased, and this will happen before it does so under the influence of the equilibrium shift of the redox couple $\text{Sb}^{3+}/\text{Sb}^{5+}$. It will therefore also increase the temperature at which this shift commences (fig. 7). In principle, such an increase in the temperature has the beneficial effect that the formation of oxygen can be 'put off' until the formation of CO_2 has been fully completed.

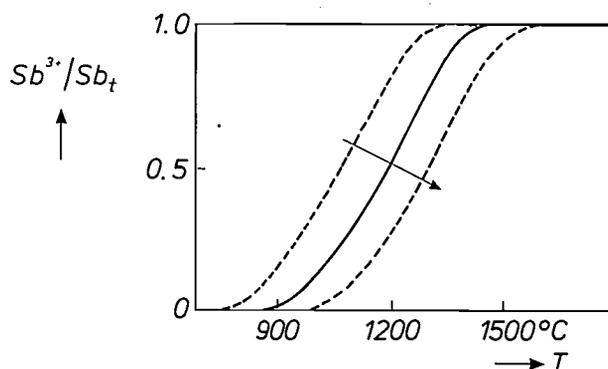


Fig. 7. The ratio $\text{Sb}^{3+}/\text{Sb}_T$ for raw-materials samples as a function of the temperature to which the samples have been subjected for 20 hours, at varying values of the initial oxygen activity (increasing in the direction of the arrow).

An increase in the initial oxygen concentration by the decomposition of NaNO_3 has the additional advantage that less oxygen has to be produced for the melt to be saturated. In this way, therefore, the NaNO_3 contributes to the formation, growth and ascent to the surface of the oxygen bubbles and consequently to the removal of the bubbles of CO_2 , N_2 , etc.

It is very significant that these are the particular bubbles that are removed, because they are not re-absorbed when the temperature is reduced. And it can be seen from the reaction equation (1) why the remaining oxygen bubbles are indeed absorbed in this case.

Our relatively simple equation therefore helps to explain the empirically well-known beneficial action of the combination of NaNO_3 and Sb on the formation of bubble-free glass.

Summary. The chemical conversions in the melting stage of glass of the composition $15\text{Na}_2\text{CO}_3\text{-}10\text{BaCO}_3\text{-}75\text{SiO}_2$ have been investigated with the aid of laser-Raman spectrometry. Below 1000 °C a number of different crystalline compounds are formed, whose silicon content increases with increasing temperature. Between 1000 and 1400 °C the melt consists of liquid disilicates or trisilicates, containing dispersed crystalline α -quartz. The oxidation/reduction equilibrium of the fining agent antimony in glass whose composition closely corresponds to that of television glass has been studied by potentiometric titration and emission spectrometry. A reaction equation can be written that gives a good description of this behaviour. The possible influence of NaNO_3 and cullet on the formation and fining of television glass has also been investigated.

A mobile system for image bulk storage

L. H. Guildford and B. D. Young

There are many kinds of video signals besides those from television. Many other kinds of systems produce images from sensors; they range from infrared surveillance systems to X-ray diagnostic equipment. Sometimes it is essential to be able to record signals for later analysis off-site. This requires far more elaborate equipment for image bulk storage than an ordinary video recorder. Recently at Philips Research Laboratories hardware and software have been developed to cope with the many possible eventualities. The instrumentation is in fact a complex system that constitutes a complete transportable laboratory.

Introduction

Television cameras, night-viewing systems, infrared surveillance systems, ultrasound systems, radar, sonar, X-ray diagnostic equipment — these all produce images from sensor systems. The images have to be processed for simple display, to activate further systems or perhaps to provide a better understanding of the observations or of the operation of the actual system.

Experiments in real-time image processing with such systems are often desirable, but there are problems. The actual scene itself may be unique and economically unrepeatable, taking place in difficult environmental conditions. The image might for example be due to a low-flying aircraft detected and tracked by infrared sensors in unusual weather conditions. It could come from an X-ray medical investigation, where the patient's discomfort should be kept as brief as possible.

It is therefore highly desirable to record the video information from the sensors on the spot. Then the images can be recreated, studied and processed at will later. The recording must be such that high-quality unadulterated video information from the sensor system is readily available.

Ultimately there are various options: the information may be viewed directly as a two-dimensional picture, plotted as a graph or perhaps used as the basic data for the recognition of a characteristic pattern (an 'object signature') to trigger automatic reaction by a system.

Collecting and recording the video signals from different types of sensor systems requires a sophisticated image bulk store with the appropriate equipment for preprocessing and postprocessing. Such a storage system including all the ancillary hardware forms the subject of this article. The system is designed for monochrome ('black and white') images only.

Image signal characteristics

Some of the sensor systems listed above will present information as a signal that is compatible with the international CCIR standards for television^[1]. Others make use of more unusual scene-scanning formats. In general these will be line-scanning structures based on rectilinear formats (as in television, but with different parameter values), radial line scans (as in radar) or vertical line scans swept horizontally to produce a 360° 'cylindrical' scan for omnidirectional surveillance.

The signals from the sensors will often be in analog form, sometimes with synchronization pulses interleaved with the video signals proper, thus producing composite waveforms. Sometimes, however, video and synchronization signals are provided on separate connections. Digitized signals supplied directly from the sensors will become more common in future. This means that our system for image bulk storage must be capable of dealing with both analog and digital signal formats and with a wide variety of scan-synchronization signals.

However the signals from the sensors are to be displayed or used, certain operational parameters always

L. H. Guildford, M.I.E.R.E., was formerly with Philips Research Laboratories, Redhill, Surrey, England; B. D. Young, B.Sc., A.M.I.E.E., is with these laboratories.

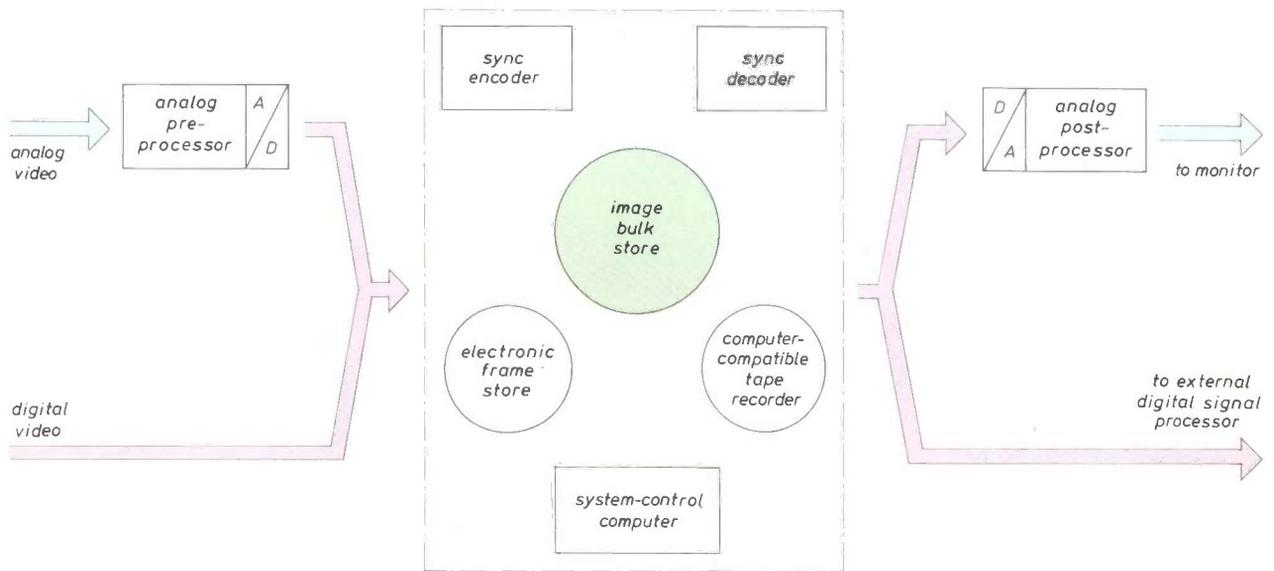


Fig. 1. Basic functional diagram of the mobile system for image bulk storage. The system consists of a number of subsystems that can be interconnected in a variety of ways to suit the particular needs of the user. All signal processing is basically digital. Analog video signals are therefore first converted to digital and back to analog after processing. The actual image bulk store is a 24-track tape recorder. An electronic frame store and a computer-compatible tape recorder are also available. For synchronization purposes special encoder and decoder circuits have been designed. The entire system is controlled from the system-control computer.

have to be established and maintained for recording and analysing the information.

- A zero reference level for the video signal (the 'black level'). This can be either absolute or related to the scene.
- The dynamic range for the signal from the sensors. This is normally expressed as the range between the black level, as a reference value, and the peak level representing white.
- The resolution of the signal. Two kinds of resolution are significant: the spatial resolution in terms of the number of picture elements (pixels) or samples within a given image region, and the amplitude resolution in terms of signal-to-noise ratio, minimum detectable signal or number of bits per sample.
- The system characteristics of the overall system from the output of the sensors to the point of measurement. They should be linear and the corresponding frequency response should be as flat as possible.

For infrared systems the dynamic range of the video signal can be greater than 60 dB (an amplitude ratio of 1000:1, requiring at least 10 bits per sample), but mostly a dynamic range of 48 dB (an amplitude ratio of 250:1, requiring at least 8 bits per sample) is sufficient.

For monochrome television signals a sampling rate of about 12 MHz is necessary. In processing the signals of other image-sensing applications of interest sampling rates as high as 18 MHz can occur, especially when two or more signals are combined in time-

division multiplex to allow a number of different sensor devices to view the same scene simultaneously.

All the normal scan formats produce some degree of 'dead time' in which no video information is produced. For instance, the scanning formats used with standard television systems allow approximately 20% of the time for scan flyback. Infrared scanners with optomechanical scanning often produce as much as 50% of dead time. By storing and reformatting the incoming information, the effective bandwidth required of a recording medium can be reduced. It is also possible to reduce the amount of storage required per frame, if the information relating to the dead times is kept to the absolute minimum.

System requirements

From the outset we wanted our system for image bulk storage to meet the following basic requirements.

- It must make a high-quality record of incoming images.
- It must enable the replay of recorded information at the original recording speed and at various much lower speeds.

[1] Report No. 308-2, Characteristics of monochrome television systems, Proc. XIIth CCIR Plenary Assembly, Vol. V, pt. 2, New Delhi 1970 (published by the ITU, Geneva 1970); Report No. 624-2, Characteristics of television systems, Recommendations and reports of the CCIR, Vol. XI, pt. 1, Broadcasting service (television), ITU, Geneva 1982.

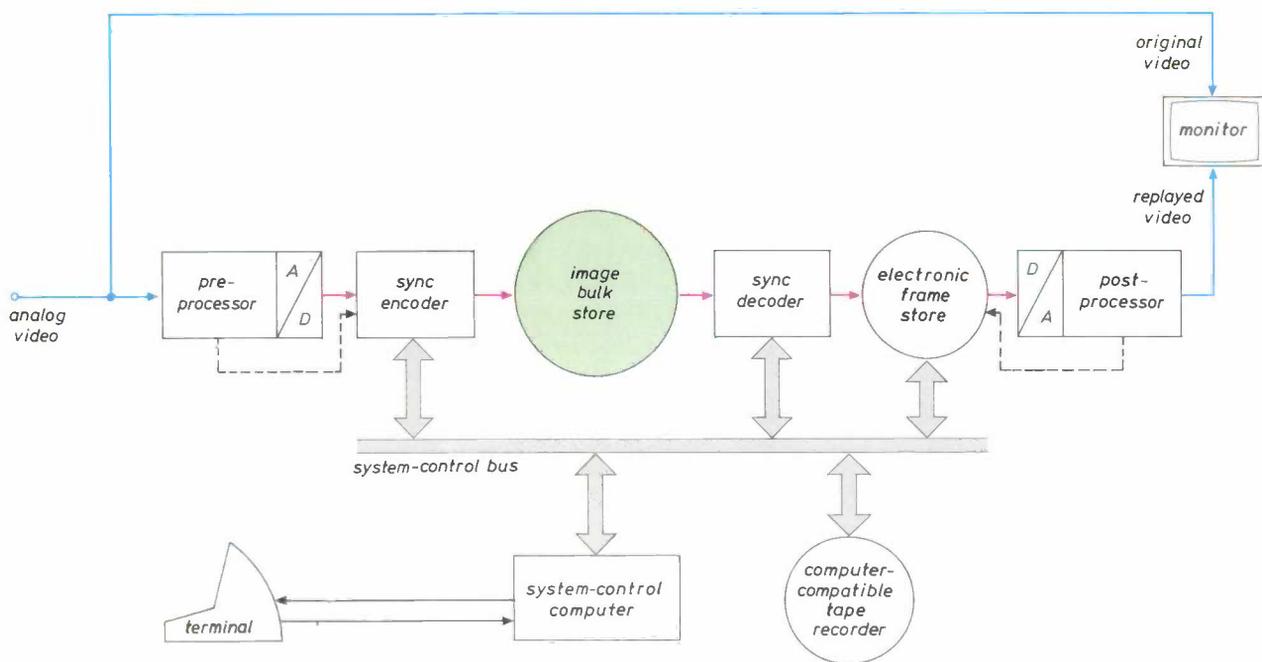


Fig. 2. System configuration for storing standard analog television signals. After analog preprocessing and A/D conversion the sync encoder replaces the conventional sync information by digital information for both synchronization and image-identification purposes. The signal is then stored in the image bulk store. On replay individual frames are stored in the frame store. They can be transferred to the computer-compatible tape recorder for further off-line processing or they can be viewed on a monitor after D/A conversion and appropriate postprocessing. Connections intended for synchronization purposes are indicated by dashed lines. The system is controlled from a computer terminal via the system-control computer and the system-control bus.

- It must be possible to obtain still pictures ('frozen frames'), so that information can be examined frame by frame.
- It must include a computer-compatible tape recorder (with 1/2-inch tape) for recording a selected still picture. Recorded images can then be analysed off-line by an external computer, manipulated as required within that computer and replayed at will for further use.
- The hardware of the complete system must be self-contained and transportable.

Now that the image signal characteristics, described in the previous section, are known the basic functional diagram of *fig. 1* can be derived; it shows the individual modules required. Because the main functions are executed digitally any analog video input signal is first preprocessed and converted to digital. Similarly any analog video output signal is obtained from the system after digital-to-analog conversion and post-processing. The heart of the system consists of three storage media:

- an electronic store that can contain two television frames at a time,
- an image tape recorder, which constitutes the actual image bulk store, and

- a computer-compatible tape recorder that serves as an output device enabling external processing later, e.g. on a mainframe computer.

As we are dealing with video signals, synchronization is an important aspect of the interaction of input devices, output devices and storage media. Two separate units indicated by 'sync encoder' and 'sync decoder' are provided for this purpose.

The system-control computer is used to set up and monitor all operations.

System architecture

The modules of *fig. 1* are deliberately shown without interconnections; depending on the particular application of the system the modules may be connected in many different ways. For instance, the frame store may be placed either before or after the bulk store or not used at all. There are many possible configurations to suit the user's needs, both for recording and replaying television or non-standard image formats. All modules have therefore been made plug-compatible at their interfaces so that they can easily be patched together to produce a great variety of recording and replay functions.

One specific system configuration with its interconnections is shown in *fig. 2*. This configuration can be used for processing a standard analog television signal. After preprocessing, A/D conversion and sync encoding the signal is recorded digitally in the bulk store. On replaying the video recording the signal passes through the sync decoder, is reformatted in the frame store and after subsequent D/A conversion and post-processing is replayed on a television monitor. Alternatively individual frames may be selected for detailed examination on the monitor or recorded on the computer-compatible tape recorder.

Some particular features of this configuration are:

- the monitor can display either the original or the replayed video signal;



Fig. 3. System hardware under test in the laboratory. From left to right: preprocessing and postprocessing rack with monitor; system-control computer and computer-compatible tape recorder; 24-track digital image tape recorder (the actual image bulk store); electronics of digital image tape recorder in combination with sync encoder and sync decoder; electronic frame store control console.



Fig. 4. The system hardware installed in a container forms a mobile laboratory.



Fig. 5. Loading the mobile laboratory on to a transporter.

- the image bulk store can be by-passed, which means that the system can be set up and tested without actually running the bulk tape;
- when incoming signals are being recorded, the replay circuits are also operational permitting the recorded information to be continuously monitored.

An impression of the actual hardware can be obtained from *figs 3 to 5*. In *fig. 3* a laboratory arrangement for testing is shown. All modules and some additional equipment for heating, lighting, power supply, dust filtering etc. have been installed in a thermally insulated container (*fig. 4*). The resulting mobile laboratory is transportable as can be seen from *fig. 5*.

Now let us consider the operation of the individual modules in some more detail.

System modules

The preprocessor

The input signal to the preprocessor (*fig. 6*) is assumed to be an analog video signal with a peak-to-peak amplitude of 1 V and with the zero level corresponding to 'black' (it is 'black-level clamped'). The chief tasks of the analog preprocessor are to remove and separate the synchronization signals, to normalize the remaining video signal to the range (-0.5 V, 0.5 V) and to generate clock pulses for the A/D conversion. (By normalizing the video-signal amplitude and removing the synchronization signals before A/D conversion, the input range of the A/D converter is fully exploited.) The preprocessor also comprises the anti-aliasing filter that is required to prevent 'fold-over' distortion in the A/D conversion^[2]. The sampling

^[2] See for example pp. 135-136 in: A. W. M. van den Enden and N. A. M. Verhoeckx, Digital signal processing: theoretical background, Philips Tech. Rev. 42, 110-144, 1985.

times for the A/D converter are determined by the clock pulses which can be made to appear regularly between 1 and 1024 times during each line time. The output samples (each corresponding to one pixel) have the form of binary words consisting of 8 bits in parallel. The maximum sampling rate is 15 MHz.

In the preprocessor the composite sync signal that is stripped from the incoming video is separated into a line sync signal that has a short pulse at the start of each line and a frame sync signal that has a long pulse at the start of each television field.

The sync encoder

The main task of the sync encoder is the insertion of digitally encoded sync information (the 'sync code block') into the stream of video samples to be pre-

The image bulk store

The heart of our system is the image bulk store. It is a commercially available 24-track digital magnetic tape recorder capable of recording and replaying binary sig-

Typical sync code block structure

A typical 16-byte sync code block is shown in fig. 7. The 16 bytes are numbered 0, 1, 2, ..., 15 and the 8 bits of each byte are numbered 0, 1, 2, ..., 7. The pattern of zeros and ones for bit 0 is the same for each sync code block; bits 1, 2 and 3 constitute the 48-bit frame title. The combination of bits 4 to 7 of bytes 2, 6, 10 and 14 yield the frame number (in this case: 0101 1000 0111 0000 = 22640). In all the other bytes the bits 4 to 7 are 'inverted reflected' versions of bits 0 to 3. They carry redundant information for error checking.

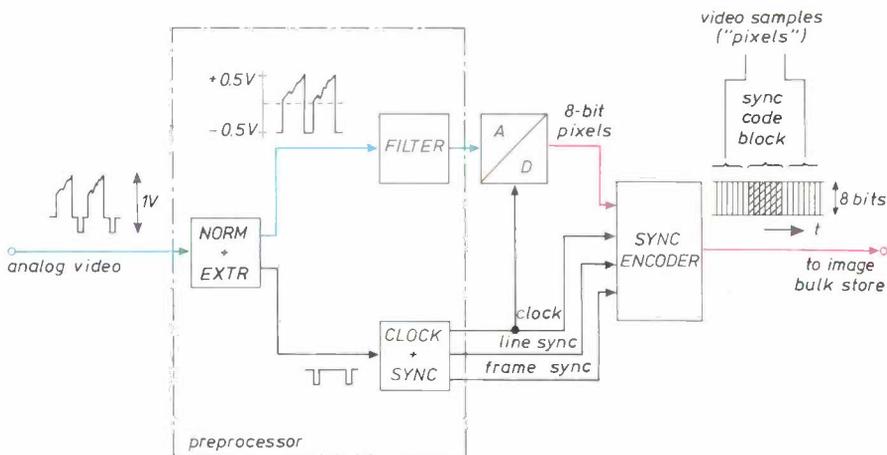


Fig. 6. Block diagram of the main circuits that are active during recording in the configuration of fig. 2. *NORM+EXTR* circuits for normalization of the video signal amplitude and extraction of sync signals. *FILTER* anti-aliasing filter inserted in front of the A/D converter. *CLOCK+SYNC* tunable clock generator determining the sampling rate of the A/D converter (= pixel rate) and generator of line-sync signals and frame-sync signals.

mented to the image bulk store. The sync code block consists of 16 successive 8-bit words ('bytes') that can easily be inserted between the 8-bit video samples. This is done at the beginning of each line period. One part of the sync code block (16 bits) is always the same and characterizes the occurrence of the block. A second part (16 bits) represents a frame number; an increase of this number by one indicates the start of a new frame. A third part (48 bits) can be used in an arbitrary way for identifying a series of frames and is therefore called the 'frame title'; these titles can be keyed in from the system-control computer. The remaining 48 bits are redundant and can provide added protection against errors.

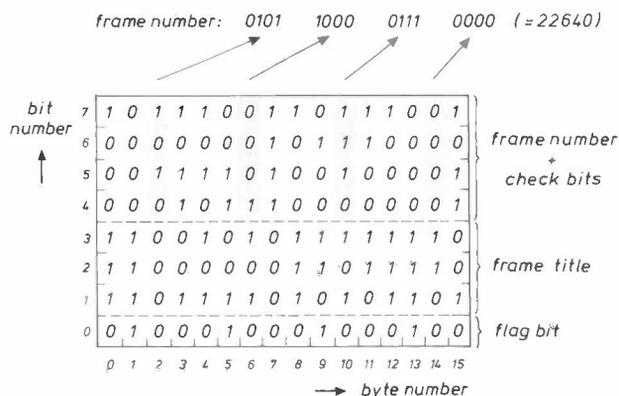


Fig. 7. Typical sync code block consisting of 16 bytes of 8 bits.

nals at rates from 0 to 3.3 Mbit/s on each track. As our video samples ('pixels') consist of 8 bits three samples can be recorded in parallel and we can accommodate 9.9×10^6 video samples per second in a straightforward way. By exploiting the 'dead times' that occur in almost any video information (p. 261) we can increase the effective recording rate. The incoming video samples are therefore stored and reformatted frame by frame in the frame store before recording. The inverse process must then be applied on replay.

The sync decoder

When the sync decoder is presented with an 8-bit parallel data stream from the image bulk store on replay, it will look for sync codes as described above. If a code block is detected the sync decoder will issue a

ceed certain preset thresholds that allow for some tape drop-outs that might occur in a code block. Once the sync code block has been detected, the detector is inactive for almost the entire estimated duration of a line period ('temporal filtering' or 'time windowing'). A time window is applied to the frame synchronization in a similar way.

Individual frame selection

To provide the ability to select any particular frame for storage in the frame store a comparator circuit in the sync decoder can be programmed to search for a specific frame number. On detection of this number, frame-sync generation is inhibited, thus preventing further updates of the frame store. This facility is particularly useful when a set of sequential images is required for off-line analysis on a mainframe computer.

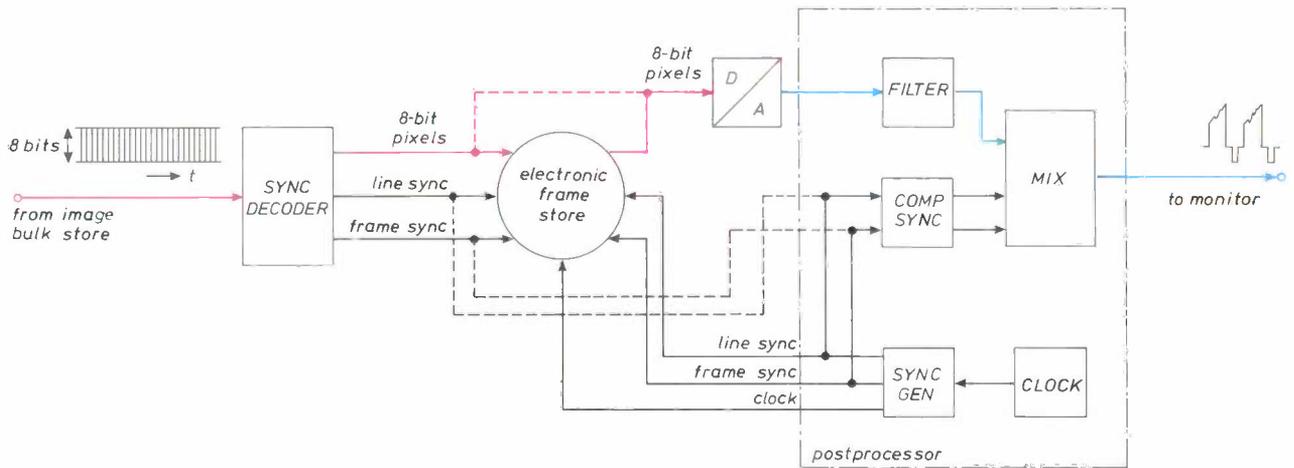


Fig. 8. Block diagram of the main circuits that are active during replay in the configuration of fig. 2. CLOCK tunable clock generator determining the pixel rate. SYNC GEN generator of line-sync signals and frame-sync signals. COMP SYNC circuit for generating composite sync signal and blanking signal. MIX circuit for combining the analog video signal (obtained after D/A conversion and lowpass filtering) with the composite sync signal and the blanking signal. The postprocessor can be connected directly to the sync decoder (dashed lines) when the frame store is not used.

line sync pulse, and if an increment of the frame number by 1 is observed a frame sync pulse is generated (fig. 8).

Initially the sync decoder is programmed from the system-control computer via the system-control bus with the frame title of the required frame. In the decoder the correlation between the actual incoming data stream and the reference title is examined continuously. This is done at the full sampling rate (hence at a rate of 9.9 MHz). The redundant bits in the sync code are also checked. The results of these two measurements need not yield 100% agreement, but should ex-

The frame store

The frame store is one of the key components that makes the overall system so versatile; it can be positioned at many different locations in the system configuration for performing many kinds of buffering and reformatting operations. As a consequence the frame store itself is a rather complicated subsystem, and only its most important features will be mentioned here.

The frame store (fig. 9) is a dedicated piece of hardware, designed at our Laboratories. It is a random-access memory (RAM) with a capacity of approximately half a million 12-bit words (as our pixels are

represented by 8-bit words, 4 bits of each word are not actually used at present). The memory is organized into two separate pages that can each store 512×512 12-bit words. Each page constitutes a complete video frame. In the fastest addressing mode (the 'linear addressing mode') data words can be written into the store at a rate of 18.6 MHz and read out at a maximum rate of 15.6 MHz.

To achieve independence between read and write accesses, the store has individual read and write ports serviced by their own address generators.

The frame store has a single-board microcomputer. Its control console provides the user with local, stand-alone operation of the frame store for manual data entry and control of the address generators. The microcomputer is also connected to the system-con-

trol sequential access to all 16 segments gives a theoretical maximum overall update rate of 20 MHz. For this sequential access 4 bits are required both at the input and at the output of the 16 segments; to compensate for the internal delay in the segments a small 'read segment delay' is provided as indicated in fig. 9.

As the operation of the dynamic memory is asynchronous to the memory accesses, extensive data buffering is employed at the input and output parts of the memory (not shown explicitly in the figure).

The postprocessor

The final operations necessary for obtaining a common composite video output signal from our system are performed in the postprocessor (fig. 8). These include both analog filtering, required for signal recon-

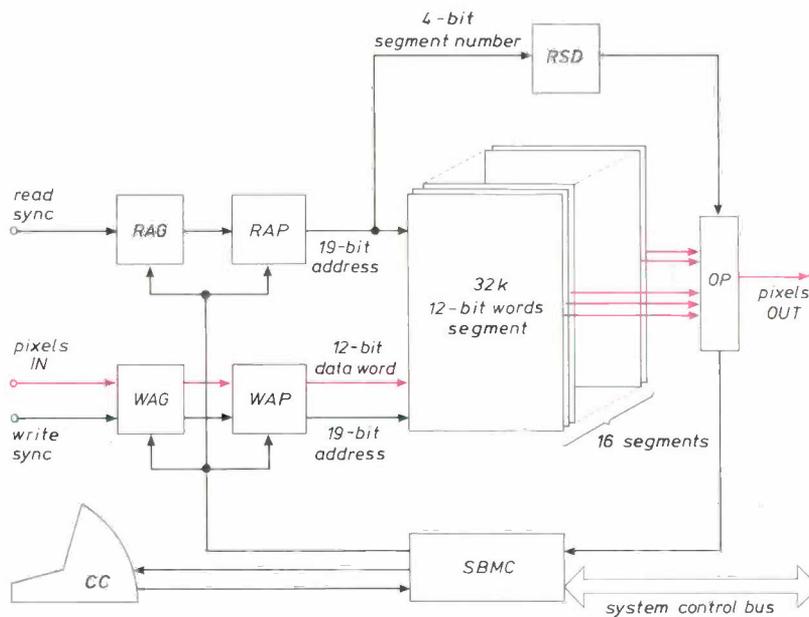


Fig. 9. Frame-store schematic diagram. The store consists of 16 parallel segments of dynamic memory, each segment comprising 32k of 12-bit words. RAG read address generator. RAP read-access port. WAG write address generator. WAP write-access port. RSD read segment delay. OP output port. SBMC single-board microcomputer. CC control console.

trol bus, interpreting and controlling communication between the system computer and the frame store. By programming the correct sequences of addresses used in writing and reading, we can let the frame store perform many types of formatting operations on the stored video information.

The frame store itself is built up from 16 parallel segments, each comprising 32k ($= 2^{15}$) 12-bit words of dynamic memory. In each segment a combined read/write cycle can be performed in 800 ns corresponding to a maximum update rate of 1.25 MHz. In-

struction after D/A conversion, and the insertion of combined line and frame-sync signals into the video signal. The postprocessor therefore has a pixel clock generator and various circuits for generating the correct sync and blanking signals. When the frame store is used, the reading process of the store is 'slaved' to the sync signals of the postprocessor (while the writing process is under the control of the sync decoder). If the frame store is not used, the synchronization of the postprocessor can be slaved to the sync decoder (indicated by dashed lines).

The system-control computer

The system-control computer is a commercially available piece of hardware (Texas Instruments 990/4) that gives the user control over many of the programmable features of our image bulk-storage system. Among these are the following facilities:

- choice of master-clock frequency between 873 kHz and 25 MHz;
- programming of address patterns for the frame store;
- choice of sync code, including the 48 bits for frame titles;
- transfer of video data from the frame store to the computer-compatible tape recorder or vice-versa;
- execution of various types of tests on parts of the system;
- storage on computer disk of particular system parameter settings that can be used for setting up the system or examined on print-out.

Application examples

From the many possible applications of our image bulk-storage system we have chosen one particular example that clearly demonstrates a number of important features. The picture information that we start with in this case originates from a system called IRSCAN. This is an omnidirectional (360°) infrared surveillance system that has a vertical angle of view of

4°. It uses a vertical line-scanning pattern comprising 16 384 lines in one complete rotation, which takes 0.5 s. This scanning format is symbolically represented in *fig. 10a*.

Our storage system can be used to replay pictures or parts of pictures from this surveillance system on a standard television monitor for analysis. Each vertical line in the IRSCAN signal is first converted into 120 pixels, each represented by an 8-bit word. Before recording on the image bulk-store tape recorder, each group of 120 pixels is augmented by a sync code block containing information about the line number and the particular azimuth scan number.

On replay, any part of an IRSCAN picture consisting of 512 lines (blue area in *fig. 10a*) is written into the frame store while a vertical line addressing format is used. For read-out from this store a horizontal line addressing format is used, corresponding to a standard television scan. In this way the contents of the frame store can be displayed in the blue area of the

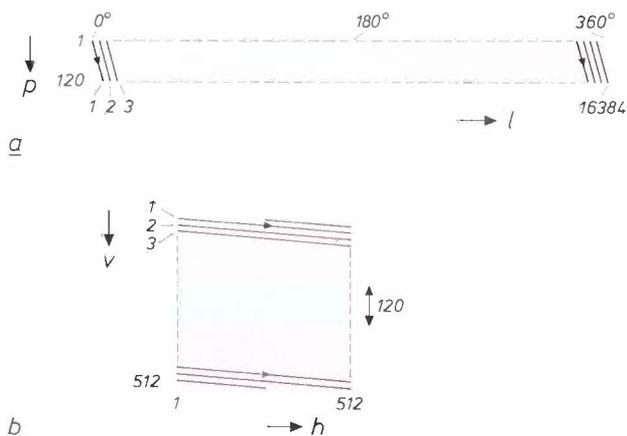


Fig. 10. *a)* Scanning format in the omnidirectional infrared surveillance system IRSCAN. There are 16 384 vertical scanning lines in each complete rotation (360°). At A/D conversion 120 pixels are derived from each line. *p* pixel number, *l* line-scan number in the IRSCAN system. *b)* After storage in the (512 × 512)-pixel frame store and reformatting, a part of the IRSCAN image can be displayed on a monitor using the standard horizontal scanning format. The blue area in the upper drawing is then represented as the blue area in the lower diagram. *v* vertical address number of the frame store corresponding to the line number in standard television, *h* horizontal address number of the frame store determining the horizontal position on the standard television monitor.

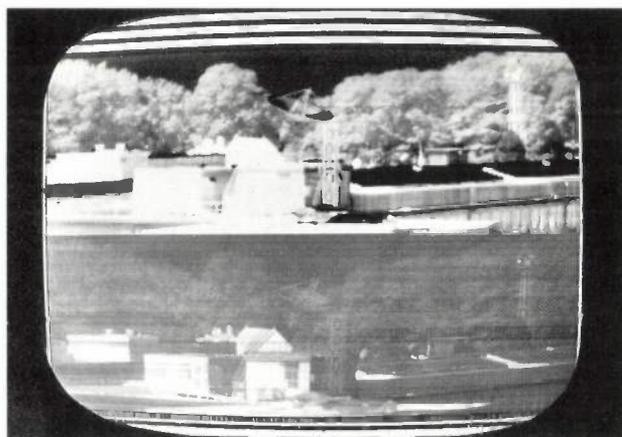


Fig. 11. An application of the image bulk-storage system. Here two optically coincident infrared images from the IRSCAN system are combined on one television screen. The upper image displays information obtained at wavelengths of 8 to 14 μm; the lower image represents wavelengths of 3 to 5 μm.

television screen, schematically represented in *fig. 10b*. (In practice certain operations related to the usual interlacing also have to be incorporated into the addressing format of the store.) In the frame store parts of different pictures can also be combined for simultaneous display, permitting the comparison of different pictorial aspects.

Some practical results are shown in *figs 11 and 12*. In *fig. 11* two coincident infrared images of the same scene are shown one on top of the other on one



Fig. 12. Sometimes it is very useful to compare images obtained from the same scene by different methods using the image bulk-storage system. Here an optical image (top) and an infrared image (bottom) of the same outdoor scene are displayed simultaneously on two monitors.

screen. The upper image has been made with an infrared detector for wavelengths of 8 to 14 μm , the lower image with a detector for wavelengths of 3 to 5 μm . The upper image is saturating on highlights due to excessive temperature differentials in the scene.

A second example is given in *fig. 12*. The upper photograph shows an ordinary video image of an outdoor scene and the lower photograph shows an infrared picture of the same scene; useful information can often be gathered from a comparison of such pairs of images.

Summary. Often the value of image signals obtained from sources such as TV cameras, night-viewing systems, infrared sensors, ultrasound systems, X-ray equipment and so on, is increased considerably if the images can be recorded faithfully for later analysis or processing. Unfortunately the signal characteristics (e.g. the scanning format) in all these cases can vary greatly. To cope with this variety a versatile system for image bulk storage has been designed around a standard 80-Mbit/s 24-track digital tape recorder. The images can be reformatted either before or after recording by means of an electronic frame store. Images can be replayed on a standard TV monitor or copied on to a computer-compatible magnetic tape for further off-line processing on a mainframe computer. A control computer manages system operation. The application of a sync-block coding scheme permits easy identification and selection of individual images. The image bulk-storage system constitutes a complete transportable laboratory housed in a standard (6.7 m \times 2.4 m) shipping container.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	A
Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	B
Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	E
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	R
Philips Research Laboratories Sunnyvale P.O. Box 9052, Sunnyvale, CA 94086, U.S.A.	S

M. L. Verheijke	E	Toepassing van neutronenactiveringanalyse, autoradiografie en radioactieve tracers in de silicium-technologie	Anal. Techn. 7	30-34	1986
J. Zhang (<i>Imp. College, London</i>), J. H. Neave, P. J. Dobson & B. A. Joyce	R	Effects of diffraction conditions and processes on RHEED intensity oscillations during the MBE growth of GaAs	Appl. Phys. A 42	317-326	1987
M. Delfino, D. K. Sadana & A. E. Morgan	S	Shallow junction formation by preamorphization with tin implantation	Appl. Phys. Lett. 49	575-577	1986
P. Rutérana, P. Friedel, J. Schneider & J. P. Chevalier (<i>C.E.C.M.-C.N.R.S., Vitry</i>)	L	High resolution electron microscopy of the GaAs/Si ₃ N ₄ interface produced by multipolar plasma deposition	Appl. Phys. Lett. 49	672-673	1986
H. van Houten, B. J. van Wees (<i>Delft Centre for Submicron Technol., Delft</i>), M. G. J. Heijman & J. P. André	E, L	Submicron conducting channels defined by shallow mesa etch in GaAs-AlGaAs heterojunctions	Appl. Phys. Lett. 49	1781-1783	1986
S. Makram-Ebeid, P. Boher & M. Lannoo (<i>ISEN, Lille</i>)	L	Interactions between bombardment-induced defects in GaAs	Appl. Phys. Lett. 50	270-272	1987
J. J. P. Bruines, R. P. M. van Hal, B. H. Koek, M. P. A. Vieggers & H. M. J. Boots	E	Between explosive crystallization and amorphous regrowth: inhomogeneous solidification upon pulsed-laser annealing of amorphous silicon	Appl. Phys. Lett. 50	507-509	1987
J. Maguire*, R. Murray*, R. C. Newman* (<i>*J. J. Thomson Phys. Lab., Whiteknights</i>), R. B. Beall & J. J. Harris	R	Mechanism of compensation in heavily silicon-doped gallium arsenide grown by molecular beam epitaxy	Appl. Phys. Lett. 50	516-518	1987
J. W. M. Bergmans & Y. C. Wong	E	A simulation study of intersymbol interference cancellation	Arch. Elektron. & Uebertragungstech. Bd. 41	33-37	1987
J. Jerrams-Smith (<i>Univ. Birmingham</i>)	R	An expert system within a supportive interface for UNIX	Behav. & Inf. Technol. 6	37-41	1987
R. Waser	A	Diffusion of hydrogen defects in BaTiO ₃ ceramics and SrTiO ₃ single crystals	Ber. Bunsenges. Phys. Chem. 90	1223-1230	1986
W. Hermann, A. Raith & H. Rau	A	Diffusion of fluorine in silica	Ber. Bunsenges. Phys. Chem. 91	56-58	1987
P. H. L. Notten	E	The role of surface charging and potential redistribution on the kinetics of hole injection reactions at n-GaAs	Electrochim. Acta 32	575-581	1987
G. M. Loiacono	N	Crystal growth of KH ₂ PO ₄	Ferroelectrics 71	49-60	1987

- | | | | | | |
|--|-----|---|--------------------------|-----------|------|
| A. Bhattacharyya, J. D. Reimer & K. N. Ritz | S | Breakdown voltage characteristics of thin oxides and their correlation to defects in the oxide as observed by the EBIC technique | IEEE EDL-7 | 58-60 | 1986 |
| F. L. van Nes (<i>Inst. Perception Res., Eindhoven</i>) | | A new teletext character set with enhanced legibility | IEEE Trans. ED-33 | 1222-1225 | 1986 |
| A. Bhattacharyya & S. N. Shabde | S | Degradation of short-channel MOSFET's under constant current stress across gate and drain | IEEE Trans. ED-33 | 1329-1333 | 1986 |
| A. W. Ludikhuizen | E | A versatile 250/300-V IC process for analog and switching applications | IEEE Trans. ED-33 | 2008-2015 | 1986 |
| P. Piret | B | Bounds for codes over the unit circle | IEEE Trans. IT-32 | 760-767 | 1986 |
| P. Delsarte & P. Piret | B | Do most binary linear codes achieve the Glick bound on the covering radius? | IEEE Trans. IT-32 | 826-828 | 1986 |
| P. van Weijer | E | Pulsed optical pumping as a tool for the determination of population mechanisms of excited states in a low-pressure mercury discharge | IEEE Trans. PS-14 | 464-470 | 1986 |
| A. G. Dirks, T. Tien & J. M. Towner | S | Al-Ti and Al-Ti-Si thin alloy films | J. Appl. Phys. 59 | 2010-2014 | 1986 |
| K. N. Ritz, M. Delfino, C. B. Cooper III* & R. A. Powell* (<i>*Varian Res. Center, Palo Alto, CA</i>) | S | Observation of slip dislocations in (100) silicon wafers after BF ₂ ion implantation and rapid thermal annealing | J. Appl. Phys. 60 | 800-802 | 1986 |
| J. Khurgin, B. J. Fitzpatrick & W. Seemungal | N | Cathodoluminescence, gain, and stimulated emission in electron-beam-pumped ZnCdSe | J. Appl. Phys. 61 | 1606-1609 | 1987 |
| G. J. van Gorp, P. R. Boudewijn, M. N. C. Kempners & D. L. A. Tjaden | E | Zinc diffusion in <i>n</i> -type indium phosphide | J. Appl. Phys. 61 | 1846-1855 | 1987 |
| P. J. Schoenmakers, F. C. C. J. G. Verhoeven & H. M. van den Bogaert | E | Application of supercritical fluid chromatography to the analysis of liquid-crystal mixtures | J. Chromatogr. 371 | 121-134 | 1986 |
| P. J. Schoenmakers & T. Blaffert | E,H | Effect of model inaccuracy on selectivity optimization procedures in reversed-phase liquid chromatography | J. Chromatogr. 384 | 117-133 | 1987 |
| C. Guedon, J. Le Bris & J. L. Gentner | L | Control of interface formation during growth of InGaAs/InP heterostructures by chloride vapour phase epitaxy | J. Cryst. Growth 79 | 909-913 | 1986 |
| E. P. Menu, D. Moroni, J. N. Patillon, T. Ngo & J. P. André | L | High mobility of two-dimensional electrons in Ga _{1-x} In _x As/InP heterostructures grown by atmospheric pressure MOVPE | J. Cryst. Growth 79 | 920-927 | 1986 |
| P. J. Dobson, B. A. Joyce, J. H. Neave & J. Zhang (<i>Imp. College, London</i>) | R | Current understanding and applications of the RHEED intensity oscillation technique | J. Cryst. Growth 81 | 1-8 | 1987 |
| J. Aarts, W. M. Gerits & P. K. Larsen | E | Monolayer and bilayer growth on Ge(111) | J. Cryst. Growth 81 | 65-66 | 1987 |
| P. F. Fewster, J. P. Gowers, D. Hilton & C. T. Foxon | R | Structural studies of GaAs-AlAs superlattices grown by MBE | J. Cryst. Growth 81 | 120 | 1987 |
| K. Woodbridge, J. P. Gowers, P. F. Fewster, J. H. Neave & B. A. Joyce | R | RHEED studies and interface analysis of GaAs grown on Si(001) | J. Cryst. Growth 81 | 224-225 | 1987 |
| E. T. J. M. Smeets | E | Solid composition of GaAs _{1-x} P _x grown by organo-metallic vapour phase epitaxy | J. Cryst. Growth 82 | 385-395 | 1987 |
| A. H. van Ommen, M. F. C. Willemsen, A. E. T. Kuiper & F. H. P. M. Habraken (<i>Univ. Utrecht</i>) | E | Etch rate modification of Si ₃ N ₄ layers by ion bombardment and annealing | J. Electrochem. Soc. 131 | 2140-2147 | 1984 |
| E. K. Broadbent, A. E. Morgan, J. M. DeBlasi, P. van der Putte, B. Coulman, B. J. Burrow, D. K. Sadana & A. Reader | E,S | Growth of selective tungsten on self-aligned Ti and PtNi silicides by low pressure chemical vapor deposition | J. Electrochem. Soc. 133 | 1715-1721 | 1986 |
| M. Delfino, D. K. Sadana, A. E. Morgan & P. K. Chu (<i>C. Evans & Associates, San Mateo, CA</i>) | S | A study of atomic and molecular arsenic ion-implanted silicon | J. Electrochem. Soc. 133 | 1900-1905 | 1986 |

- | | | | | | |
|--|-----|---|-------------------------------------|---------|----------------|
| P. H. L. Notten & J. J. Kelly | E | Evidence for cathodic protection of crystallographic facets from GaAs etching profiles | J. Electrochem. Soc. | 444-448 | 1987 |
| | | | 134 | | |
| D. M. de Leeuw, T. Kovats & S. P. Herko | N | Kinetics of photostimulated luminescence in BaFBr:Eu | J. Electrochem. Soc. | 491-493 | 1987 |
| | | | 134 | | |
| P. K. Bachmann, P. Geittner, D. Leers & H. Wilson | A | Loss reduction in fluorine-doped SM- and high N.A.-PCVD fibers | J. Lightwave Technol. | 813-817 | 1986 |
| | | | LT-4 | | |
| P. Geittner, H. J. Hagemann, J. Warnier & H. Wilson | A | PCVD at high deposition rates | J. Lightwave Technol. | 818-822 | 1986 |
| | | | LT-4 | | |
| P. K. Bachmann, D. Leers, H. Wehr, D. U. Wiechert, J. A. van Steenwijk, D. L. A. Tjaden & E. R. Wehrhahn
(Philips Kommunikations Industrie, Nürnberg) | A,E | Dispersion-flattened single-mode fibers prepared with PCVD: performance, limitations, design optimization | J. Lightwave Technol. | 858-863 | 1986 |
| | | | LT-4 | | |
| H. J. G. Draaisma*, W. J. M. de Jonge* (*Univ. of Technol., Eindhoven) & F. J. A. den Broeder | E | Magnetic interface anisotropy in Pd/Co and Pd/Fe multilayers | J. Magn. & Magn. Mater. | 351-355 | 1987 |
| | | | 66 | | |
| A. Bhattacharyya & C. Vorst | S | Effect of addition of TCA (trichloroethane) on the electrical properties of thin oxides processed by a two-step oxidation technique | J. Phys. D | 19 | L161-L166 1986 |
| M. J. Keesman, P. H. G. Offermans & E. P. Honig | E | Silica gel formation followed by dynamic shear experiments | Mater. Lett. | 5 | 140-142 1987 |
| M. Delfino & P. K. Chu (C.Evans & Associates, San Mateo, CA) | S | Donor creation during oxygen implanted buried oxide formation | Mater. Res. Soc. Symp. Proc. | 53 | 245-250 1986 |
| G. J. Campisi, H. B. Dietrich (Naval Res. Lab., Washington, DC), M. Delfino & D. K. Sadana | S | High dose implantation of nickel into silicon | Mater. Res. Soc. Symp. Proc. | 54 | 747-752 1986 |
| W. G. M. van den Hoek | S | The etch mechanism for AL ₂ O ₃ in fluorine and chlorine based RF dry etch plasmas | Mater. Res. Soc. Symp. Proc. | 68 | 71-78 1986 |
| W. G. M. van den Hoek | S | Characterization of plasma-enhanced chemical vapour deposition of silicon-oxy-nitride | Mater. Res. Soc. Symp. Proc. | 68 | 335-342 1986 |
| J. M. DeBlasi, D. K. Sadana & M. H. Norcott | S | Interfacial tunnel structures in CMOS source/drain regions following selective deposition of tungsten | Mater. Res. Soc. Symp. Proc. | 71 | 303-308 1986 |
| S. Aronowitz* & M. Delfino (*Fairchild Res. Center, Palo Alto, CA) | S | Laser activated glass flow modeling | Mater. Res. Soc. Symp. Proc. | 71 | 479-485 1986 |
| A. L. J. Burgmans | E | IC-technologie op submicron-niveau; een CMOS-proces | Ned. T. Natuurk. A | 52 | 99-102 1986 |
| F. M. Klaassen | E | Fysica van kleine IC-componenten | Ned. T. Natuurk. A | 52 | 107-112 1986 |
| P. C. Zalm, A. W. Kolfshoten, F. H. M. Sanders & P. Vischer | E | Surface processes in ion-induced etching | Nucl. Instrum. & Methods Phys. Res. | B18 | 625-628 1987 |
| H. J. Ligthart, E. Gerritsen, P. J. van den Kerkhoff, S. Hoekstra, R. E. van de Leest & E. Keetels | E | Aluminium implantations in copper | Nucl. Instrum. & Methods Phys. Res. | B19/20 | 209-212 1987 |
| A. H. van Ommen, M. F. C. Willemssen & R. A. M. Wolters | E | The effect of ion beam mixing on MoSi ₂ formation | Nucl. Instrum. & Methods Phys. Res. | B19/20 | 742-745 1987 |
| C. W. J. Beenakker | E | Evolution of two-dimensional soap-film networks | Phys. Rev. Lett. | 57 | 2454-2457 1986 |
| M. A. Brummell*, R. J. Nicholas*, M. A. Hopkins* (*Clarendon Lab., Oxford), J. J. Harris & C. T. Foxon R | R | Modification of the electron-phonon interactions in GaAs-GaAlAs heterojunctions | Phys. Rev. Lett. | 58 | 77-80 1987 |
| F. Berz & J. A. Morice | R | Accuracy of the transmission coefficient across parabolic barriers as obtained from a generalized WKB approach | Phys. Stat. Sol. b | 139 | 573-582 1987 |
| J. J. Goedbloed | E | Electromagnetic compatibility | Phys. Technol. | 18 | 61-67 1987 |
| H. W. A. M. Rompa, R. Eppenga & M. F. H. Schuurmans | E | Ab initio determinations of oscillator strengths of metals and semiconductors from ASW | Physica | 145B | 5-15 1987 |

W. F. Druyvesteyn, A. J. M. Kaizer, D. J. Verschuur* & D. de Vries* (*Univ. of Technol., Delft)	E	Wigner representation of loudspeaker responses in a living room	Proc. AES Conv., London 1987 (preprint)	12pp.	1987
V. Pauker	L	GaAs monolithic microwave active gyrator	Proc. GaAs IC Symp., Grenelefe, FL, 1986	82-85	1986
D. Meignant (RTC Compelec, Limeil-Brévannes), E. Delhay & M. Rocchi	L	A 0.1 - 4.5 GHz, 20 mW GaAs prescaler operating at 125 °C	Proc. GaAs IC Symp., Grenelefe, FL, 1986	129-132	1986
T. Ducourant, M. Binet, J.-C. Baelde, C. Rocher, J.-M. Gibereau	L	3 GHz, 150 mW, 4 bit GaAs analogue to digital converter	Proc. GaAs IC Symp., Grenelefe, FL, 1986	209-212	1986
P. A. Gough, M. R. Simpson & V. Rumennik	N	Fast switching lateral insulated gate transistor	Proc. IEDM 86, Los Angeles, CA, 1986	218-221	1986
S. Mukherjee, C. J. Chou, K. Shaw, D. McArthur & V. Rumennik	N	The effects of SIPOS passivation on DC and switching performance of high voltage MOS transistors	Proc. IEDM 86, Los Angeles, CA, 1986	646-649	1986
C. A. M. Mulder & G. Frens	E	The incorporation of fluorine in vitreous silica	Proc. 1st Int. Workshop on Non-crystalline solids, San Feliu de Guixols 1986	259-265	1986
J. Hightower, M. Blanco, M. Cagan & K. Monahan	S	An application of statistical experimental design: comparison of process latitudes for photoresist and cel	Proc. Kodak Microelectron. Seminar, San Diego, CA, 1985	39-49	1985
U. Schiebel, W. Hillen & T. Zaengel	A	Image quality in selenium-based digital radiography	Proc. SPIE 626	176-184	1986
P. J. Severin	E	Passive components for multimode fibre-optic networks	Proc. SPIE 630	98-109	1986
M. Erman, N. Vodjdani, P. Jarry, P. Stéphan, J. L. Gentner & C. Guedon	L	III-V semiconductor waveguides and phase-modulators: the localized vapor phase epitaxy approach	Proc. SPIE 651	75-82	1986
W. S. Newman & N. Hogan (Massachusetts Inst. of Technol., Cambridge, MA)	N	The optimal control of balanced manipulators	Robotics 3	177-184	1987
A. G. Dirks & J. J. van den Broek	E	Precipitation from metastable solid solutions in aluminum rich AL-V thin films	Scr. Metall. 21	175-180	1987
D. E. Lacklison & P. Capper	R	Hall effect measurements on Bridgman-grown $Cd_xHg_{1-x}Te$ and their analysis	Semicond. Sci. Technol. 2	136-144	1987
A. F. Deutz*, H. B. Brom*, H. Deelen*, L. J. de Jongh*, W. J. Huiskamp* (*Univ. Leiden) & K. H. J. Buschow	E	Induced-moment ferromagnetic ordering in the singlet ground state compound $TmNi_2$	Solid State Commun. 60	917-921	1986
P. W. J. M. Boumans & J. J. A. M. Vrakking	E	Detection limits of about 350 prominent lines of 65 elements observed in 50 and 27 MHz inductively coupled plasmas (ICP): effects of source characteristics, noise and spectral bandwidth—'Standard' values for the 50 MHz ICP	Spectrochim. Acta 42B	553-579	1987
C. W. T. Bulle-Lieuwma & P. C. Zalm	E	Suppression of surface topography development in ionmilling of semiconductors	Surf. & Interface Anal. 10	210-215	1987
E. E. Havinga & L. W. van Horssen	E	Influence of chain length and structural order on the electrical conductivity of poly- <i>p</i> -phenylenes heavily doped with potassium	Synth. Met. 17	623-628	1987
A. G. Tangena	E	Tribology of thin film systems	Thesis, Eindhoven	1-129	1987

Contents of Philips Telecommunication and Data Systems Review 45, No. 1, 1987

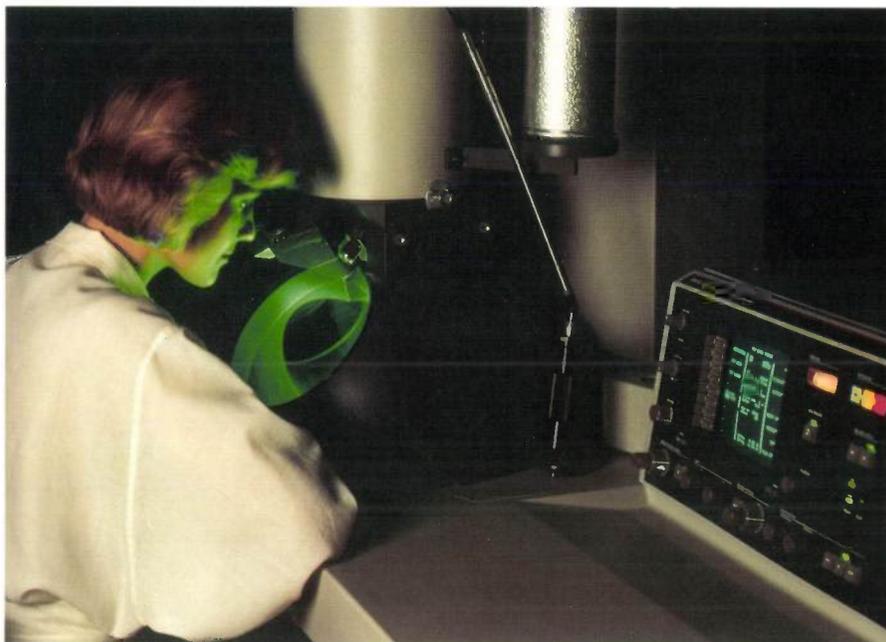
A. Brinkman: Transmatron: Electronics improving public transport (pp. 2-15)

J. Goris & J. Veldman: The Rotterdam Vessel Traffic Management System (pp. 16-33)

C. Viret: Trends in civil avionics for the near future (pp. 34-41)

The microprocessor-controlled CM12/STEM scanning-transmission electron microscope

U. Gross, F. J. M. Mescher and J. C. Tiemeijer



Continuous technological innovation in electron optics, mechanical engineering, electronics and vacuum engineering have earned Philips a leading position in the market for transmission electron microscopes. A recent invention is the Twin objective lens, which makes it possible to switch quickly from a transmission image to an image produced by scanning the specimen with the electron beam. In the latest microscope, the CM12/STEM scanning-transmission electron microscope, ease of operation has been considerably increased by the addition of microprocessor control. This microscope gives images that show that the microscope has a line resolution of 0.14 nm in the transmission mode and a point resolution of 1 nm in the scanning mode.

Introduction

The classical transmission electron microscope, as developed in the early thirties by Ernst Ruska and others, is comparable, in the geometry of the ray pattern, with a photographic enlarger or a slide projector. The extremely thin specimen is illuminated from

an electron source by means of one or more condenser lenses. The subsequent lenses produce a magnified image on a fluorescent screen, which converts the electron image into a visible image. The electron image can be recorded on photographic material.

The motivation for using electrons instead of visible light is of course to improve the resolution. Electrons at an energy of 120 kV are equivalent to electro-

Dr U. Gross, Ir F. J. M. Mescher and Ing. J. C. Tiemeijer are with the Philips Industrial and Electro-acoustic Systems Division, Eindhoven.

magnetic radiation at a wavelength λ of about 3.4×10^{-3} nm.

The wavelength λ is given by the relation

$$\lambda = \frac{h}{\sqrt{2em_0V^*}},$$

where h is Planck's constant, e the charge and m_0 the rest mass of the electron, and V^* a corrected value for the accelerating voltage V , taking account of the relativistic change of mass. V^* is given by the relation

$$V^* = V \left(1 + \frac{e}{2m_0c^2} V \right),$$

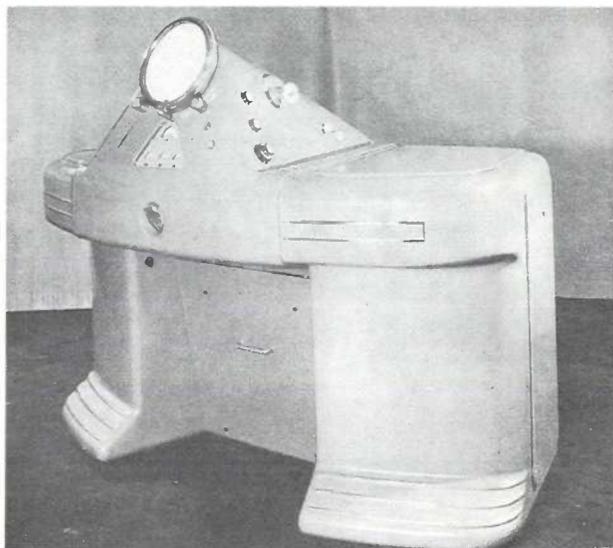
where c is the velocity of light. A useful expression, which does not take the relativistic change of mass into account and is valid for accelerating voltages lower than about 100 kV, is:

$$\lambda = \frac{1.23}{\sqrt{V}} \text{ nm},$$

where the accelerating voltage V is in volts.

The wavelength in an electron microscope is thus about 100000 times smaller than in an optical microscope. The resolution of an electron microscope, however, is not 100000 times better [1]. This is because the resolution is not merely proportional to the wavelength; it is also inversely proportional to the aperture angle of the objective lens. In an electron microscope the aperture angle has to be very much smaller than that of an optical microscope. The main reason for this is the need to minimize spherical aberration. The magnetic lenses that are combined to form the column of a modern electron microscope always have a positive spherical aberration, which means to say that the lenses are stronger for peripheral rays than for rays at the centre. Spherical aberration, unlike most other lens errors, cannot therefore be corrected.

It is perhaps interesting to illustrate the structure of a classical electron microscope by referring to a photograph and diagrams of the first microscope that Philips put on the market: the EM 100 of 1950; see *fig. 1* [2]. This microscope had one condenser lens and four image-forming lenses. These were a considerable improvement at the time, because the addition of a 'diffraction lens' made it possible to observe an ordinary image and an electron-diffraction pattern one after the other. This diffraction pattern related to a very small part of the specimen, selected by means of an SA diaphragm (SA for 'selected area'). The principle was due to Prof. J. B. le Poole and was for a long time the only way in which information about local crystal structures could be obtained in a transmission electron microscope [1].



a

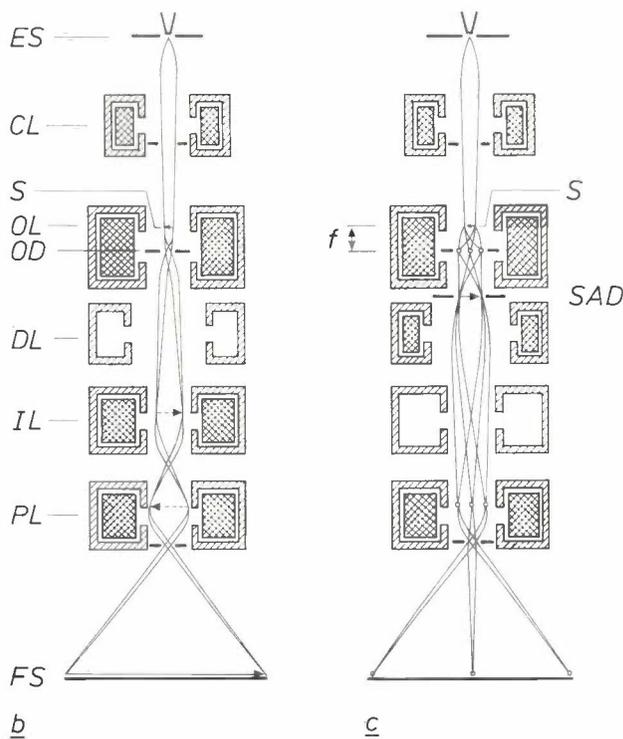


Fig. 1. *a*) The EM 100 transmission electron microscope, which Philips put on the market in 1950. (This microscope had a horizontal column to facilitate observation of the fluorescent screen.) *b*) The ray diagram for observing a highly magnified image of the specimen. *ES* electron source. *CL* condenser lens. *OL* objective lens. *DL* diffraction lens. *IL* intermediate lens. *PL* projector lens. *FS* fluorescent screen. The lenses with the cross-hatched coils are in operation. *S* specimen. *OD* objective diaphragm. *c*) The ray diagram for observing an electron-diffraction pattern. *SAD* selected-area diaphragm. *f* focal length of the objective.

Fig. 1b shows the important part played by the objective diaphragm *OD*, which is located at the back focal plane of the objective lens. It intercepts electrons that are scattered in the specimen, and is a useful aid for varying contrast in the image. **Fig. 1c** shows

that, with parallel illumination of the specimen, a diffraction pattern is formed at the back focal plane of the objective. To obtain an image of this diffraction pattern the objective diaphragm must be removed. This diaphragm does have to be used, however, for producing an ordinary image. If the aperture of the diaphragm coincides with the image of the electron source formed by the electrons transmitted through the specimen, not scattered in the specimen (fig. 1*b*), the image formed on the screen is called a bright-field image. This image, which is the one normally used, has contrasts different from those in a 'dark-field' image. A dark-field image is obtained by tilting the illuminating beam in such a way that the aperture of the objective diaphragm coincides with the image of the electron source formed by electrons that have been scattered over a defined angle in the specimen.

Philips transmission electron microscopes have undergone continuous improvement in the last 30 years^[3]. While there have been a number of more gradual modifications and refinements, the emergence of another type of electron microscope, the scanning electron microscope, has considerably accelerated the evolution of transmission electron microscopes. Experience gained with scanning electron microscopes^[4] — developed mainly for scanning the surface of large specimens with a narrow electron beam — and with electron microprobes^[5] — developed primarily for wavelength-dispersive analysis of X-radiation from large specimens with a crystal spectrometer — has led to the realization that additional information could be obtained in a transmission electron microscope by scanning the specimen.

In the seventies this resulted in the development of a 'STEM' accessory unit for the Philips transmission electron microscope (STEM is an acronym for Scanning-Transmission Electron Microscope). With this electronic unit the electron beam can be made to 'write' a rectangular raster across the specimen. The electron gun of a cathode-ray tube in the STEM unit is controlled by a signal from a detector, and produces an image on its screen in which the contrast is obtained by techniques different from the conventional methods. Fig. 2 shows the detectors that are available for these techniques. The signals from the different detectors originate from:

- high-energy electrons back-scattered from the specimen (*BSD*),
- low-energy electrons produced by secondary emission in the specimen (*SED*),
- X-rays emitted by excited atoms in the specimen (*EDX*),
- non-scattered transmitted electrons (*BFD*),
- scattered transmitted electrons (*DFD*),

- transmitted electrons of specific, selected energy (*EELS*).

It should be noted that this evolution from a classical transmission electron microscope (TEM) to a scanning-transmission electron microscope (STEM) does not mean that the 'ordinary' scanning electron microscope (SEM) has become redundant. An SEM has its own specific range of applications, for the investigation of much larger and thicker specimens. Very often the SEM can also be used for wavelength-dispersive analysis of X-radiation, whereas a STEM

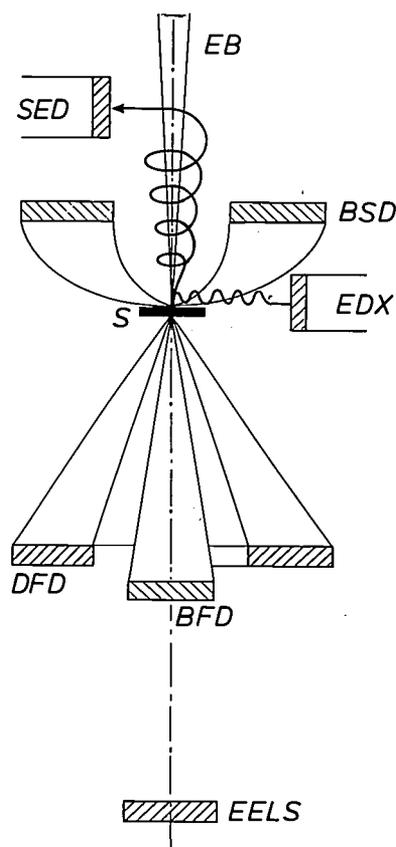


Fig. 2. Schematic representation of the detectors around the specimen and in the projection chamber. The detectors are mainly used when the instrument is used as a scanning-transmission electron microscope (STEM). *EB* electron beam. *SED* secondary electron detector. *BSD* back-scattered electron detector. *EDX* detector for energy-dispersive X-ray analysis. *BFD* bright-field detector, for non-scattered transmitted electrons. *DFD* dark-field detector, for scattered transmitted electrons. *EELS* electron energy loss spectroscopy detector, for transmitted electrons sorted by energy; see also fig. 9.

- [1] J. B. le Poole, A new electron microscope with continuously variable magnification, Philips Tech. Rev. 9, 33-45, 1947/48.
- [2] A. C. van Dorsten, H. Nieuwdorp and A. Verhoeff, The Philips 100 kV electron microscope, Philips Tech. Rev. 12, 33-51, 1950/51.
- [3] C. J. Rakels, J. C. Tiemeijer and K. W. Witteveen, The Philips electron microscope EM 300, Philips Tech. Rev. 29, 370-386, 1968.
- [4] W. Kuypers and J. C. Tiemeijer, The Philips PSEM 500 scanning electron microscope, Philips Tech. Rev. 35, 153-165, 1975.
- [5] M. Klerk, The electron microprobe, Philips Tech. Rev. 34, 370-374, 1974.

can only perform an energy-dispersive analysis with a solid-state detector. In scanning a specimen, however, a STEM will give a better resolution than that of most SEMs, since it has a finer beam, which is not broadened so much in the thin specimen.

The first Philips transmission microscopes of the EM 400 series appeared on the market in 1975. These microscopes have a column that facilitates the addition of a STEM accessory. The column has two condenser lenses and five image-forming lenses. The large number of lenses provided more facilities for eliminating chromatic magnification errors, radial distortion and image rotation. In addition, a cross-over point could be created for all the geometric ray patterns at a fixed position just beneath the final lens. It was then possible to use a diaphragm with an aperture of 200 μm cross-section, to separate the column from the projection chamber, where the fluorescent screen and the photographic material are located. In this arrangement an oil-free ion-getter pump can maintain a clean vacuum in the column and in the emission chamber, containing the electron source. This, and the use of

metal bellows instead of rubber rings, meant that a very clean environment could be created for the specimen and the filament. Consequently there is very little contamination of the specimen by hydrocarbons and the filament has a long life. Another advantage of these microscopes is that they no longer have to be mechanically aligned, since the alignment is done with correction coils.

In 1978 the Twin objective lens was introduced [6]. A microscope with this objective can easily be switched from a broad beam for normal images to a narrow beam for scanning images and analyses. Switching between broad beam and narrow beam does not change the heat flow in the objective, so that there is hardly any change in its dimensions. Nor is it necessary to change the polepieces. The objective has so much space between the polepieces that there is room for detectors and for a maximum specimen tilt through an angle of $\pm 60^\circ$. Since the distance between the specimen and the X-ray detector is therefore small, more of the X-radiation emitted by the specimen can be intercepted and measured. The large tilt angle is useful for crystallographic electron-diffraction analyses.

The Twin objective is a modified immersion objective of the Riecke and Ruska type, with the specimen half-way between the polepieces. The specimen is located at the common focal plane of the two lenses formed by the magnetic field between the polepieces. Added to the immersion objective is a 'minilens', which acts as a third condenser lens. A special magnetic configuration allows the minilens to be switched off optically by simply reversing the current in the associated coil. This explains the high stability of the Twin objective. By making the mini-condenser lens optically inactive or active the minimum cross-section of the 'spot' at the specimen can be made either 1.5 nm or 0.04 μm . The smaller spot is called the 'nanoprobe' spot and the larger one the 'microprobe' spot.

With its Twin objective, a STEM unit and various detectors a Philips transmission-electron microscope can be used for a wide range of analyses. It can also produce transmission and scanning images of high resolution. *Fig. 3* shows the EM 420 microscope complete with all the accessory units for various purposes. Although the microscope is highly versatile and is virtually a complete laboratory in itself, its use requires a thorough knowledge of electron optics and analysis techniques. The complexity of an electron microscope with all these additional facilities is clear from the number of controls — 250 in all.

Recently Philips have introduced the scanning-transmission electron microscope type CM12/STEM,



Fig. 3. The EM 420 transmission electron microscope complete with all accessories. The Dewar vessel immediately to the right of the column contains liquid nitrogen for cooling the energy-dispersive X-ray detector (*EDX* in *fig. 2*); below this and to the right is the electronic unit for displaying the X-ray spectrum. On the extreme right is the STEM unit. This, like most other accessory units, is integrated into the latest electron microscope, CM12/STEM; see *fig. 4*.



Fig. 4. *a*) The CM12/STEM scanning-transmission electron microscope, which is much easier to operate than the microscope in fig. 3. On the left in the lower panel is the control screen for communication between user and the microprocessor — the heart of the instrument. The upper control panel contains two monitor screens, which display the STEM images. On the extreme right is the 'videoscope', an accessory that displays the brightness variation for one image line. *b*) Close-up of the control screen, with 'soft keys' on either side; the user uses these to communicate his requirements to the control software.

shown in the photograph of *fig. 4a*. In this instrument all the features of the EM 420 in *fig. 3* are *integrated*, and all the electronic circuits are controlled by a microprocessor. This microprocessor guides the user to the various features available — a very useful facility. The microprocessor first offers the user, via the control screen (see *fig. 4b*), a choice of the modes he can work with, where a mode is a setting with a particular ray pattern. All the necessary operations are then carried out in an interactive dialogue with the user. The user requires less expertise in electron optics, and can now concentrate more closely on his investigations. In each mode the values of currents in the correction coils can remain stored in the microprocessor's memory. If the user wishes to depart from the programmed settings, he can adjust the lens currents to

suit his own requirements. Since there are no longer any fixed electrical connections between the controls and the electronics, there are far fewer manual controls. Consequently the CM12/STEM electron microscope is almost as easy to operate as the EM 100 in *fig. 1*, but it offers many more facilities and the resolution is incomparably better.

We shall now take a closer look at the electron-optical aspects of the instrument, dealing first with the Twin objective and then with the various modes. Then we shall touch briefly on the microprocessor control system. Finally, we shall demonstrate the features of the CM12/STEM with some micrographs of very dissimilar specimens.

[6] The Twin lens is patented; the number of the patent in the United States is 4306149.

The electron optics

Fig. 5 shows a schematic cross-section of the column with emission chamber and projection chamber. The illumination system consists of the electron gun with filament and (electrostatic) Wehnelt lens, two condenser lenses, the minilens and the upper half of the objective lens. The last two lenses also act as a condenser lens and form part of the Twin objective,

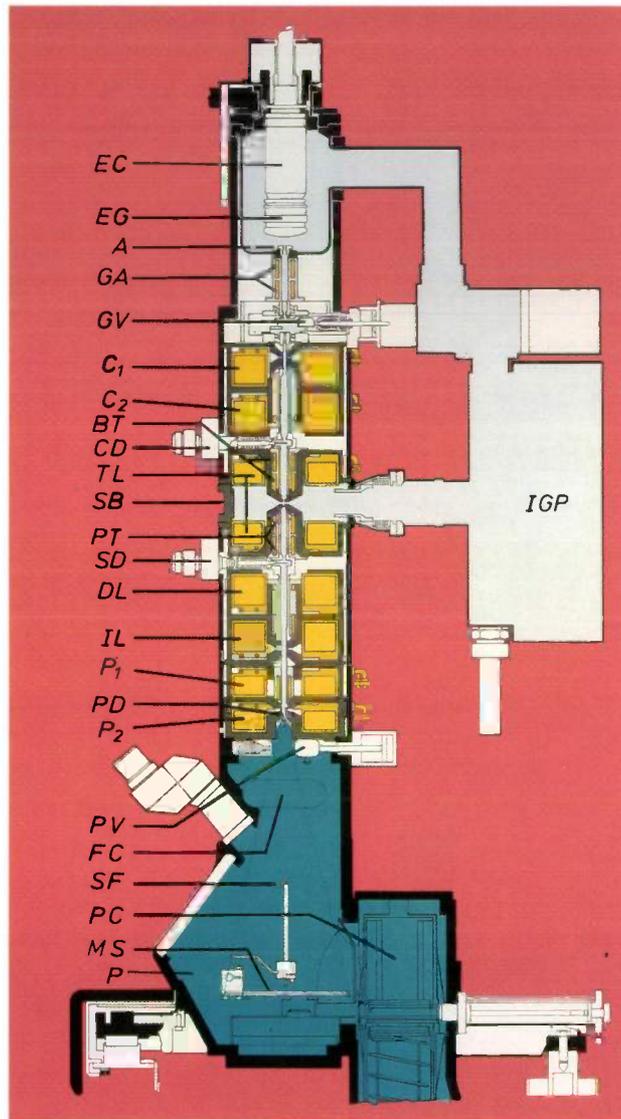


Fig. 5. Cross-section of the microscope column. *EC* emission chamber. *EG* electron gun. *A* anode. *GA* deflection coils for electron-gun alignment. *GV* valve for separately venting the emission chamber. *C₁* first condenser lens. *C₂* second condenser lens. *BT* beam-deflection coils above the specimen. *CD* condenser-diaphragm holder. *TL* Twin-objective lens. The mini-condenser lens is not shown (see fig. 7). *SB* block for accommodating the 'goniometer' (not shown); the goniometer permits the specimen in the specimen holder to be displaced. *PT* deflection coils below the specimen. *SD* SA diaphragm holder. *DL* diffraction lens. *IL* intermediate lens. *P₁* first projector lens. *PD* partition diaphragm. *P₂* second projector lens. *PV* valve for separately venting the projection chamber. *FC* film camera. *SF* fluorescent screen, for accurate focusing with a binocular magnifier. *PC* plate camera. *MS* fluorescent screen for observing the image through lead-glass windows in the projection chamber *P*, see also the title photograph. *IGP* ion-getter pump.

which we shall return to presently. The image-forming system consists of the lower half of the objective lens, the diffraction lens, the intermediate lens and two projector lenses.

At various positions in the column there are correction coils, which correct the beam for errors due to geometric imperfections and material inhomogeneities. The correction coils consist of:

- stigmator coils, for eliminating beam astigmatism, and
- deflection coils; these correct errors in deflection but their main function is to alter the direction of the beam so that it writes a rectangular raster across the specimen, for example.

The stigmator coils are oriented radially and arranged in sets of eight with an angle of 45° between them, as shown in fig. 6*a*. The deflection coils are combined in groups of four at 90° spacing. Two of these groups placed one above the other can produce a wide range of deflections in two perpendicular planes through the optical axis, as illustrated in fig. 6*b*, *c* and *d*. Two coils are shown from each group of four; the beam deflection is only shown in one plane through the optical axis. The correction coils free the user from the need for mechanical adjustment of the lenses.

At various positions in the column there are diaphragms that are adjustable in position and interchangeable, with a choice of four different sizes. These are:

- the condenser diaphragm, in the illumination system;
- the objective diaphragm, just beneath the specimen in the back focal plane of the lower half of the objective lens; and
- the SA diaphragm, in the image-forming system.

The column also contains fixed diaphragms, whose aperture cannot be changed. There is also a partition diaphragm (*PD* in fig. 5) of 200 μm cross-section, as mentioned earlier, which separates the vacuum in the column from that in the projection chamber. Water vapour from the photographic material and oil vapour from the diffusion pump and backing pump, which are connected to the projection chamber, are thus prevented from contaminating the specimen and reducing the life of the filament. The vacuum in the emission chamber and specimen chamber are maintained by a 'dry' ion-getter pump, i.e. one that requires no oil for its operation.

The specimen is placed in a specimen holder, which is slid through a vacuum lock into the specimen stage or 'goniometer'. The specimen can then be displaced in two directions at right angles to each other over an area of 2mm by 2mm and also tilted through a maximum angle of ± 60° about an axis in the plane

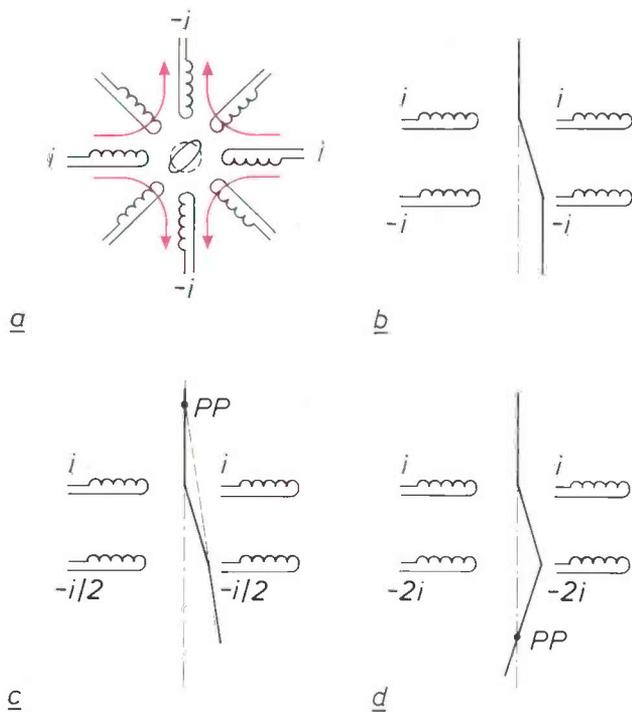


Fig. 6. The correction coils. *a*) A set of eight stigmator coils, used for correcting beam astigmatism. The diagram shows how an elliptical cross-section of the beam is changed into a circular cross-section by currents i and $-i$ in four of the eight coils. The magnetic lines of force are shown red. *b*) Deflection coils that not only correct errors in the beam direction but also cause the illuminating beam to describe a line across the specimen, as illustrated here schematically. (The deflection of the beam shown corresponds to a direction of the field and hence of the coils perpendicular to the plane of the drawing; the coils as shown are thus rotated through 90° ; this also applies to *c* and *d*.) *c*) In this mode of energization the beam rotates about an imaginary pivot point PP above the coils. *d*) The same situation, but now with the pivot point below the coils. In certain modes, e.g. the dark-field mode with conical illumination, PP is in the specimen.

of the specimen. The special feature here is that this axis intersects the optical axis in every position of the specimen. This is referred to as a 'eucentric' movement of the specimen. The advantage is that the tilt axis always appears to go through the centre of the image on the user's screen. In addition to the standard specimen holder there are special specimen holders in which the specimen can be rotated, stretched, cooled or heated, and can always be tilted.

The Twin objective lens

In modern electron microscopes the specimen is not situated above the objective lens but inside it. The lens system is known as an immersion objective. In the symmetrical immersion objective proposed by Riecke and Ruska the specimen is situated exactly half-way between the polepieces^[7]. Such a lens is a 'condenser objective'^[8], with the part of the magnetic field above the specimen acting as a condenser lens and the part of the field below it as an image-forming lens. In an analogy with ordinary optics we can think

of a strong biconvex glass lens, which is cut in half perpendicular to the optical axis. The two identical halves are then moved away from each other through a distance of twice their focal length, and the specimen can then be considered to be located in the common focal plane between them.

The Twin objective lens is a symmetrical immersion objective to which a condenser lens, the minilens, has been added; see fig. 7^[9]. The minilens is situated im-

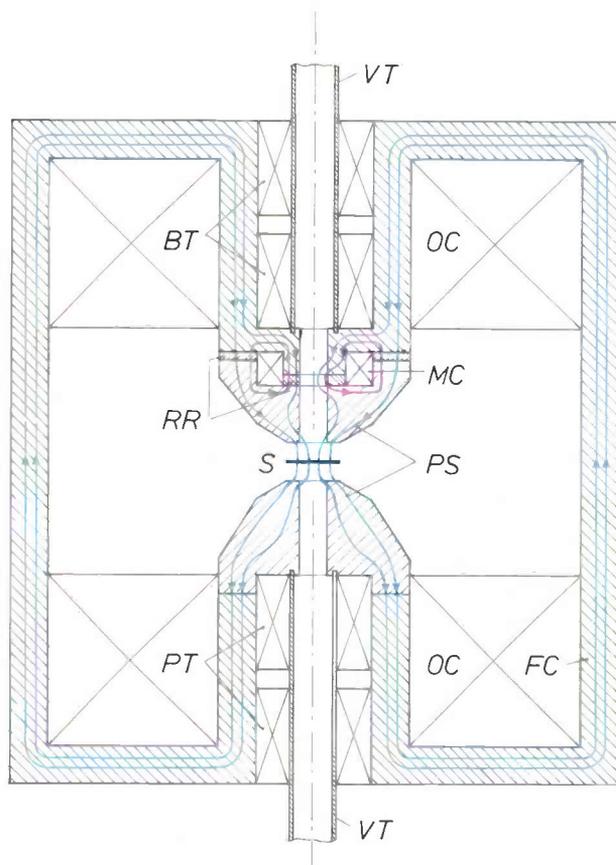


Fig. 7. The Twin objective lens^[6]. PS polepieces. FC ferromagnetic core, including the coils OC . BT and PT deflection coils above and below the specimen S . MC mini-condenser coil. RR reluctances. VT stainless-steel tubes, enclosing the vacuum for the electron beam. The goniometer (not shown) is located between the coils OC ; there is therefore another partition between the vacuum and these coils. In the left-hand half of the figure the minilens is not active; in the right-hand half it is. The blue magnetic lines of force are produced by a current in the coils OC . Because of their curvature, these lines of force form a positive lens both below and above the specimen; see also fig. 8. The red lines of force are the result of a current in the coil MC . In the right-hand half the red lines of force and the blue lines of force together form the mini-condenser lens, because of their curvature. In the left-hand half the current in MC is reversed, so that the red and blue lines of force counteract each other, causing the mini-condenser lens to become optically inactive here.

[7] W. D. Riecke and F. Ruska, A 100 kV transmission electron microscope with single-field condenser objective, Proc. Int. Cong. on Electron Microscopy Vol. 1, Kyoto 1966, pp. 19-20.
 [8] J. R. A. Cleaver, The choice of polepiece shape and lens operating mode for magnetic objective lenses with saturated polepieces, Optik 57, 9-34, 1980.
 [9] K. D. van der Mast, C. J. Rakels and J. B. le Poole, A high quality multipurpose objective lens, Proc. Cong. on Electron Microscopy Vol. 1, The Hague 1980, pp. 72-73.

mediately above the polepieces of the objective and can be switched on or off (i.e. made optically active or inactive), depending on the required mode of operation. Normally, switching a lens off interrupts the current through the lens, which means that the heat-flow conditions are changed, and hence the dimensions. Near the specimen, however, the dimensions have to be highly stable, so that changes in heat flow are not permitted. The minilens is switched on and off optically without causing fluctuations in heat flow, by means of a special magnetic configuration in which the minilens is de-activated by *reversing* the current in the appropriate coil.

In this magnetic configuration the magnetic flux of the objective lens is split in such a way that about 50% of the action of the minilens is based on this flux and 50% on the flux from its own coil. In the left-hand half of fig.7 the two fluxes in the 'air gap' of the minilens oppose one another. The resultant flux in the air gap is therefore zero and the minilens is optically inactive. In the situation shown in the right-hand half of fig.7 the current in the coil of the minilens has been reversed, and the two fluxes in the air gap reinforce, so that the minilens is optically active. The flux of the objective lens is split by including an additional reluctance around the coil of the minilens. This reluctance and the 'air gap' of the minilens are formed by rings of non-magnetic material, such as aluminium or copper.

The unmodified symmetrical immersion objective is particularly suitable for making very small spots on the specimen. *Fig. 8a* shows the ray diagram for the smallest spot, the nanoprobe spot, used for scan images and for analyses. Here the electrons move parallel to one another between the second condenser lens (C_2) and the condenser part (O_u) of the objective. An incidental advantage of this objective is that there is a relatively large amount of space between the polepieces. A disadvantage of a symmetrical immersion objective without an additional condenser lens, however, is that when the instrument is used as a TEM the objective is less suitable for parallel illumination and defocused illumination, i.e. with the spot on the specimen slightly out of focus.

Fig. 8b shows that these disadvantages have been overcome by the addition of the minilens (C_m). When this lens is energized so that its focus coincides with that of the condenser part (O_u) of the objective, the result is a 'telecentric' optical system. The special feature of such a system is that a parallel incident beam leaves the system still as a parallel beam. *Fig. 8b* shows that, in spite of the presence of the converging lens O_u just above the specimen, parallel illumination of the specimen is nevertheless obtained. If the con-

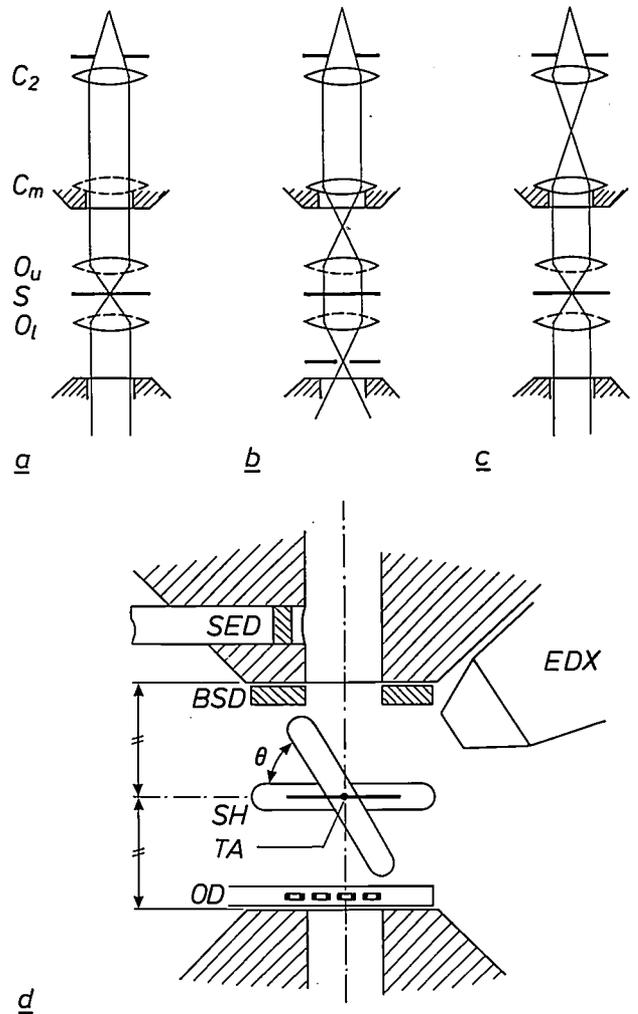


Fig. 8. The main operating conditions of the Twin objective lens. C_2 second condenser lens. C_m mini-condenser lens. O_u condenser lens formed by the magnetic field above the specimen S . O_l image-forming lens produced by the field below the specimen. *a*) Ray diagram when the smallest possible spot — the nanoprobe spot — is formed on the specimen. C_m is not in operation. *b*) Ray diagram for parallel illumination of the specimen. C_m and O_u together form a telecentric system. *c*) Ray diagram when C_2 is energized in another way, producing a larger spot — the microprobe spot — on the specimen. *d*) The detectors and holders that have to be accommodated between and in the polepieces near the specimen; see also caption to fig. 2. SH specimen holder, which can be tilted through an angle θ about an axis TA perpendicular to the plane of the drawing. OD holder for four different objective diaphragms.

denser lens C_2 is energized in another way a larger spot than the nanoprobe spot is produced on the specimen — the microprobe spot; see fig. 8c. *Fig. 8a, b* and *c* show three important modes of operating with the Twin lens. The ray diagrams are limiting cases for a large number of variations, e.g. with a defocused spot or with incompletely parallel illumination of the specimen.

The specimen is thus placed half-way between the polepieces of the objective lens, since with the easily switched mini-condenser lens a strong lens immediately above the specimen is no handicap. *Fig. 8d* shows that the relatively large free distance above and below

the specimen is used for accommodating a holder for four different objective diaphragms, the detectors for back-scattered electrons and emitted X-rays, and also for tilting the specimen holder through the largest possible angle. The detector for secondary electrons, which has an electrode at a potential of 10 kV to collect electrons with an energy of less than 50 eV, is accommodated in a radial hole in the upper polepiece. The specimen must be tilted:

- for investigating crystal structures with the aid of diffraction patterns;
- for making stereo micrographs, i.e. micrographs with opposite angles of tilt, but otherwise identical;
- for making series of micrographs with different tilt angles, from which a computer can reconstruct three-dimensional images, e.g. of biological macromolecules; and
- for obtaining maximum signals from the X-ray detector.

There is sufficient space for the X-ray detector to be mounted close to the specimen. The solid angle at which the detector 'sees' the specimen is therefore relatively large, and this also means that a strong signal is obtained from the X-ray detector, so that selective energy-dispersive analysis of the specimen is possible.

The designers of the objective described here were confronted with two conflicting requirements: they could have either a large space between the polepieces or a high resolution. Scaling down the polepieces reduces the space available, but gives a higher resolution because the spherical and chromatic aberrations will be smaller. On the other hand, a relatively large space between the polepieces gives a slightly lower resolution. The dilemma was resolved by making two versions of the objective lens, the Twin lens and the Super-Twin lens. In the first version the emphasis is on a large tilt angle, in the second on a high resolution with high magnification. *Table I* lists the focal distance, spherical and chromatic aberration, point and line resolution, minimum spot diameter and tilt angle for the two objectives.

Table I. Focal length f , constants C_s for spherical aberration and C_c for chromatic aberration^[10], minimum spot diameter d , point resolution R_p , line resolution R_l and maximum tilt angle θ for the Twin and Super-Twin objectives.

	Twin	Super-Twin
f	2.7 mm	1.7 mm
C_s	2.0	1.2
C_c	2.0	1.2
d	2.0 nm	1.5 nm
R_p	0.34 nm	0.30 nm
R_l	0.20 nm	0.14 nm
θ	$\pm 60^\circ$	$\pm 15^\circ$

Operating modes

With the versatile Twin lens and the deflection coils above and below the specimen the CM12/STEM microscope offers a large number of optical settings or modes. In this section we shall look at the most important of these modes, classified by the nature of the energization of the deflection coils.

The illumination of the specimen corresponds in general to fig. 8*a*, *b* or *c*, i.e. illumination with the nanoprobe spot, a parallel beam or the microprobe spot. The spots may be slightly defocused to cover a larger area of the specimen. We note in passing that the deflection coils above the specimen can also be used as a focusing aid with a standard TEM image, by energizing the upper coils with a square-wave voltage and situating the pivot point PP in the specimen, as shown in fig. 6*d*. The user then sees two images, and the image-forming system is correctly focused when the two images coincide. This 'wobbling' illumination procedure considerably simplifies the focusing^[10].

Modes with no dynamic energization of the deflection coils

These modes include the conventional TEM modes. They consist of a high-magnification mode TEM-HM and a low-magnification mode TEM-LM. In the TEM-HM mode the objective is always energized and most of the other image-forming lenses are energized; the image is focused by varying the current through the objective. In the TEM-LM mode the objective is only weakly energized, so that really the diffraction lens is the first image-forming lens. The image is therefore focused by varying the current through the diffraction lens. With central illumination a bright-field image is obtained, and with tilted illumination a dark-field image. In the latter case the electrons that are not scattered in the specimen are intercepted in TEM-HM by the objective diaphragm. In tilted illumination the deflection coils above the specimen are therefore *statically* energized.

An image of the electron diffraction pattern at the back focal plane of the first image-forming lens can be obtained both in the TEM-HM and in the TEM-LM mode. The other image-forming lenses are then energized so as to produce an image corresponding to this focal plane on the fluorescent screen. In the TEM-HM diffraction mode the SA diaphragm (see fig. 5) can be used to select a part of the ordinary image of the specimen for closer investigation of the crystal structure.

Element analysis can be carried out in the nanoprobe mode. Since the spot at the specimen is extremely small (fig. 8*a* and Table I) it is then possible to ana-

^[10] C. E. Hall, Introduction to electron microscopy, McGraw-Hill, New York 1953.

lyse the X-radiation originating from a very small part of the specimen. Here the energy-dispersive X-ray detector *EDX* is used; see fig. 8*d*. The energy distribution of the X-radiation permits qualitative and quantitative analysis of the elements in the specimen.

Sensitive specimens can be analysed in the low-dose mode. This can be done for ordinary micrographs, diffraction micrographs and even with STEM, which will be discussed below. The procedure is as follows. At low electron density and low magnification an interesting part of the specimen is identified. The specimen is not displaced further and with high magnification and at the normal electron density — which might cause local damage to the specimen — a part *next to* the interesting part is observed; the image is then focused. During this procedure the deflection coils above and below the specimen are energized statically (fig. 6*b*), in such a way that the lateral displacement of the beam above the specimen is compensated beneath it. Finally, without energizing the deflection coils and without refocusing, the actual exposure on the photographic plate is made. The interesting part of the specimen that produces the image then lies on the optical axis of the microscope again.

Diffraction patterns can be obtained not only in the conventional way by illuminating the specimen with a parallel beam but also by using a convergent beam. According to W. Kossel and G. Möllenstedt this gives a CBED pattern (CBED is the abbreviation for convergent-beam electron diffraction) built up from 'patches' instead of points^[11]. The patches have a structure that carries additional crystallographic information. The use of a convergent beam in diffraction has the further advantage of providing information about a very small part of the specimen.

An important mode for users who wish to depart from the programmed settings for electron-optical experiments is the free-control mode. In this mode the user can vary the currents through the separate lenses as he wishes.

Modes with the upper deflection coils dynamically energized

The upper deflection coils are mainly used when the instrument is used as a scanning-transmission electron microscope (STEM). The spot then describes a rectangular pattern across the specimen, and the video signal supplied to one of the monitors on a control panel (fig. 4) is modulated by a detector signal. The detectors situated above the specimen can be seen in fig. 8*d*. Their output signals correlate with the number of secondary electrons, the number of back-scattered electrons or the X-ray intensity. Fig. 9 shows that a signal can also be used that correlates with the number

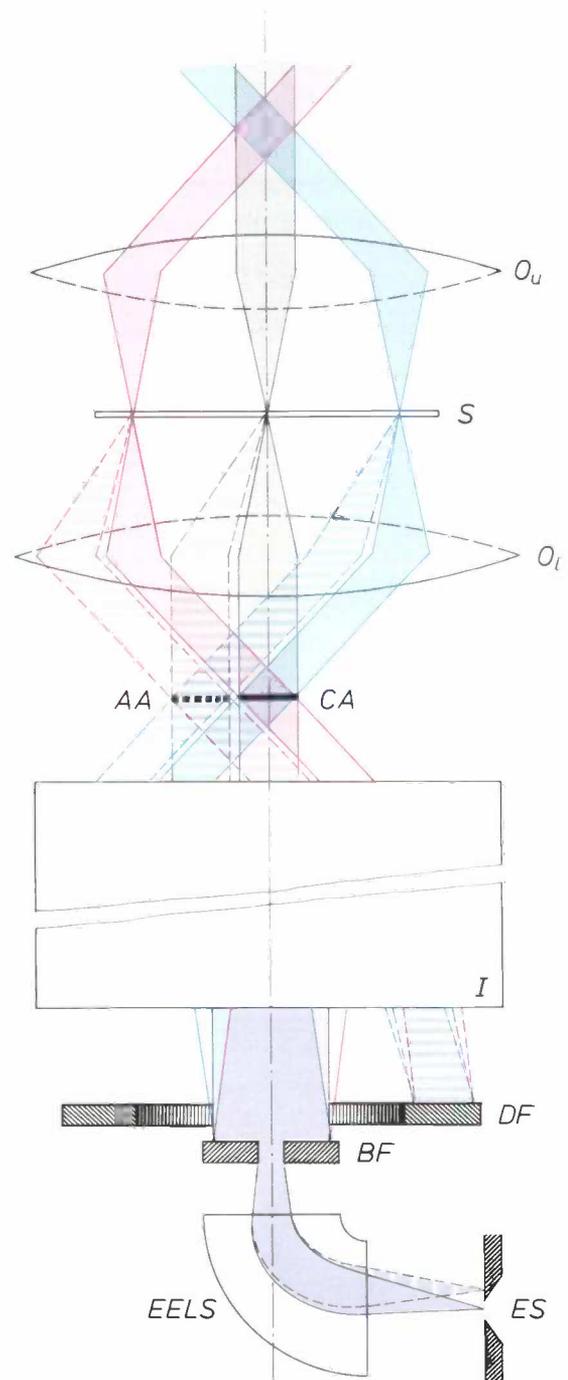


Fig. 9. Ray diagram for use of the microscope as a STEM, using the detectors in the projection chamber; see also figs 2 and 5. I imaging lenses below the objective. The grey, red and blue beams correspond to different positions of the spot, which describes a line across the specimen S . The beams defined by the continuous lines are not scattered in the specimen; they all pass through a circular area CA in the back focal plane of O_l . The beams defined by the dashed lines are scattered at a particular angle in the specimen; they all pass through an annular area AA , which is concentric with CA . A STEM bright-field image can thus be obtained with the signal from detector BF , and a STEM dark-field image with the signal from detector DF . In the $EELS$ detector electrons with a high kinetic energy are not deflected so much in a magnetic field perpendicular to the plane of the drawing as electrons of low energy. The electrons can thus be sorted by energy. ES exit slit of the detector.

of electrons transmitted through the specimen. Fig. 10a shows the corresponding ray diagram when the nanoprobe spot is used.

Fig. 9 demonstrates the advantages of a symmetrical immersion objective when the instrument operates as a STEM. When a parallel beam is incident on the condenser part (O_u) of the objective a parallel beam emerges from the objective part (O_l), since the specimen is located in the common focal plane. In addition to non-scattered beams there are also scattered beams emerging from the objective. All the non-scattered beams pass through the same circular area in the back focal plane of the objective part. All the beams scattered at a particular angle pass through an annular area concentric with the first one. The special feature is that these areas do not move when the spot travels across the specimen.

Since the lens system forms an image of both the circular area and the annular area, it is in general possible to measure the intensities of transmitted non-scattered and scattered beams. This is done with the detectors *BF* and *DF* in the projection chamber (fig. 9). The corresponding video signals produce a STEM bright-field image or a STEM dark-field image on a monitor. The transmitted electrons can also be sorted by energy with an EELS detector (EELS stands for electron energy-loss spectroscopy); we shall return to this presently.

In the SEM and STEM modes a choice can be made between high-magnification (HM) and low-magnification (LM) modes. In the S(T)EM-HM mode the nanoprobe spot (fig. 10a) is used, with the fineness of the raster pattern on the specimen adapted to the size of the spot. In the S(T)EM-LM mode a coarser pattern and a much larger spot are used. This is done by switching off the objective and focusing the illuminating beam with the condenser lens C_2 (fig. 5).

Another mode is the TEM dark-field mode with conical illumination. The incident beam is tilted over a preset angle (fig. 6d) and then rotated continuously about the optical axis. On a photographic plate a kind of superimposition of dark-field images is obtained, since all the azimuthal angles of the illuminating beam are described. The advantage of this is that more details can be seen in the image than in an ordinary dark-field image, since the details are illuminated from many more directions.

Modes with the lower deflection coils dynamically energized

In the SCIM mode (SCIM: scanning in imaging) the image of the specimen is scanned with the aid of the lower deflection coils, see fig. 10b, making use of detector *BF* in the projection chamber. In the SCID

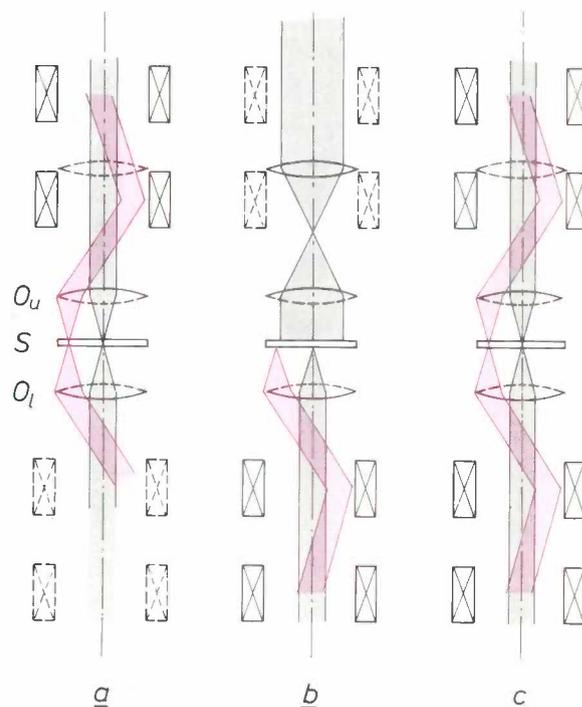


Fig. 10. Ray diagrams in different modes, when a) only the deflection coils above the specimen b) only the deflection coils below the specimen, and c) the deflection coils above and below the specimen are used. The ray diagram in (a) is the one used for the nanoprobe spot in STEM. The ray diagram in (b) is the one used for scanning the specimen from below; the signal from detector *BF* (fig. 9) can then be converted into a video signal for one of the monitors. The ray diagram in (c) is the one used in the low-dose mode, when the coils are statically energized. When the coils are dynamically energized, this is the ray diagrams for use as a STEM with the *EELS* detector; see also fig. 9. The deflection of the beam above the specimen is then compensated below it.

mode (SCID: scanning in diffraction) the image of the diffraction pattern is treated in the same way. In these forms of scanning, the image of the specimen or the diffraction pattern moves across the detector. The result presented on a monitor is comparable with an ordinary TEM micrograph. The advantage of this procedure is that the video signal can be processed electronically. If the analog video signal is converted into a digital signal with an A/D converter, the signal can also be processed numerically.

Modes with the upper and lower deflection coils dynamically energized

In the EELS detector (fig. 9) the electrons transmitted through the specimen are deflected through an

[11] G. Möllenstedt, Über die chromatischen Verluste von Elektronen beim Durchtritt durch Materie, *Optik* 9, 473-480, 1952; G. Thomas and M. J. Goringe, *Transmission electron microscopy of materials*, Wiley, New York 1979; M. N. Thompson, A review of TEM microdiffraction techniques, *Philips Electron Opt. Bull. EM* 110, 31-39, 1977; J. W. Steeds, Convergent beam electron diffraction, in: J. J. Hren, J. I. Goldstein and D. C. Joy (eds), *Introduction to analytical electron microscopy*, Plenum, New York 1979, pp. 387-422.

angle of about 90° in a magnetic field. The angle of deflection depends on the energy of the electrons: high-energy electrons are deflected less than those with low energy. At a particular current through the magnetic coil the electrons passing through the output slit have a particular energy, so that the electrons at the input of the detector can be sorted by energy by varying the current. These electrons always have to pass through the aperture in the partition diaphragm (*PD* in fig. 5) below the second projector lens, which separates the vacuum of the column from that of the projection chamber. At the input of the EELS detector, however, the electrons must all have the same direction. If the EELS detector is to be used for making a STEM micrograph, the deflection due to the upper deflection coils must therefore be compensated by an equal and opposite deflection due to the lower deflection coils; see fig. 10c. The ray diagram is comparable with that in the low-dose mode, but there the deflection coils are not used dynamically but statically.

The control system

The user communicates with the instrument by means of controls that adjust quantities such as currents in lenses, deflection coils and stigmator coils and d.c. voltages in the electron gun. As a result a picture appears on the screen in the projection chamber, photographic material is exposed or an image is formed on a monitor screen. In earlier microscopes there were fixed connections between almost every coil or electrode in the microscope column and one or more controls on the control panels. This meant that there were very many manual controls, and the user had to be thoroughly familiar with the consequences of varying any particular current or voltage.

The new microscope discussed here has some manual controls and a control screen for communication with the user; see fig. 11a. The screen shows the mode the instrument is operating in and the settings of various optical parameters. On each side of the screen there are eight controls, called 'soft keys', whose function is determined by the control software. The information on the screen, which is called a 'page', indicates the function of these soft keys; see fig. 11b. The other controls on the control panels do not in general have any fixed connection with the microscope column either. When a control is operated, it delivers one or more electrical pulses, which are converted into a command to the control software. The controls are therefore usually associated with a function. The principal functions are magnification, image brightness and focus. In addition to the soft keys there are two 'ordinary' controls whose function is not fixed



Fig. 11. The user communicates with the instrument by means of a) the controls on the panels and b) the soft keys on either side of the control screen. Information displayed on the screen includes the name of the mode in which the microscope is operating and the parameter settings.

but can be changed. These are the multifunction controls, which are mainly used for alignment of the microscope column; their function is indicated by the pages relating to the alignment procedure.

The control software guides the user to his goal — a micrograph or an analysis — via the pages on the screen and the controls on the panels. The user does not therefore need to know as much about the optical structure of the instrument as he did with previous instruments. Fig. 12 gives an idea of the hierarchy of the pages that the user can call to the screen. The figure gives four levels. (The particular pages shown refer to a simpler version of the microscope, with a limited range of mode options.)

The page from level 1 shows a menu with four options. If the user selects the option MODES, the screen then shows the page MODE SELECTION from level 2. The menu now gives options for the modes that are available with the particular version of the instrument, and the option CONFIGURATION, which shows a page that displays the features of the instrument. (This page shows here that the user has placed an electron gun with a tungsten cathode, not

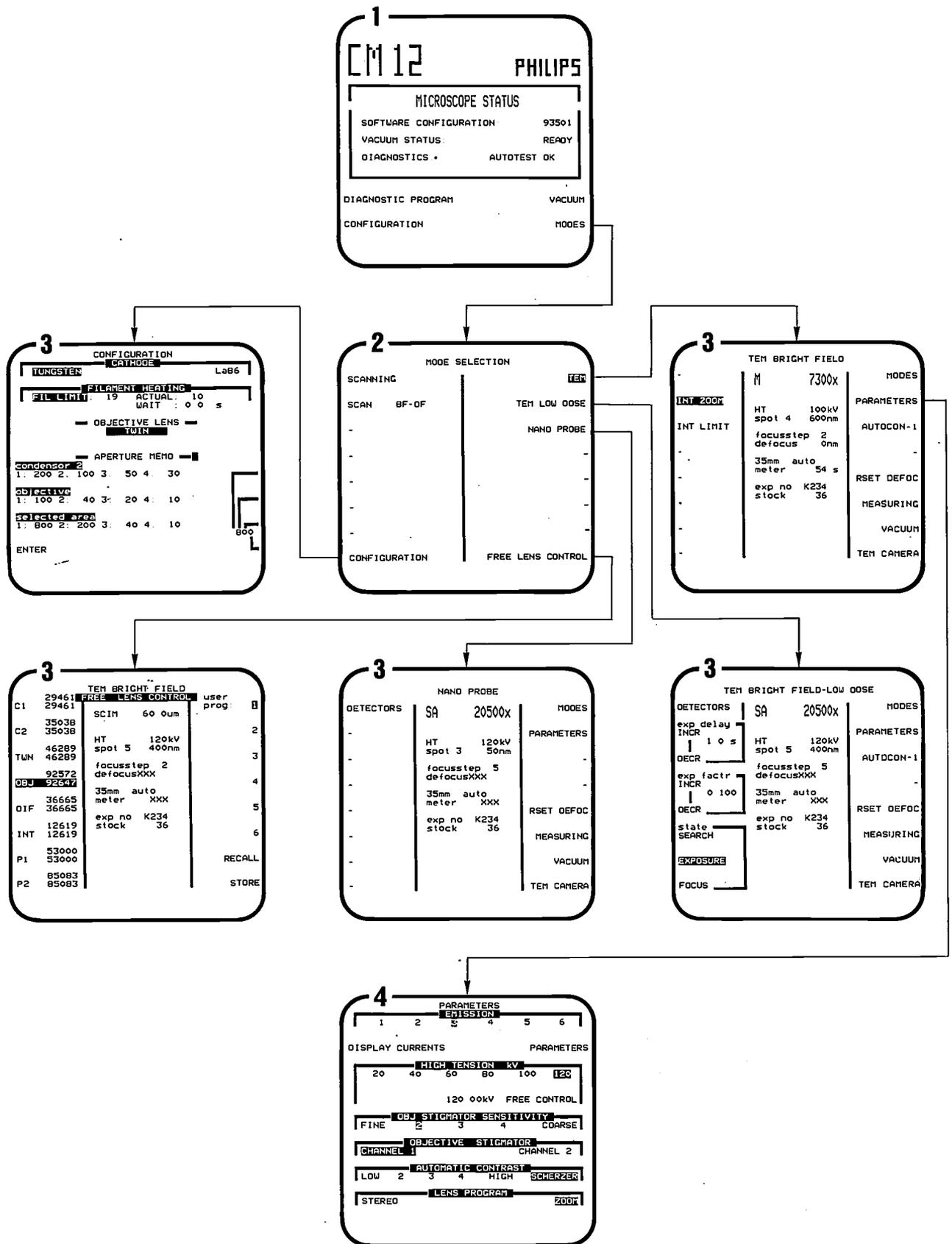


Fig. 12. The hierarchy of some of the many 'pages' that can be displayed on the control screen. The number at the upper left of each page indicates the level. The mode is selected from the options on a page at level 2. The mode to which the instrument is set can be seen from a page at level 3. The field in the centre then gives the settings of the parameters; the fields on either side indicate the functions of the soft keys. Parameters are changed via a page at level 4.

an LaB₆ cathode, in the microscope.) The pages at level 3 associated with the various modes show at the centre the settings of parameters such as the 'high tension' and spot size. The functions of the sixteen soft keys are indicated on the left and right. If the user presses the soft key PARAMETERS, the screen displays the relevant page of level 4. This indicates the options selected for the parameter values (dark against a light background) and the options still open (light against a dark background). Options previously selected can be changed by pressing one of the two soft keys on the appropriate line. The user can return to a page at the preceding level by pressing the key READY at the bottom of the screen; see fig. 11*b*. Light-emitting diodes (LEDs) indicate whether this key or other keys are 'active'.

The heart of the control system is a 16-bit microprocessor, the INTEL 8086. To gain computing speed and accuracy this microprocessor works in conjunction with a numeric data processor, the INTEL 8087. This can process floating-point numbers between 8.43×10^{-37} and 3.37×10^{38} to an accuracy of seven decimal places. The 8086 microprocessor processes the commands given by the user or by sensors with the RMX 88 operating system. The operating system ensures that the commands are dealt with virtually in real time. All the computed results are transmitted every 15 ms in order of importance, under the control of an external clock signal. The two processors and the input and output units interact via an IEEE 796 bus system. The 8086 processor communicates with the various electronic circuits through a bus system that we have designed. This bus system is insensitive to interference (from high-voltage flashover for example). This insensitivity is largely due to the use of optical couplings.

As an example we shall show, with reference to fig. 13, how the current in the deflection coils above and below the objective is controlled. In general this current is equal to the sum of a direct current used for correcting the beam direction and a sawtooth alternating current used for scanning. The current in the deflection coils must be very accurately controlled to avoid unacceptable movement of the image.

The circuit for generating the current through the deflection coils is designed in such a way that a digital signal from the control system is multiplied by an analog signal in two digital/analog converters. In one digital/analog converter the digital signal D_{DC} is multiplied by an accurate direct voltage V_{ref} . In the other digital/analog converter the digital signal D_{AC} is multiplied by an accurate sawtooth voltage V_{ST} . Multiplication can be performed in four quadrants, so that for example two negative signals multiplied together

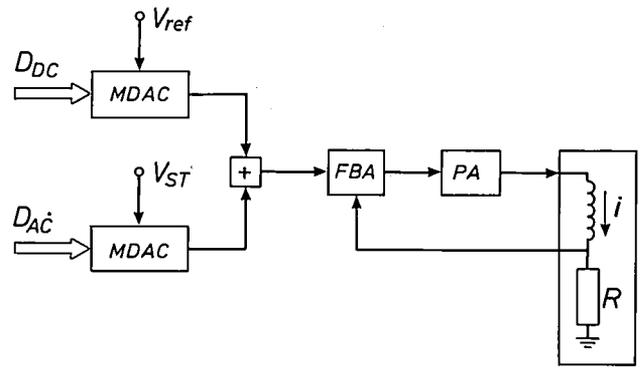


Fig. 13. Diagram of the current control in a set of deflection coils above or below the objective. MDAC multiplying digital/analog converters. D_{DC} and D_{AC} digital signals from the control system. V_{ref} accurate direct voltage. V_{ST} accurate sawtooth voltage. FBA feedback amplifier. PA power amplifier. Both amplifiers are part of a control network that includes the deflection coils. The value of the current i in these coils is fed back as a voltage across a resistor R .

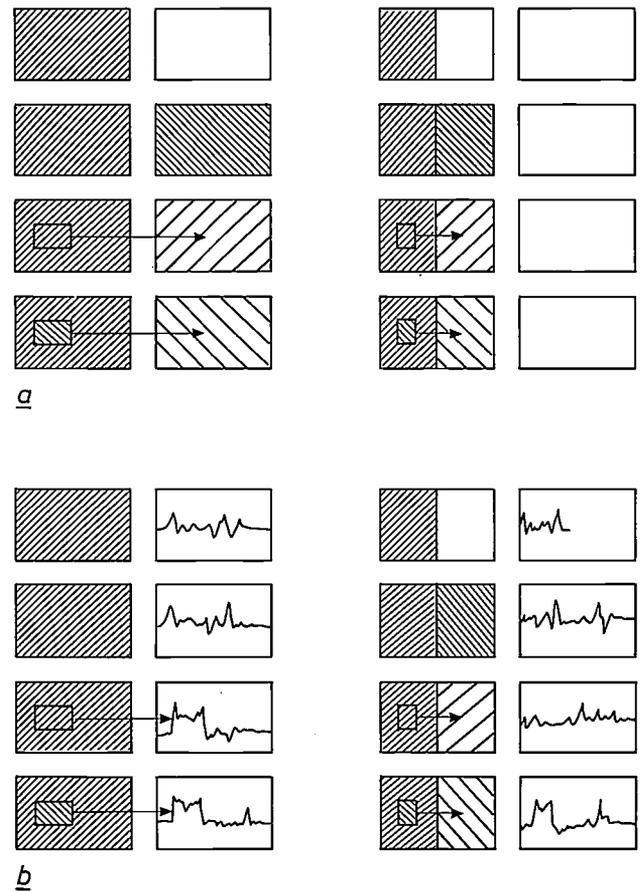


Fig. 14. STEM results displayed on the two monitor screens. Images are indicated by hatching. The same images have the same direction of hatching; coarser hatching indicates higher magnification. *a*) One image or two images on each screen or on each half of the left-hand screen. An image can be projected into the other image at a different magnification. Fig. 19 shows an example for the choice given schematically at the lower right. *b*) One image or two images on the left-hand screen only. The right-hand screen displays the brightness variation of one line as a curve. The curve may also be displayed on the 'videoscope', fitted to the microscope as an accessory; see fig. 4*a*.

give a positive signal. The product signals are then added. The result acts as a reference signal in a control network, which includes a feedback amplifier and a power amplifier. The power amplifier supplies a current i to the deflection coils. The value of the current is fed back into the control network as a voltage across a resistor.

The period of the sawtooth voltage V_{ST} in fig. 13 is between 1 and 1000 ms, or is equal to the line period in one of the television standards. The amplitude of the analog sawtooth signal is digitally calibrated by the control system every time the microscope is switched on. Because of this automatic calibration the advantage of digital electronics — accuracy — is combined with that of analog electronics — speed.

In the STEM modes the video signals for the two monitors (see fig. 4a) are derived from one or more detector signals. The video signals are the result of such operations as amplification, switching, mixing and filtering. These operations have to take place in a large bandwidth (up to 10 MHz) without the occurrence of noise or distortion due to overloading. This difficult problem has been solved by means of constant-current-source technology, also used in high-frequency oscilloscopes. The signal from the detectors is converted into a current that appears to come from

a source with an infinitely large internal impedance. There are two channels for video signals, so that the two monitors can be operated separately. The channels can also be combined in a variety of ways, and different magnifications can be selected for the monitors. One screen may be divided into two halves, so that each half shows the signal from one channel; see fig. 14a. If required, a monitor can display the video signal of one line as a curve; see fig. 14b.

Some results of investigations made with the microscope

To conclude, we shall illustrate the versatility of the CM12/STEM microscope with a number of micrographs of different specimens, made in different modes. The corresponding page on the screen is shown with each picture. Fig. 15 is a TEM micrograph of a gold particle vacuum-evaporated on a film of amorphous carbon. Twinning structures are visible in the gold particle. The micrograph was made with the Super-Twin objective at a magnification of about 8×10^6 . More important than the magnification is the resolution achieved. The spacings of the (111) planes and the spacings of the (11 $\bar{1}$) planes in the crystal lattice (the value is 0.235 nm) appear to be visible. Since

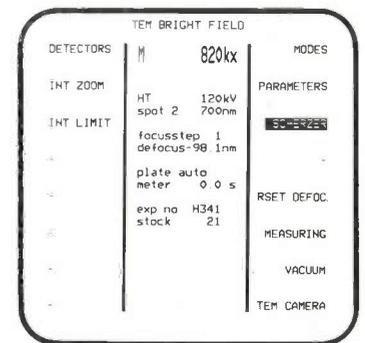
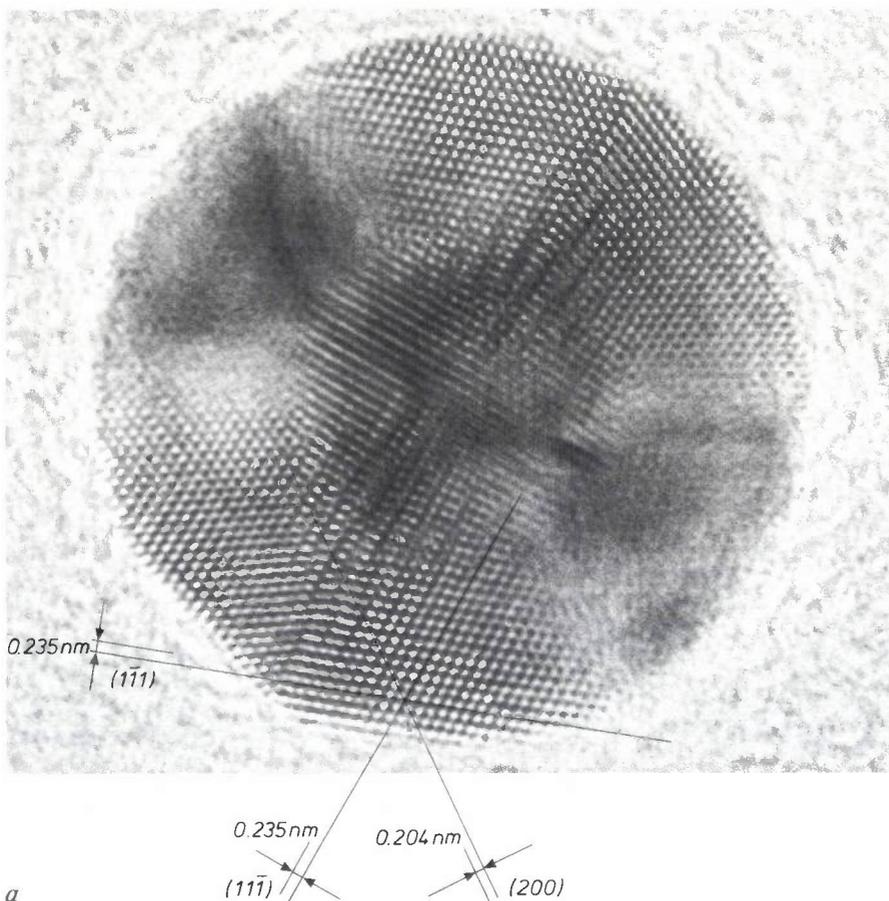
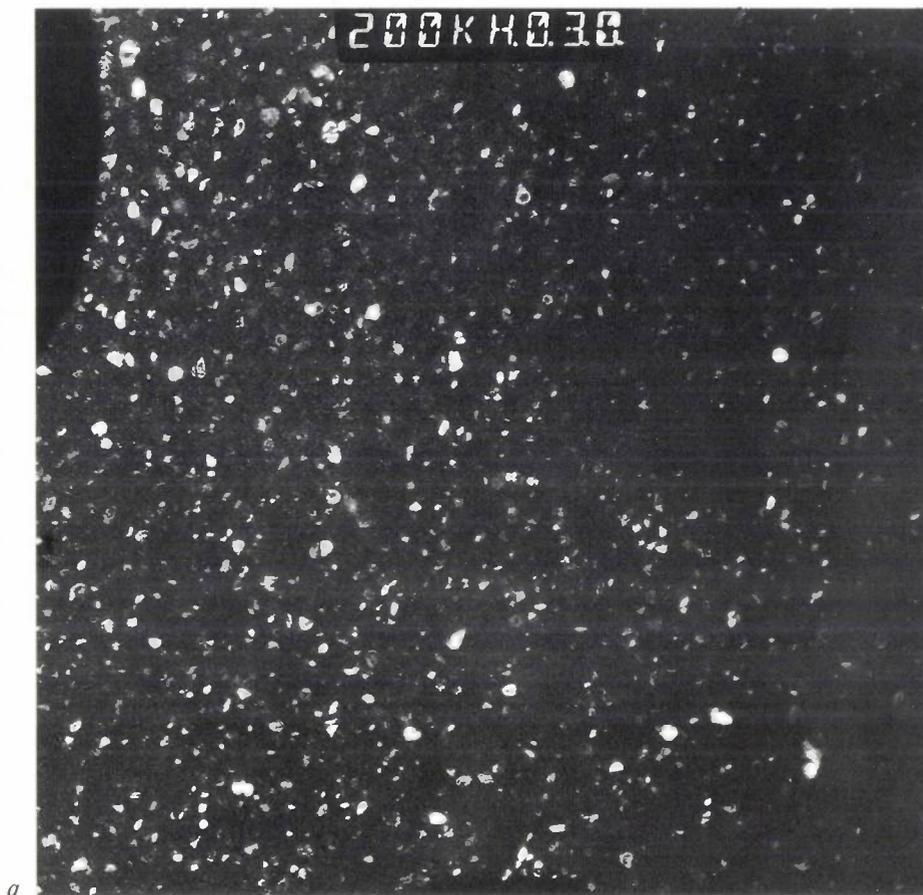
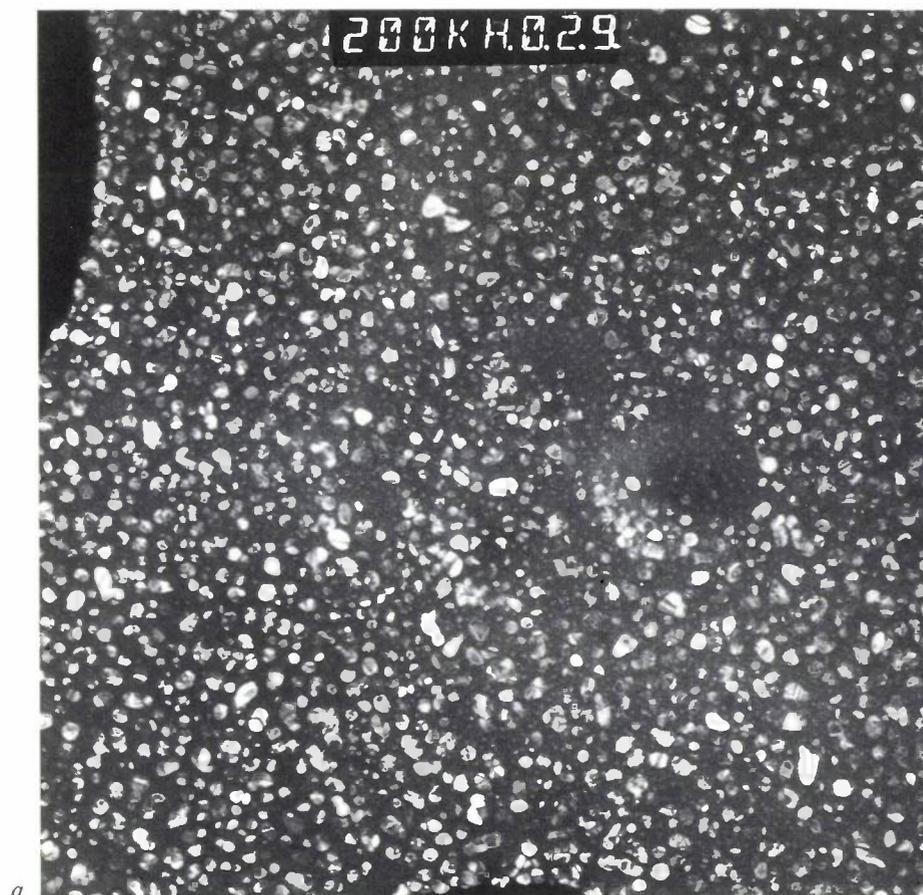


Fig. 15. a) TEM bright-field micrograph of a gold particle evaporated on to a film of amorphous carbon. Images of this kind are generally used for demonstrating the resolution of a transmission electron microscope. The spacings of some of the lattice planes are indicated in the picture; the series of points are interference patterns, so that it cannot be said that atoms are visible. The total magnification factor of the micrograph here is about 8×10^6 . b) The corresponding page. The magnification shown corresponds to the magnification on the screen in the projection chamber. A 'Scherzer contrast' has been used here, i.e. the image is defocused until the phase contrast reaches a maximum as judged by a particular criterion written into the control program.



TEM CONICAL DARK FIELD			
DETECTORS	SA	200kx	MODES
INT ZOOM	HT	120kV	PARAMETERS
	spot	3	AUTOCOR-1
INT LIMIT	focusstep	2	
	defocus	0pm	
	plate auto		
df-mode	meter	0.0 s	RSET DEFOC
<input checked="" type="checkbox"/> CONE	exp no	A001	MEASURING
	stock	35	
circle			VACUUM
<input checked="" type="checkbox"/> DYN			TEM CAMERA
INCR			
↓	DF1tx	-0.24 d	
DECR	DF1ty	0.44 d	

Fig. 16. *a*) TEM dark-field micrograph of the same type of gold particles as in fig. 15, but now at reduced magnification. Not all of the particles are visible (see fig. 17) because the illumination was from one side. *b*) The corresponding page. The tilt angle of the illumination can be calculated from the tilt angles in the *x*- and *y*-directions, stated at the bottom of the page. The user can select a different azimuthal angle if desired.



TEM CONICAL DARK FIELD			
DETECTORS	SA	200kx	MODES
INT ZOOM	HT	120kV	PARAMETERS
	spot	3	AUTOCOR-1
INT LIMIT	focusstep	2	
	defocus	0pm	
	plate auto		
df-mode	meter	0.0 s	RSET DEFOC
X-Y <input checked="" type="checkbox"/> CONE	exp no	A000	MEASURING
	stock	36	
circle			VACUUM
<input checked="" type="checkbox"/> MAX			TEM CAMERA
INCR			
↓	DF1tx	0.50 d	
DECR	DFrotXXX		

Fig. 17. *a*) Micrograph comparable with the one in fig. 16 but now with conical illumination. The tilted illuminating beam rotated continuously about the optical axis. Since the illuminating beam passed through all azimuthal angles, nearly all the gold particles are visible. *b*) The corresponding page.

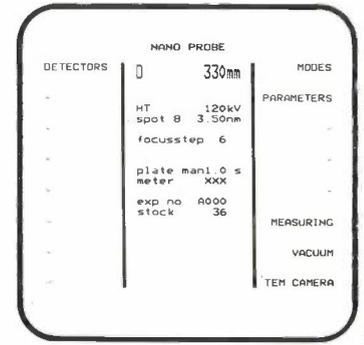
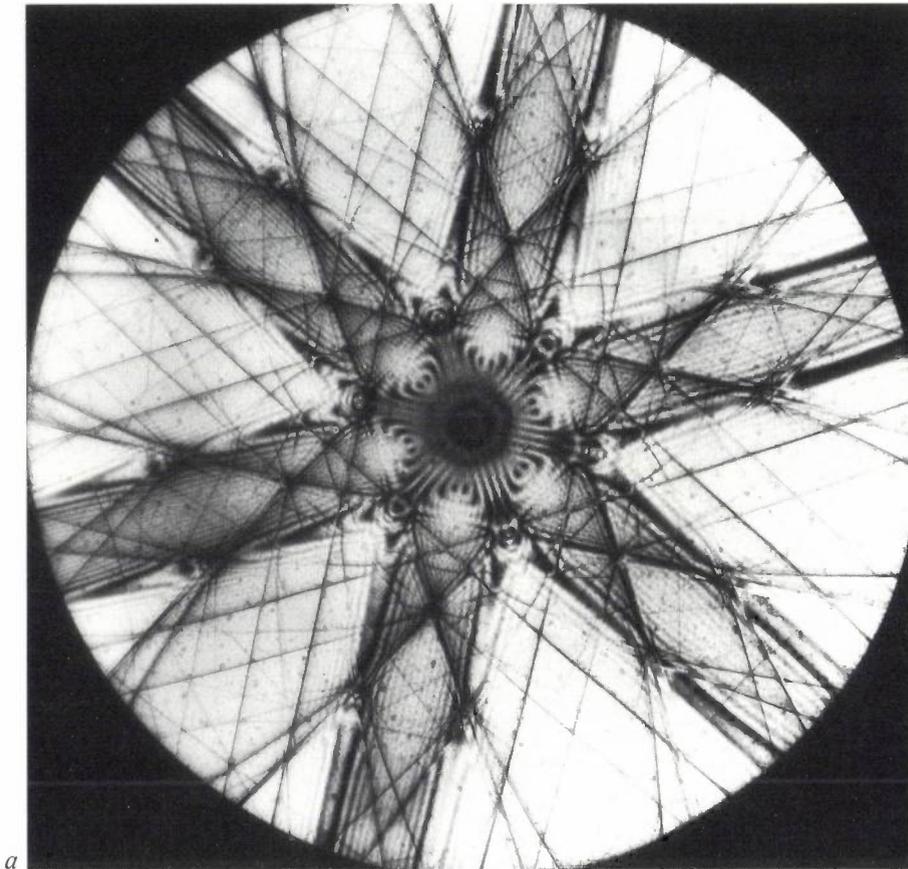


Fig. 18. a) Micrograph of a CBED pattern (CBED: convergent-beam diffraction pattern^[11]) of single-crystal silicon. Because of the [111] orientation of the specimen surface a 120° symmetry can be observed. b) The corresponding page.

the intersection of these planes with the plane of the micrograph is 'visible' — in fact, interference patterns are observed — there is also an indication of the intersection with the (200) planes, whose spacing is 0.204 nm.

Fig. 16 shows a TEM dark-field micrograph of the same amorphous carbon film with gold particles, but now at a lower magnification. For comparison fig. 17 shows a TEM dark-field micrograph with conical illumination of the same specimen, made with the tilted beam rotating continuously about the optical axis. More particles are visible in fig. 17 than in fig. 16, because for fig. 16 the specimen was illuminated from one direction only. In the case of fig. 17 the illuminating beam passed through all possible azimuthal angles.

Fig. 18 shows a CBED pattern, where the illuminating beam converges so as to make the spots of the diffraction pattern overlap completely. The resulting complicated pattern^[11] gives information about the symmetries in the crystal structure, here of single-crystal silicon with [111] orientation. Observation of the cubic structure of silicon in the [111] direction reveals a 120° symmetry. The specimen is locally 100 to 200 nm thick. This was achieved by ion-etching a silicon wafer a few microns thick until it was so thin that a small hole appeared in it. The diffraction micrograph corresponds to a 'wedge' of material next to the

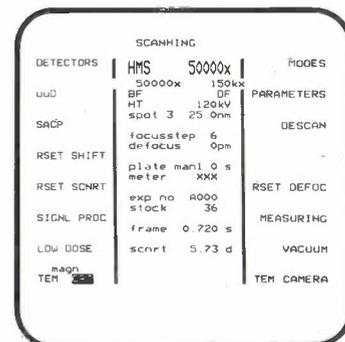
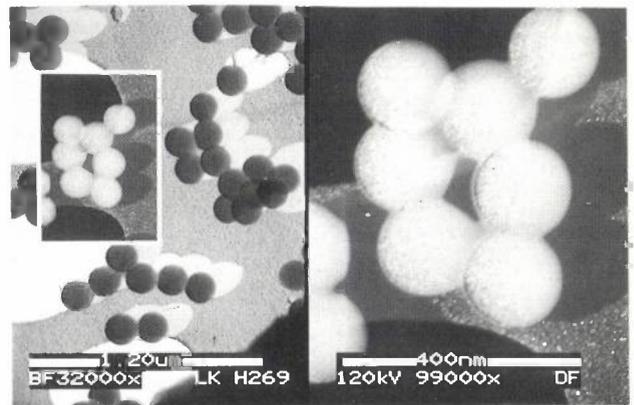
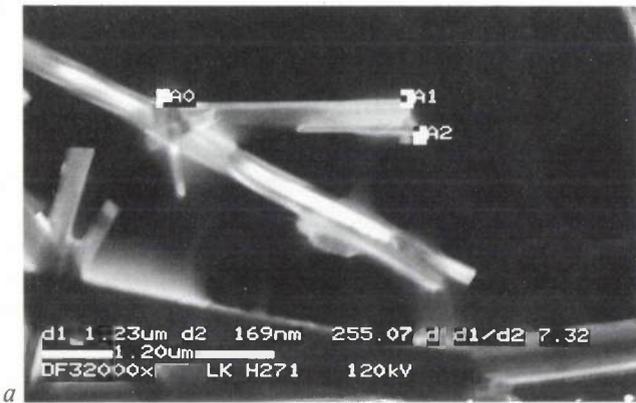


Fig. 19. a) STEM-HM micrograph of latex spheres. Gold has been evaporated on to the specimen from one direction. On the left-hand half of the screen there is a bright-field image, on the right-hand half a dark-field image. The image on the right is projected into the left-hand image at a different magnification. b) The corresponding page.

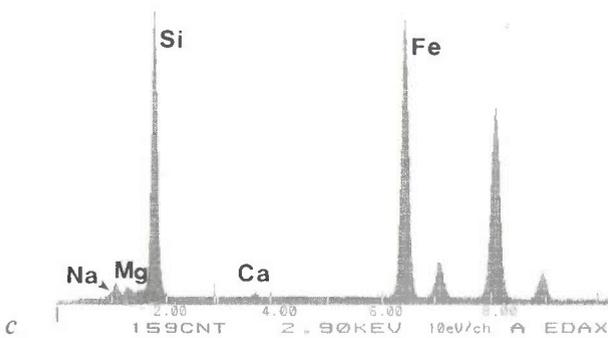


a

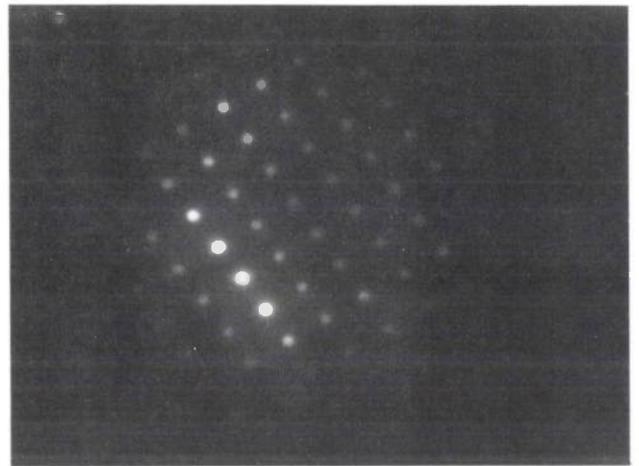
SCANNING			
DETECTORS	HMS	50000x	MODES
UUD	BF	DF	PARAMETERS
	MT	120kV	
SACP	spot 3	25.0nm	DESCAN
RSET SHIFT	focusstep	6	
	defocus	0pm	
RSET SCHRT	plate man1	0 s	
	meter	xxx	measuring
SIGNL PROC	exp no	A000	STEREO
	stock	36	cal
LOW DOSE	d1	1.23um	
	d2	169nm	
magn	angle	255.07 d	
TEM	d1/d2=	7.32	ENTER

b

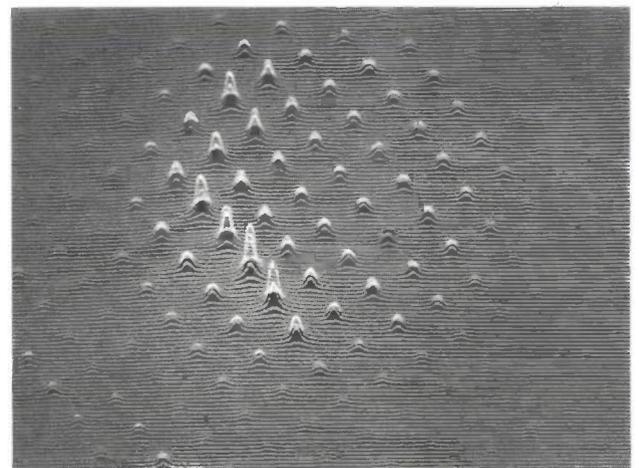
23-JAN-87 08:43:45 EDAX READY
 RATE= 610CPS TIME= 400LSEC
 FS= 1049ICNT PRST= 100LSEC
 A=CROCIDOLITE ASBESTOS



c



a



b

POST SPECIMEN SCANNING			
DETECTORS	SCID	70.0mr	MODES
INT ZOOM	BF	DF	PARAMETERS
	HT	120kV	
INT LIMIT	spot 3	500nm	
	focusstep	1	
RSET SCHRT	plate man1	0 s	
	meter	0 0 s	
SIGNL PROC	exp no	A000	MEASURING
	stock	36	VACUUM
	frame	0 720 s	
	scart	28 19 d	
magn			
TEM			TEM CAMERA

c

Fig. 20. a) STEM-HM micrograph of crocidolite fibres. To confirm that we are dealing with asbestos, the length/breadth ratio d_1/d_2 was determined. This ratio was calculated by the software and shown in the micrograph. A0, A1 and A2 are reference points for the measurement. The micrograph also gives the measured lengths of d_1 and d_2 , corrected for the magnification, and the angle between d_1 and d_2 . b) The corresponding page, also with the measured results. c) Spectrum taken with the energy-dispersive X-ray spectrometer, showing the peaks for the elements that are characteristic of crocidolite. The three peaks on the right are a second Fe peak and Cu peaks, originating from copper in the specimen gauze.

Fig. 21. SCID micrograph (SCID: scanning in diffraction) of a MoO_3 crystal. a) Image on the monitor screen, which is comparable with an electron-diffraction micrograph taken with one of the two cameras in the projection chamber. b) Almost the same image, but now with an extra vertical deflection for each line, which is proportional to the brightness in (a). The relative brightness, and hence the relative value of the detector signal, can thus be measured in the image. c) The corresponding page.

hole. Micrographs of this kind are very useful in the study of silicon since dopants show up as changes in the pattern.

Fig. 19 shows a STEM-HM micrograph of tiny spheres of latex. To obtain a shadow effect, gold was evaporated on to the specimen from one direction. The gold particles on the specimen are comparable with those in figs 15, 16 and 17. Specimens of this type

are widely used in scanning electron microscopy for checking the resolution. The left-hand half of the monitor screen shows a bright-field image, the right-hand half a dark-field image. The image on the right was projected at a different magnification into the image on the left. The method of splitting the screen corresponds with that shown on the lower right in fig. 14a.

Fig. 20 shows how the presence of asbestos, whose harmful effects are well known, can be demonstrated. *Fig. 20a* is a STEM-HM dark-field micrograph of fibres of blue asbestos (crocidolite), one of the most harmful forms of asbestos. Crocidolite has the formula $\text{Na}_{2-x}\text{Ca}_{0.5x}(\text{Mg},\text{Fe}^{2+})_{3+1.5x}\text{Fe}_{2-x}^{3+}\text{Si}_8\text{O}_{22}(\text{OH})_2$. The micrograph also presents the result of a method of measuring the dimensions of one of the fibres. A distance in the image can be defined by positioning a cross-wire, pressing a button, then repeating this at another setting. In this micrograph, and also on the page in *fig. 20b*, the length and breadth of the fibre are displayed, with their ratio and the angle in degrees between the directions of measurement. The length/breadth ratio, which should be at least 3, is extremely important in the characterization of asbestos. The measurement of this ratio and *fig. 20c*, which shows the spectrum measured with the energy-dispersive X-ray detector, confirm that this is a micrograph of crocidolite. (The three unnamed peaks on the right of *fig. 20c* are a second Fe peak and two peaks due to Cu in the gauze supporting the specimen.)

Fig. 21 shows SCID micrographs of single-crystal molybdenum oxide (MoO_3). The monitor image in *fig. 21a* is comparable with a normal diffraction pattern. One of the bright spots is the zero-order image

of the electron source; the other spots are images of higher order. The distance between the spots is a measure of the lattice constant of the crystal. *Fig. 21b* shows the monitor image obtained when the detector signal is not used for modulating the intensity of the electron beam in the monitor tube but for producing extra vertical deflection of the beam. In this way the relative value of the detector signal can be measured in the monitor image. The changeover between this picture and the previous one is not made via the page on the control screen, as in *fig. 21c*, but by means of a key.

Summary. The CM12/STEM electron microscope has a Twin objective and microprocessor control. With this special objective it is possible to switch over quickly from a conventional TEM mode, with a direct image of the specimen or diffraction pattern on a fluorescent screen in the projection chamber, to a STEM mode, where the specimen is scanned by the illuminating beam and a detector signal is converted into a video signal for one of the monitors. Because of the special magnetic configuration of the objective, which includes a mini-condenser lens, this switch-over produces hardly any change in the heat-flow conditions in and around the polepieces. The CM12/STEM is very much easier to operate than other microscopes because the microprocessor guides the user to the various operating modes and features by providing information on a control screen. The use of 'soft keys' and multi-functional controls on the panels keeps the number of manual controls relatively small. Micrographs demonstrate the resolution and some of the many features of the instrument.

Noise due to optical feedback in semiconductor lasers

B. H. Verbeek, D. Lenstra and A. J. den Boef

Wave reflection of signals into a device in which they are amplified can have considerable effects on its operation. A typical example is acoustic feedback with a microphone. An effect not so well known, but just as undesirable, is the fluctuation in the optical power of a laser because of reflected radiation. Since reflected radiation is an essential feature in many applications of semiconductor lasers, that radiation may lead to problems in the laser cavity.

Introduction

After the construction of the first working lasers in the early sixties, there was much speculation about the potential uses of the new radiation source. They ranged from 'science-fiction' applications to machining aids and distance measurements. While many of the uses envisaged at that time have since materialized, no-one could then have foreseen the scale on which lasers would be used in a wide variety of consumer products as well as in scientific equipment. Semiconductor lasers are used today as light sources in optical communications and in Compact Disc and LaserVision players. A problem that can arise in such applications is the fluctuation in intensity that can occur if radiation is reflected back into the laser cavity. Some aspects of this optical feedback will be dealt with in this article.

Laser action is based on stimulated emission, a process in which a photon of energy $h\nu$ stimulates electrons that are in an energy state E_1 to return to a state of lower energy E_2 , with the emission of a photon such that $E_1 - E_2 = h\nu$; see fig. 1. The emitted photon has the same energy and phase and moves in the same direction as the original photon, and in its turn can cause further stimulated emission. To obtain the

highest possible stimulated emission the photons are reflected back and forth through the material a number of times from two reflecting surfaces. These reflecting surfaces define the boundaries of the laser cavity. (Since we shall refer to an external feedback reflector later, we shall call this cavity the internal cavity.) The two reflecting surfaces allow some of the radiation to pass through, giving a beam of coherent radiation that only contains energy at the frequencies that lead to constructive interference on repeated reflection (the longitudinal modes). A necessary condition is that there should be more electrons available in the high energy level than in the low one (population inversion). This condition can be met in a number of ways. In a solid the population inversion can be brought

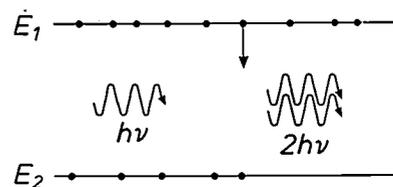


Fig. 1. Laser action can occur when the energy level E_1 of the atoms or molecules of a material contains more electrons than the lower energy level E_2 . An incident photon of energy $h\nu = E_1 - E_2$ causes the emission of a second photon with the same energy and phase as the original photon.

about by intense light flashes. In a gas it results from collisions with accelerated electrons or with fast molecules or atoms. The inversion in semiconductor lasers is achieved by injecting charge carriers into a pn junction.

Briefly, the operation of a semiconductor laser is as follows; see fig. 2. A current in a pn junction in a direct [1] semiconductor causes electrons to flow from the conduction band in the n-type material towards the conduction band in the p-type material. These electrons can recombine with holes in the p-type material, with the emission of photons. Above a certain current stimulated emission occurs. The optical power of the laser can be directly controlled by varying the current.

Semiconductor lasers are now more and more widely used in various kinds of business and professional applications. This is because they are small and robust, are not so expensive as other types of lasers, and are easy to use.

An important application of semiconductor lasers is in information read-out in Compact Disc and Laser-Vision players. In these players the laser light is focused on the reflecting layer of the disc. The information is contained in the disc in the form of pits with a depth of about $\frac{1}{4}\lambda$ (λ is the wavelength of the laser light). The intensity of the reflected light is detected by one or more photodiodes. The pattern of the pits on the disc is represented in the detector signal by the difference in intensity between the light reflected from the bottom of the pits and the light reflected from the

surface of the disc. Because of interference, the light reflected by the bottoms of the pits is less intense than the light reflected by the surface of the disc. The detector signal therefore approximates to a rectangular waveform with two levels. In spite of precautions, some of the reflected light arrives at the output reflector of the laser (optical feedback) rather than the detector. This radiation can affect the laser action in the internal cavity, producing fluctuations in the intensity of the laser light. These fluctuations can be large enough to distort the intensity pattern of the light reflected from the disc, introducing read-out errors. To avoid these problems it is necessary to understand the mechanisms involved in this unwanted feedback.

This article describes our experimental and theoretical work on these problems. First we shall discuss the structure and characteristics of the laser without feedback. We shall then deal with the subject of optical feedback, making a distinction between coherent feedback and incoherent feedback. We make this distinction because the feedback radiation may not necessarily be coherent with the radiation in the laser cavity. If the feedback is coherent (at low feedback levels) a good theoretical description of the behaviour of the laser can be obtained. In this article we shall mainly confine ourselves to coherent feedback, with a brief mention of some of the effects encountered in incoherent feedback.

Laser structure

In the introduction we touched briefly on the principle of the semiconductor laser with a pn junction in a semiconducting material ('homojunction laser'). Structures of this type give laser action above a relatively high threshold current, which gives problems with overheating. Laser action can be obtained at a much lower current when the pn junction consists of two different semiconducting materials. For this reason lasers have been developed that consist of two layers of different composition [2], known as 'double-heterojunction lasers'. There are many structures and materials in use, depending on the characteristics required for a particular laser application. The structure of the laser we have used in our investigations is the one used in Compact Disc and LaserVision players. In

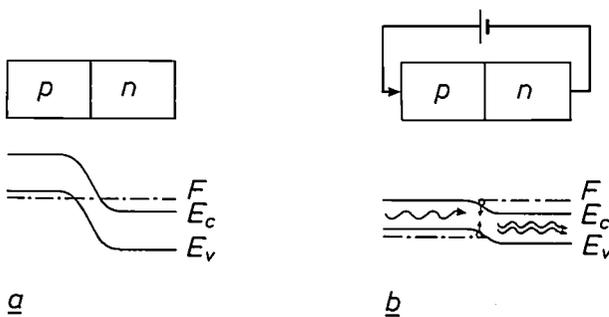


Fig. 2. The position of the energy level with the lowest energy in the conduction band E_c and the energy level with the highest energy in the valence band E_v , showing the population of these bands near a pn junction (schematic). *a*) When there is thermal equilibrium between the conduction and valence bands the population of the energy states by electrons and holes can be indicated by a single Fermi level F . *b*) If a voltage is applied across the junction, the corresponding energy levels in the p-type and n-type material come closer together. The energy barrier is now so small that the charge carriers can cross the junction. The thermal equilibrium between the bands is now disturbed and the population of the bands can no longer be indicated by a single Fermi level. If the current that starts to flow through the junction is not too high, an equilibrium still exists in each separate band and the population of the levels can be represented as shown. This causes population inversion in a region around the pn junction and stimulated emission can occur.

[1] Semiconductors can be either direct or indirect. The direct type is appropriate for laser action. This terminology indicates that the optical transitions between energy bands can take place directly, i.e. without the need for an extra impulse quantum of momentum, as required in indirect semiconductors.

[2] See J. C. J. Finck, H. J. M. van der Laak and J. T. Schrama, A semiconductor laser for information read-out, Philips Tech. Rev. 39, 37-47, 1980, and G. A. Acket, J. J. Daniele, W. Nijman, R. P. Tijburg and P. J. de Waard, Semiconductor lasers for optical communication, Philips Tech. Rev. 36, 190-200, 1976.

this section we shall take a brief look at the operation and characteristics of a laser with this structure. The problems that arise due to optical feedback are also to be found in many other laser structures, however. The results and the descriptions of these problems have a much wider relevance than for this special laser structure alone.

A cross-section through this structure is shown in *fig. 3*. The laser action is excited in layer *A*, the active layer, which consists of $p\text{-Al}_{0.14}\text{Ga}_{0.86}\text{As}$. This layer is sandwiched between two $\text{Al}_{0.45}\text{Ga}_{0.55}\text{As}$ layers, one p-type and the other n-type. (Since the Ga atoms in a GaAs lattice can be replaced by Al atoms without causing distortions or stresses in the lattice, it is relatively easy to make perfectly matching layers of different material.) Because of the difference in the aluminium content of the active layer and the boundary layers, the refractive index in the active layer is larger and the band gap is smaller than in the two adjacent layers. The band gap is the difference in energy between the lowest level of the conduction band and the highest level of the valence band. The energy-band diagram for this layer structure can be seen in *fig. 4*.

This structure of layers of semiconducting material is essentially that of an ordinary diode. When a voltage is applied across the diode in the forward direction the electrons of the n-type material and the holes of the p-type material move towards the active layer. Here both types of charge carriers are confined by the energy barriers ΔE_v and ΔE_c . Holes and electrons can only disappear from this active layer by radiative recombination (ΔE_v and ΔE_c are larger than the mean thermal energy of holes and electrons). Laser action in the diode starts when the current through the active layer is high enough to cause population inversion and the photons resulting from the radiative recombination can cause sufficient stimulated emission to compensate for the optical losses in the laser (the 'lasing' condition). To meet this condition the radiation has to be enclosed in a resonant cavity. This is effected in the *y*-direction by the difference in refractive index mentioned earlier. In the *z*-direction there is a step in the refractive index at the cleavage planes of the crystal, where about 30% of the radiation is reflected back into the laser cavity. The cavity boundary in the *x*-direction is formed by inserting a current-isolating layer of n-GaAs. The pnp structures on either side of the groove in this layer ensure that the current is limited to the groove, so that it is the only place where population inversion occurs. Since the laser light is present not only in the active layer but also in both boundary layers and even in a small zone of the current-isolating layer, there is a small step in the refractive index (5×10^{-3}) in the *x*-direction at the position

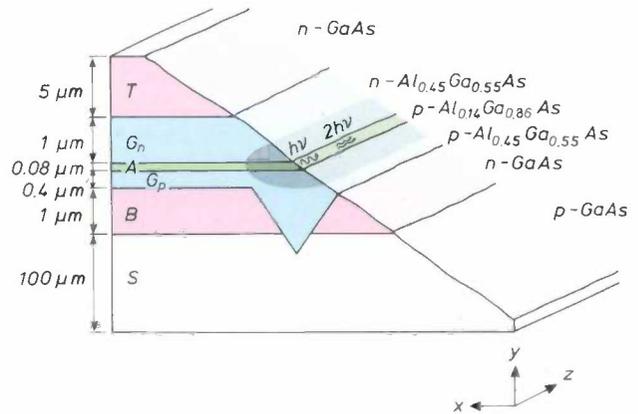


Fig. 3. Cross-section through the layer structure in a laser with a double heterojunction. A substrate *S* of p-GaAs contains, one above the other, a current-isolating layer *B*, a cladding layer G_p , an active layer *A* in which the laser action is produced, a cladding layer G_n and a top layer *T*, to which contacts can be applied. The V groove in the current-isolating layer ensures that the current injected into the structure in the *y*-direction is confined in the *x*-direction. Dimensions and compositions of the different layers are indicated in the figure. The laser light is generated in the grey area and propagates in the *z*-direction. The length of the layer structure in the *z*-direction is much greater than the transverse dimensions of the structure.

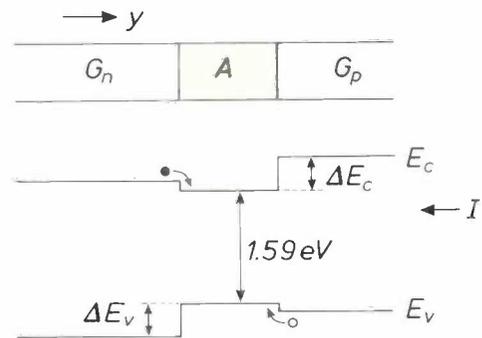


Fig. 4. Schematic representation of the band structure in the layers G_p , *A* and G_n of the laser. A current in the *y*-direction causes electrons in the conduction band of the n-type material to be pumped towards the active layer. Because of the energy difference ΔE_c between the conduction bands of the active layer and the boundary layer G_p , the electrons are confined. In an analogous way, holes collect in the active layer. Holes and electrons can only disappear from the active layer by recombination. Under the action of incident photons a coherent beam of radiation is thus produced.

of the V groove. As a consequence, light is optically confined in the *x*-direction as well.

The four layers are grown on a substrate of p-GaAs and this structure is then covered with a top layer of n-GaAs. These layers facilitate the application of contacts. As aluminium oxidizes easily, it is not easy to apply good metal contacts to AlGaAs.

A current in this structure in the *y*-direction produces radiation. *Fig. 5* shows the light output in the active layer as a function of the injection current. Up to a certain current value there is only spontaneous emission. The radiated power, the number of photons formed by recombination of electrons and holes, only

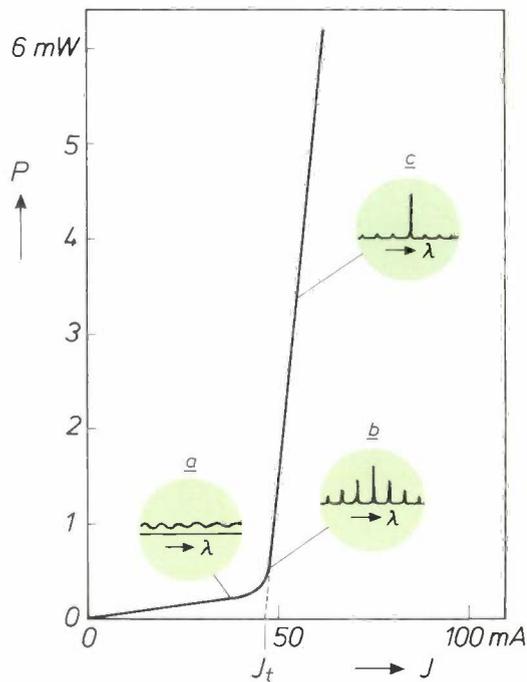


Fig. 5. Optical power P of the laser as a function of the current J . The spectrum of the laser is shown for a number of current values. *a*) Spontaneous-emission spectrum, *b*) multimode spectrum, *c*) single-mode spectrum.

increases slowly with the number of electrons injected into the active layer. When the injection current reaches a threshold value J_t , there is population inversion and stimulated emission can occur. Above this current the power increases strongly with current. At low currents the power is emitted in a fairly broad wavelength range (about 1 nm) close to 780 nm. When the current reaches a value well above the threshold current, the laser emits at a single sharply defined wavelength (single-mode behaviour with coherent radiation).

Optical feedback, experiments

In the applications of semiconductor lasers we have mentioned, the intensity or the laser power plays an important part, because it is the variation in the laser power reflected by the disc that contains the information to be read out or transmitted. This is why it is so important that the intensity of the laser is sufficiently constant for the intentional variations to be observed.

In a laser with no feedback there are various processes that cause fluctuations in the power. In the first place, there are small fluctuations in the laser power that are inherent in the process of radiative recombination (spontaneous emission). This noise contribution is fairly small and causes few problems, because it is constant above the threshold current. In addition, small changes in the input signal, the injection cur-

rent, affect the output signal or the laser power. Variations in the temperature of the laser structure also affect the laser power. It is therefore essential to keep the injection current and the temperature of the laser structure as constant as possible [3]. But even if all these conditions are satisfied, there may still be marked fluctuations in the laser power, if some of the laser radiation is reflected back into the internal cavity. This can happen to some extent when information is read from a Compact Disc or LaserVision disc and when a laser is connected to an optical fibre for communications purposes. Depending on the amount of radiation reflected back to the internal cavity and on the phase difference between the reflected radiation and the radiation in the laser, the optical field can become perturbed, giving large fluctuations of the laser power. We have studied this effect in the experimental arrangement described below (see fig. 6).

An image of the laser radiation (the 'spot') is formed by two lenses L_1 and L_2 on the feedback reflector M , which simulates a reflecting element such as a plate or the end of a fibre. The amount of feedback (the feedback fraction r) can be determined with the aid of filters located between the laser and the feedback reflector. The phase of the feedback radiation can be varied by moving the feedback reflector M and the lens L_2 . A part of the laser beam is split off by a beam splitter B for measuring the spectrum, the power, the noise and the coherence of the laser radiation.

The output power and the spectrum of the laser depend on the current through the laser structure (see fig. 5). We shall have to take this current dependence

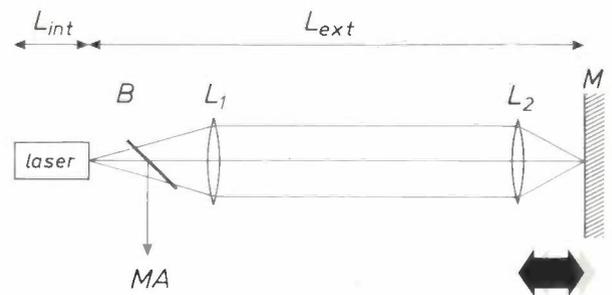
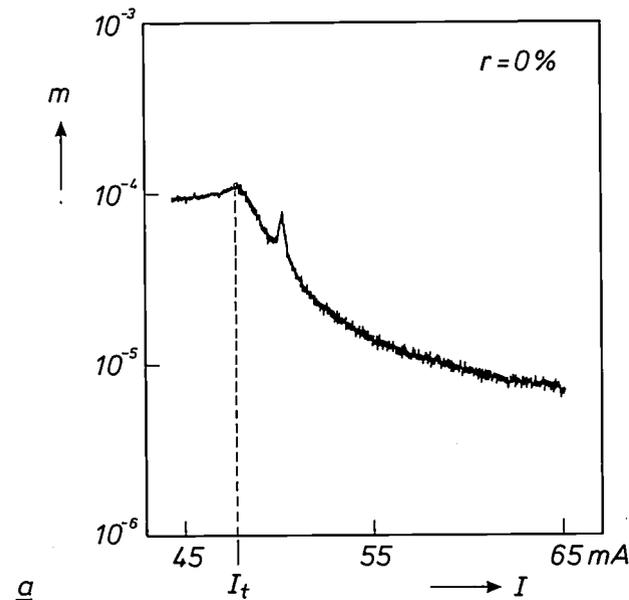


Fig. 6. Diagram of an experimental arrangement for determining fluctuations in the laser power. An image of the laser light is produced by two lenses L_1 and L_2 on the feedback reflector M . The feedback fraction can be varied with the aid of filters placed between the two lenses. The phase of the feedback radiation can be varied by moving the feedback reflector M and the lens L_2 . A part of the laser beam is extracted through a beam splitter B for measuring the spectrum, the noise, the coherence, the wavelength and the power (MA).

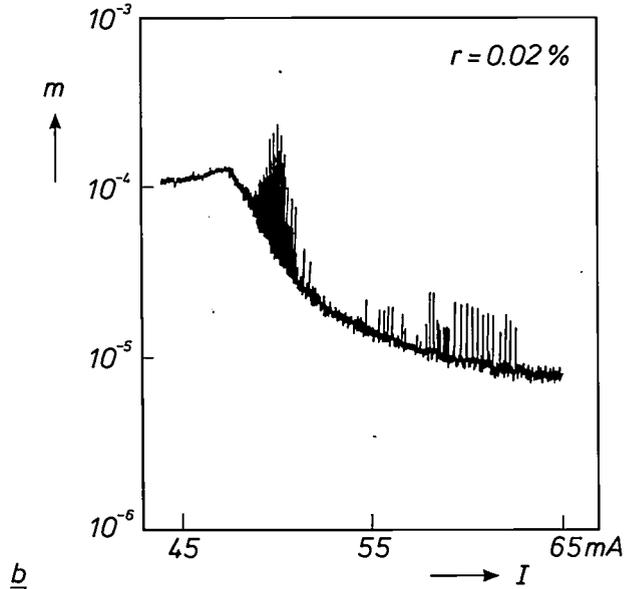
[3] In our experiments we do indeed keep the temperature of the laser structure as constant as possible. In systems containing lasers a feedback loop with a photodiode controls the injection current of the laser to give constant output power.

into account if we want to study the fluctuations in the laser power as a function of the feedback parameters.

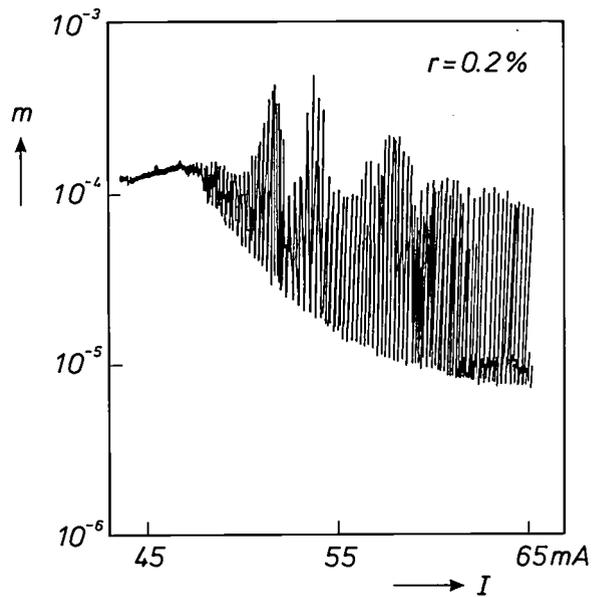
The laser power is measured with a P-I-N diode. The output signal of the diode is split to give a mean-value term ($\langle P \rangle$) and a fluctuating term ($\Delta P(t)$). In the experimental arrangement the root-mean-square value of $\Delta P(t)$, $\sqrt{\langle (\Delta P(t))^2 \rangle}$, is determined in a band of about 3 kHz at a frequency of 200 kHz. The modulation index ^[4] m is defined as the ratio of this r.m.s. value to the mean value $\langle P \rangle$. In *fig. 7* this modulation index is shown as a function of the current through the laser structure for a number of values of the feedback fraction. The phase of the feedback radiation is varied by moving the feedback reflector through a range of a few wavelengths. In this way minimum and maximum values for the modulation index are found. *Fig. 7a* shows the modulation index for a laser with no feedback. In general terms, the modulation index decreases with the injection current above the threshold value. (With no feedback the laser power, the d.c. component of the diode signal, increases while the fluctuations, the a.c. component of the diode signal, remain constant.) The abrupt change in the laser power at the threshold current appears here as a change in slope. Just above the threshold current there is a peak in the modulation index. It follows from wavelength measurements that this occurs at the current at which the laser ceases to emit at more than one wavelength (multimode behaviour) and starts to emit at a single wavelength (single-mode behaviour).



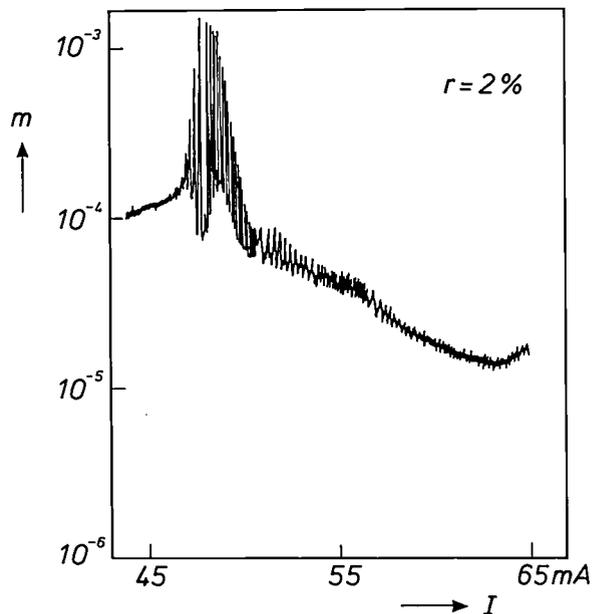
a



b



c



d

Fig. 7. The modulation index as a function of the injection current Δ for different values of the feedback coefficient ($r = 0, 0.02, 0.2$ and 2%). In these measurements the length of the external cavity is varied through a range of a few wavelengths. In this way maximum and minimum values are found for the modulation index.

As feedback is increased (figs 7b to d) we see that the modulation index at this current value increases. Further above the threshold current the noise also increases with increasing feedback (figs 7b and c). Wavelength measurements show that the laser does not really operate in a single mode in these circumstances but that radiation at a single wavelength is alternated with radiation at other wavelengths. These are not the wavelengths determined by the length of the internal cavity. The separations between the possible wavelengths are in fact much less than the separations between the laser modes. The fluctuations in the intensity of the laser radiation arise because the laser jumps from mode to mode, and these fluctuations account for the higher modulation index.

At even higher feedback fractions ($r \approx 2\%$) the modulation index above the threshold current, although larger than with no feedback, is no longer so strongly phase-dependent. Wavelength measurements show that the laser spectrum has a number of stable modes. At this feedback value, however, there is a distinct increase in the noise near the threshold current. In the applications we are concerned with here, in which lasers are operated well above the threshold current, this contribution to the noise is not significant. In what follows we shall not therefore consider the threshold noise, and will take a closer look at the noise at low feedback levels, as can be seen in fig. 7b and c.

At a current above the threshold value the laser without feedback emits virtually monochromatic radiation. The wavelength of the radiation is one of the possible wavelengths determined by the length of the internal resonant cavity. The introduction of the feedback reflector has the effect of producing a second resonant cavity that also has a number of laser modes. The wavelengths of these modes are determined by the length of the external cavity and come between the wavelengths of the internal modes (see fig. 8). The presence of the external modes may cause laser light to be generated at a different wavelength from the wavelength corresponding to the internal mode. What happens depends on two factors, the amount of feedback and the difference in wavelength between an internal mode and the nearest external mode. Roughly speaking, as the feedback increases, the closer the ultimate wavelength becomes to one of the external modes. The wavelength difference between internal and external modes can be very precisely controlled by small variations in the position of the feedback reflector. With a symmetrical adjustment — i.e. with the internal mode exactly midway between two external modes — and sufficient feedback, the laser wavelength can go towards that of either of the modes. The

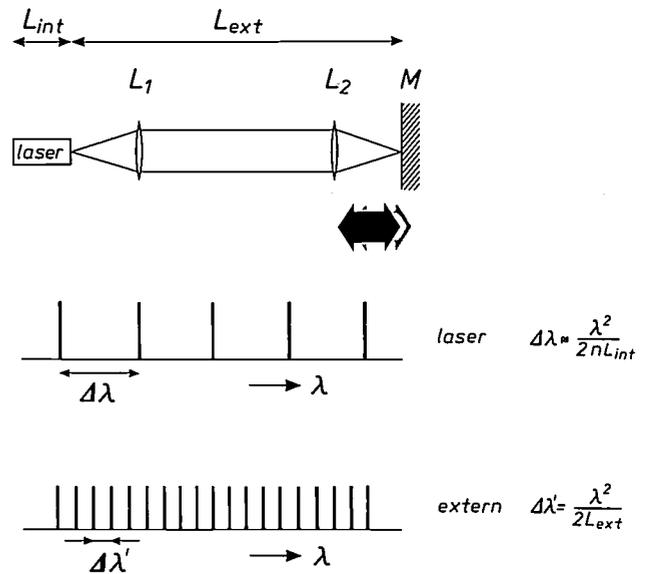


Fig. 8. Diagram showing the relative positions of the possible wavelengths of a laser with feedback. With no feedback the laser operates at one of the wavelengths indicated as laser modes. These wavelengths are fairly widely spaced. The effect of the feedback is that the wavelength at which the laser ultimately operates is a suitable compromise between a wavelength close to a laser mode and one that is close to an external mode.

actual 'choice' is not important here. We should note, however, that there is a particular feedback value at which the laser finds it hard to 'choose' between the two alternatives. The wavelength then oscillates at random between the two extremes, thus producing a great deal of noise.

Fig. 9 shows how the wavelength and power of the laser radiation change as a function of the phase of the feedback radiation, for three values of the feedback. The amount of radiation fed back is indicated by a dimensionless parameter C , which will be discussed later. If little radiation is fed back (fig. 9a) the variation in wavelength is small. If there is more feedback (fig. 9b) the effects are more pronounced. At even higher feedback values the laser wavelength may change so much that there is laser action on external modes (fig. 9c).

The effects that occur (corresponding to the results of fig. 9a, b or c) when there is feedback of laser radiation in the internal cavity are determined by the value of the feedback and the matching of the internal and external cavities (i.e. how close the external modes are to the wavelength of the laser without feedback).

[4] Determining the modulation index is a method of measuring the noise in a band of a particular bandwidth B . Another possible measure of the noise is the Relative Intensity Noise (RIN), which is the amount of noise per unit frequency, defined as $\langle \Delta P^2(\nu) \rangle / \langle P \rangle^2$. If the noise is white, i.e. independent of frequency, the relation between the two quantities is: $m = \sqrt{RIN \times B}$.

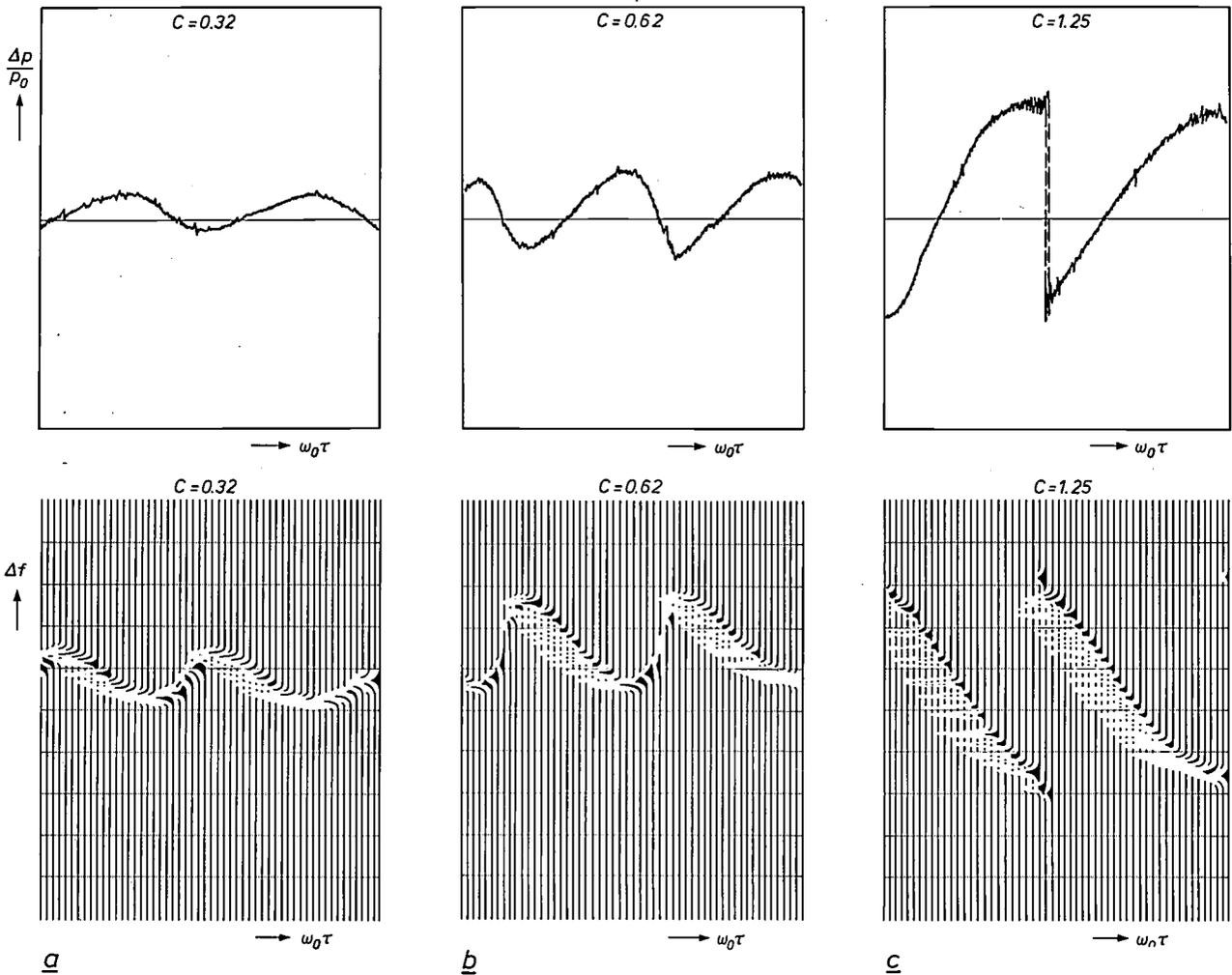


Fig. 9. The change in frequency (Δf) and the relative power change ($\Delta P/P_0$) of the laser radiation as a function of the phase ($\omega_0\tau$) of the feedback radiation for three values of the feedback. The factor C takes account of both the feedback fraction and the matching between the external and internal laser cavities. As C increases the feedback increases. a) $C = 0.32$. b) $C = 0.62$. c) $C = 1.25$.

Both conditions can be described by a single parameter C , which takes account of the amount of feedback and the length of the external cavity. (It will be shown in the next section that there is a theoretical basis for the definition of such a parameter.)

The system parameter C contains the product of the feedback factor γ and the transit time τ of the radiation in the external cavity, and also a term connected with the change in the refractive index of the laser material because of the change in the concentration of charge carriers. The feedback factor is defined as:

$$\gamma = f \frac{c}{2l_{d, \text{eff}}} (1 - R) \sqrt{\frac{r}{R}}, \quad (1)$$

where c is the velocity of light, $l_{d, \text{eff}}$ the effective length of the internal cavity, R the fraction of the radiation reflected by the laser reflectors in the internal cavity, r the fraction of the radiation reflected by the feedback reflector M , and f the fraction of the feedback radiation that actually arrives in the internal cavity of the laser.

It appears that there can be a maximum in the noise of the laser power for three values of this parameter (see *fig. 10*). The first maximum appears at $C \approx 1$. (This corresponds to a feedback fraction of only about 0.01% in the geometry of the arrangement used in *fig. 8*.) This maximum appears because the laser cannot choose between two symmetrical alternative wavelengths, as described above.

The second maximum appears because the laser jumps between more than two external laser modes ('mode hopping'). If the feedback is increased again there is also mode hopping between the internal laser modes (the third maximum). In this case the noise level is much higher than at lower feedback values because more modes are involved in the process. Finally, there is a stable situation in which the laser emits simultaneously at several wavelengths and a higher but constant noise level is present.

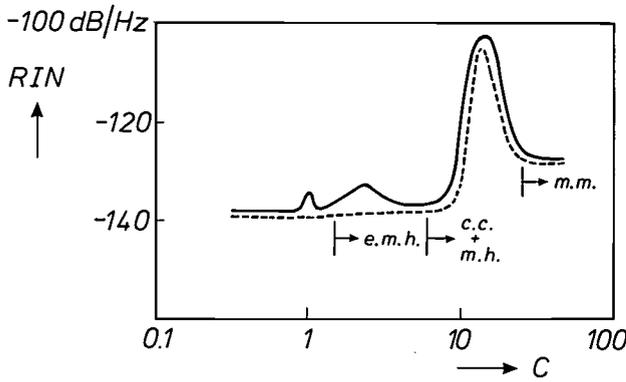


Fig. 10. The amount of noise^[3] as a function of the parameter C (which takes account of the feedback and the matching between the external and internal cavities). Varying the phase of the feedback radiation gives a maximum value of the noise (continuous curve) and a minimum value (dashed curve). The figure also shows the regions where there is mode hopping between several external laser modes (e.m.h.), the regions where there is coherence collapse (c.c.), and the regions where there is mode hopping between laser-cavity modes (m.h.). Beyond the third maximum the laser gives a multi-mode spectrum.

In principle the effects that cause the noise peaks can be described in the same way. In all cases the noise arises as a result of mode hopping in the laser. In the situation where the second maximum appears, there are only external modes, whereas in the situation where the third maximum appears there are both internal and external modes.

In the next section we shall give a theoretical description of the first situation (corresponding to the appearance of the first maximum). It will be seen that the results calculated from this theory for the phase-dependence of the wavelength and the power agree with the experimental results.

The behaviour of the laser with high feedback is in reality more complicated than indicated above. The wavelength measurements show that in addition to the mode hopping of the laser between more than two wavelengths determined by the external and internal cavities, there is a marked broadening of the laser line. The linewidth, which amounts to a few tens of MHz at low feedback, increases to about 20 GHz. The coherence length of the laser (which is inversely proportional to the linewidth) has therefore decreased enormously (by a factor of 1000). We call this effect 'coherence collapse'. The feedback radiation is no longer coherent with the radiation produced in the internal cavity. Both the change in the linewidth, and the fact that the noise is independent of the length of the external cavity lead to the conclusion that the phase relation between the feedback radiation and the radiation in the internal cavity has been completely distorted. The parameter C has no further significance in this region, and strictly speaking the scale in fig. 10 ought to be replaced from $C = 5$ by a scale that only shows the feedback fraction.

Optical feedback, theory

Two important quantities for laser action in the active layer are the light intensity and the associated concentration of charge carriers. Statements about laser

power and changes in it can be made if we know how these two quantities depend on the parameters of the laser cavity (refractive index, dimensions) and the current injected into the active layer. We shall examine the situation in a laser that emits monochromatic radiation, i.e. where the current in the active layer is greater than the threshold current.

The concentration of electron-hole pairs $n(t)$ in the active layer may change for three different reasons. In the first place the number of electron-hole pairs decreases as a result of spontaneous recombination. The number of electron-hole pairs that disappear per unit time in this way is inversely proportional to their mean lifetime T_1 . The second process that causes electron-hole pairs to disappear is stimulated recombination. This number (of pairs that disappear in this way) is proportional to the intensity P of the radiation field in the active layer. Finally, electron-hole pairs are generated by the passage of a current through the active layer. The resulting change in concentration is J/qd , where J is the current, q the electronic charge and d the thickness of the active layer. Equation (2) expresses the change in concentration of electron-hole pairs due to all three processes:

$$\frac{dn(t)}{dt} = -\frac{n(t)}{T_1} - G_A(n(t))P + \frac{J}{qd}. \quad (2)$$

The second term in this equation contains the intensity gain G_A , which in turn depends on the concentration of electron-hole pairs.

The change in the radiation field in the active layer can be described by considering the amplitude E of the optical field as a function of time. The light intensity in the cavity is equal to $|E|^2$. There are three contributions to the change in the field amplitude. The first is proportional to the gain $G_A(n(t))$, the second to the constant losses that arise because radiation leaks from the cavity, Γ_0 , and the third is proportional to a dispersion factor $G_D(n(t))$ that takes account of the refractive-index variations due to variations in the concentration of electron-hole pairs. If there is any feedback radiation returning to the internal cavity after a time τ , the equation contains a further term that takes the feedback factor (γ) and the phase ($\omega_0\tau$) of this radiation into account. The equation for the change in the amplitude of the optical field is:

$$\frac{dE(t)}{dt} = \frac{1}{2} [G_A(n(t)) - \Gamma_0 + iG_D(n(t))] E(t) + \gamma E(t - \tau) e^{-i\omega_0\tau}. \quad (3)$$

The factor $\frac{1}{2}$ appears in this equation because we are interested in the amplitude of the field here, whereas the expressions between brackets relate to the power.

The angular frequency of the laser with no feedback is ω_0 . The quantity $\omega_0\tau$, which occurs in the exponent in (3), is the feedback phase referred to earlier. In this treatment we neglect multiple reflections. This is justified because the feedback factor γ is very small in the experiments that our model is intended to describe.

We assume that the system without feedback is in a stable state characterized by an intensity P_0 , an electron-hole pair concentration n_0 and an angular frequency ω_0 . We thus have:

$$G_A(n_0) = \Gamma_0, G_D(n_0) = 0 \text{ and } \frac{n_0}{T_1} = -\Gamma_0 P_0 + \frac{J}{qd}. \quad (4)$$

In a situation where there is feedback we can linearize the gain G_A and the dispersion G_D around n_0 :

$$G_A(n) = \Gamma_0 + \xi(n - n_0), \quad G_D(n) = \eta(n - n_0). \quad (5)$$

This linearization permits us to express equations (2) and (3) that describe our model in terms of parameters that represent the situation with no feedback:

$$\frac{d}{dt} [n(t) - n_0] = -\left(\frac{1}{T_1} + \xi P\right)[n(t) - n_0] - \Gamma_0 [P(t) - P_0], \quad (6)$$

$$\frac{dE(t)}{dt} = \frac{1}{2}(\xi - i\eta)[n(t) - n_0]E(t) + \gamma E(t - \tau)e^{-i\omega_0\tau}. \quad (7)$$

We can thus describe the interaction between the optical field and the charge carriers by two inter-related differential equations, (6) and (7). If these can indeed describe our experimental results as discussed in the previous section, then they must permit single-frequency solutions that are stable under certain conditions. Closer examination^[5] of our equations shows that such solutions do exist and that they are of the form $E(t) = \sqrt{P} \exp(i\Delta\omega t)$ and $n(t) = \bar{n}$ where P , \bar{n} and $\Delta\omega$ are independent of time. The frequency shift $\Delta\omega$ resulting from the feedback, the power P and the concentration of charge carriers \bar{n} are characterized by:

$$\Delta\omega\tau = C \sin[\phi_0 - (\omega_0 + \Delta\omega)\tau], \quad (8)$$

$$P = P_0 + 2\gamma \left(P_0 + \frac{1}{\xi T_1} \right) \frac{\cos(\omega_0 + \Delta\omega)\tau}{\Gamma_0 - 2\gamma \cos(\omega_0 + \Delta\omega)\tau} \approx P_0 + \frac{2\gamma P_0}{\Gamma_0} \left(\frac{J/J_t}{J/J_t - 1} \cos(\omega_0 + \Delta\omega)\tau \right), \quad (9)$$

$$\bar{n} = n_0 - \frac{2\gamma}{\xi} \cos(\omega_0 + \Delta\omega)\tau. \quad (10)$$

The constant $C (= \gamma\tau/\cos\phi_0$, where $\tan\phi_0 = \eta/\xi$) in these solutions indicates the extent to which the intrinsic behaviour of the laser (frequency etc.) changes because of the feedback. We find in this parameter the feedback factor γ , the transit time τ of the radiation through the external cavity, and a quantity ϕ_0 that depends on the laser material and the structure of the laser. This justifies the introduction of the parameter C in the description of the experimental results (see page 297). We can compare the measured and theoretical behaviour of the frequency and power as a function of the phase of the feedback radiation if the equations (8), (9) can be solved for the quantities $\Delta\omega$ and P as a function of the phase $\omega_0\tau$. The power P has a cosine dependence on the phase of the feedback radiation (see eq. (9), subject also to the condition that the diode losses Γ_0 are much larger than the value of the feedback factor γ). A closer examination of the stability of the solutions shows that there are two different regions. If $C < 1$, i.e. if there is little feedback, there is only one solution for $\Delta\omega$, but if $C > 1$ there are more possibilities, which can be related to external modes.

In *fig. 11* the solutions of equations (8) and (9) are shown as $\Delta\omega\tau$ (a measure of the frequency shift) and $\cos(\Delta\omega + \omega_0)\tau$ (the power variation apart from a constant pre-factor) as a function of the phase $\omega_0\tau$ for a number of values of C . The continuous lines give the stable solutions and the dashed lines give the unstable solutions. In these calculations the ratio of the imaginary part to the real part of the refractive index, η/ξ , is equal to -4 , a value close to the last known experimental value and characteristic of semiconductor lasers. This results in a value of -76° for ϕ_0 , so that $\Delta\omega$ and P are represented as a function of $\omega_0\tau$ by practically the same curves, but with opposite sign. Low values of C give a single stable solution ($C = 0.5$). If $C = 1$, both curves have a vertical slope at a particular value of the phase. At higher feedback levels ($C = 3$) there are three possible values for the power and the frequency shift in a particular phase range. As the phase goes through this range in both directions a hysteresis loop is described as indicated by the arrows. Laser action (stimulated emission) only occurs if the solution is stable. If small temperature and current fluctuations lead to an unstable solution, the laser jumps to the next stable solution. If C increases, more stable solutions become possible and the behaviour of the laser becomes more complicated.

We can determine the fluctuations in the power as a function of the phase by calculating the derivative of the power with respect to that phase. This is justified for the low-frequency fluctuations we are concerned with here (the frequency of the fluctuations in our ex-

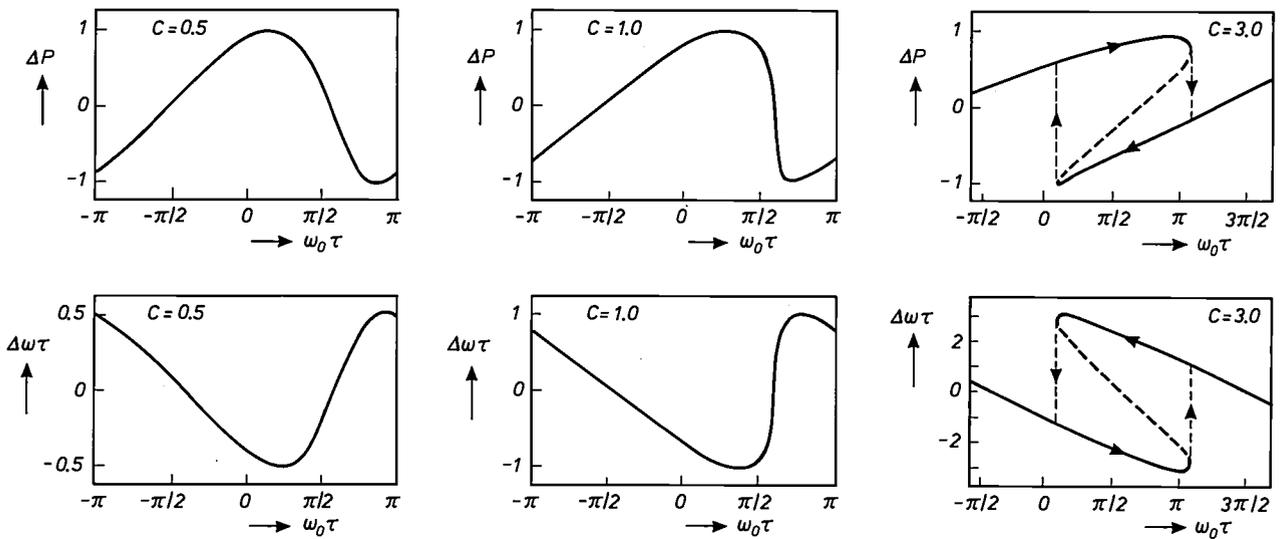


Fig. 11. The frequency shift $\Delta\omega\tau$ and the power P as a function of the phase $\omega_0\tau$ of the feedback radiation for different values of the parameter C .

periments is less than about 10^6 Hz. The calculation shows that most of the fluctuations will be in the neighbourhood of $C=1$.

The calculated and measured results agree well. The feedback in the experiment can be calculated from a comparison of measured and calculated curves by determining C for the experimental curve. This is found to agree well with the feedback value as determined from the physical characteristics of the experimental arrangement (such as reflection coefficients of the laser and the feedback reflector and the length of the external cavity). It follows from the theoretical analysis that we can expect a large amount of noise in the vicinity of $C=1$, which corresponds to the results of the experiments. The theory does not account for the large amount of noise indicated in fig. 10 for higher values of C ; the parameter C cannot be used in the same way in this range since the linearization that we performed in equations (4) and (5) is no longer justified.

Tackling the noise problems

The good agreement found between the results of the experiments and the theoretical calculations indicate that we have gained a better understanding of the processes that are responsible for noise in the intensity of the laser radiation. But this does not mean that we have solved the problem introduced by this noise in the applications mentioned here. We can however indicate the conditions in which the operation of the laser will be least affected by noise. This is the case, for example, when the laser shows multimode behaviour (see fig. 10), which can arise in various ways. The noise level is then constant, although a little higher

than if there is no feedback. This can be taken into account by making the differences in intensity that contain the information substantially larger than this constant noise level.

In the first place, multimode behaviour of the laser can be obtained by providing high feedback. The effect of this is to reduce the threshold current and to increase the noise around the threshold (see fig. 7d). This presents no new problems because the lasers are operated well above the threshold current. Secondly, multimode behaviour can be induced by high-frequency modulation of the injection current. A third way of inducing multimode behaviour in the laser is to modify the structure of the laser. The laser we have described here is one whose characteristics are between those of a 'gain-guided' laser (with an optical gain profile), and an 'index-guided' laser (with a step in the refractive index). This implies that the laser cavity is partly formed by differences in the refractive index of the materials forming the layers of this structure. A laser of the gain-guided type, in which the cavity in the x -direction is defined by making only a restricted region suitable for the passage of current, gives multimode behaviour and therefore no noise problems arise. The threshold current of a gain-guided laser is higher, however, than in the laser structure discussed in this article. This higher threshold current introduces further problems connected with cooling and laser life. The beam quality of a gain-guided laser is also poorer because of astigma-

[6] More information about these calculations can be found in G. A. Acket, D. Lenstra, A. J. den Boef and B. H. Verbeek, The influence of feedback intensity on longitudinal mode properties and optical noise in index-guided semiconductor lasers, IEEE J. QE-20, 1163-1169, 1984.

tism. Consequently there is sufficient inducement in practice to use lasers with characteristics between those of the two types.

Another way of minimizing noise is to design a laser structure in which the system parameter C has the lowest possible value. Small values of C (< 1) result in a low noise level; see fig. 10. The small value of C can be obtained, for example, by increasing the feedback of the laser reflectors (C is proportional to $1 - R$; see equation (1)).

In some applications noise is largely avoided by using lasers with reflectors that return nearly all the

radiation into the internal cavity ($R \approx 1$), or lasers that are forced into multimode behaviour by modulation of the injection current.

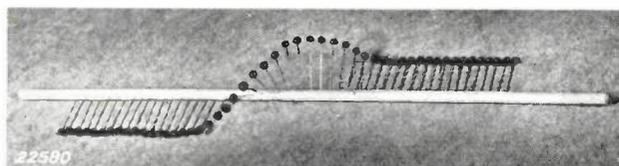
Summary. Fluctuations in the optical power of semiconductor lasers can arise because of radiation returning to the laser cavity after reflection from a surface outside the cavity. These fluctuations are undesirable in many applications in which information is optically transmitted as variations in laser power. Experiments and calculations have been carried out to gain a better understanding of the ways in which these fluctuations are affected by the various laser parameters. Experimental and theoretical results agree well. Conditions can now be indicated in which there is very little variation in laser power as a result of feedback of optical power.

1937

THEN AND NOW

1987

Permanent magnets

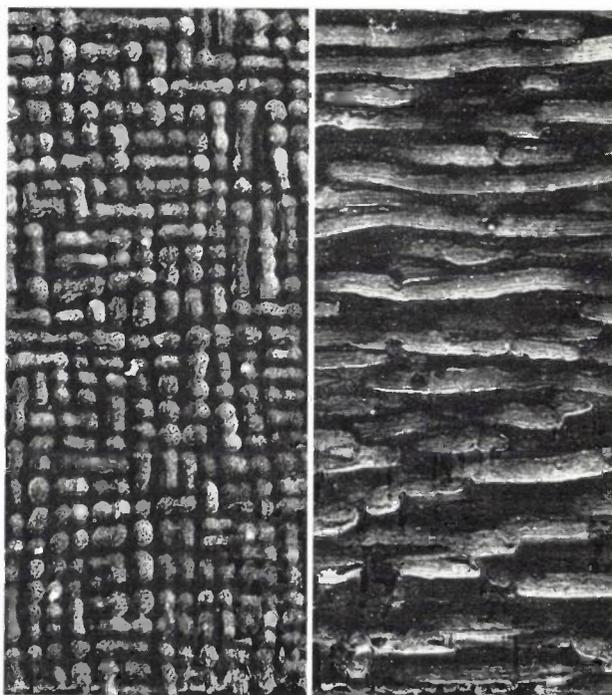


Although permanent magnets have been used for many centuries — in the compass — a good model for explaining ferromagnetic behaviour has only been available since the beginning of the twentieth century. This was when it was first understood that magnetic material consists of small domains of magnetic moments all pointing in the same direction, and that the domains could change their orientation by the movement of a 'Bloch wall' (see the illustration^[*] at the top left). With a better understanding of magnetism, it became possible to control the microstructure of alloys and oxides of the ferromagnetic elements — iron, cobalt, nickel — so that their $(BH)_{\max}$ -values were ten or twenty times those of the classic magnet steels. 'Ticonal G' from 1937, 'Ferroxdure II' from 1954 and 'Ticonal XX' from 1956 were all results of such work at Philips on the perfection of magnetic materials (see the two electron micrographs on the right, taken in perpendicular directions in 'Ticonal XX').

Later on, elements from the rare-earth group were applied, as well as the traditional ferromagnetic elements. This led to the discovery of SmCo_5 in 1968 and $\text{Nd}_2\text{Fe}_{14}\text{B}$ in 1983. Philips are about to apply the latter material in the pick-up of the CD player.

A striking picture of the progress during the last twenty or thirty years is obtained by comparing the load (the brass rods, colour photo) that the various kinds of magnets, of the same dimensions, can support if they levitate at the same distance above an 'opposite pole'.

'Ticonal' does not appear in this picture, since its best features, the combination of a high saturation magnetization and a small coercivity (just the opposite of 'Ferroxdure', which has a low saturation magnetization and a high coercivity), do not show to their best advantage in such an experiment.



[*] From Philips Technical Review, August 1937.

Metastable phases and thermodynamic equilibrium

J. J. van den Broek and A. G. Dirks

The hardening of steel is one of the oldest and best-known processes in which a metastable crystalline phase occurs. A supersaturated solution of carbon in iron is formed, and the associated crystalline structure is called martensite. Metastable amorphous phases in metal alloys have been the subject of keen interest in recent years because of their interesting properties, which arise because they have no regular atomic order and no grain boundaries. The occurrence of metastable phases, either amorphous or crystalline, can in many cases be predicted from the tendency of every system to seek a minimum of the Gibbs free energy. For predictions of this kind the familiar phase diagram, which applies to a state of stable equilibrium, can be a valuable aid.

Introduction

Until about 1934, when amorphous metal was first successfully produced from the vapour phase, only the crystalline form of metals was known. Technical applications of amorphous metal did not emerge until the seventies, when it became possible to produce this material in the form of a metallic ribbon. The method consisted in rapidly cooling metal in the liquid state by squirting it on to a copper wheel rotating at high speed. It seemed likely that the absence of preferred directions for the magnetic moment would give metals with better magnetic properties, and that the absence of grain boundaries would yield metals of much greater strength [1].

Amorphous metals have not yet become as widely applied as first seemed likely. In the past fifteen years, however, there has been renewed interest. This is partly due to the emergence of a new application in the form of thin films for magneto-optical recording [2]. These films are made by evaporating the components on to a cooled substrate, so that atomic disorder is 'frozen in' [3]. Another innovation has been the introduction of mechanical alloying, which can be used, it appears, to permit the economical use of amorphous metal for exceptionally strong mechanical compo-

nents. In this method the constituents are ground to a fine powder and then intimately mixed [4]. An industrial application announced recently is the use of amorphous metal ribbon in transformer cores. The losses are lower, so that energy can be saved and the transformers can be much smaller.

Amorphous metals are always metastable, which means that a fairly significant disturbance, such as a sudden increase in temperature, can make the metal become crystalline. Evaporation from the vapour phase can produce both metastable *amorphous* alloys and metastable *crystalline* alloys. We have even been able to make homogeneous mixtures of metals that are normally so immiscible that two phases occur both in the solid state and even in the liquid state (like oil in water) [5].

According to the laws of thermodynamics a system at constant temperature and pressure tends towards an equilibrium state with a minimum of the Gibbs free energy [6]. The Gibbs free energy G is defined as:

$$G = U - TS + pV, \quad (1)$$

where U is the internal energy, T the absolute temperature, S the entropy, p the pressure and V the volume. The quantity $U + pV$ is called the enthalpy (H). In all

known metal alloys at room temperature the Gibbs free energy of the amorphous phase is greater than that of the most stable crystalline phase, i.e. the phase with the lowest Gibbs free energy. This explains why an amorphous phase is always metastable. If the mobility of the atoms at room temperature is insufficient to cause the amorphous structure to become crystalline, the amorphous structure can maintain itself.

The phase diagram, an invaluable aid in phase theory, shows the phases that occur at different temperatures and compositions (fractions) of the components. These phases are the various crystalline phases, the liquid phase or phases and sometimes the gas phase. The phase diagram is directly related to the (G, x) diagram, which gives the Gibbs free energy G at a given temperature for the various phases as a function of the mole fraction x ^[7]. N. Saunders and A. P. Miodownik have also pointed out that the (G, x) diagram can often be used for predicting whether metastable phases will occur^[8], especially for alloys that have been made by depositing the components on to a cold substrate from the vapour phase. The mobility of the atoms is then usually so low that the homogeneous distribution due to the evaporation remains unchanged. At a particular composition the phase obtained is the one that has the lowest Gibbs free energy, provided that the crystal structure is not too complicated to form on the cold substrate. As a rule, then, two separate phases do not occur, although this often happens when a state of stable equilibrium is reached.

The difficulty with predicting metastable amorphous and crystalline phases from the (G, x) diagram is that for most combinations of metals the diagram is not known and is also difficult to calculate^[9]. A method of approximation for predicting amorphous phases, which requires no calculations, is known from phase theory^[7]. In this method the location of a composition region where amorphous alloys may occur is given by extrapolating liquidus curves in the phase diagram. If the phase diagram of an alloy system is known, this method soon gives a result, and is just as useful for making predictions as the (G, x) diagram.

In the rest of the article we shall first consider the phase diagram and its relation to the (G, x) diagram. We shall then show how metastable phases can be predicted with the (G, x) diagram. Next, we shall deal with the prediction of amorphous phases with the aid of the phase diagram. Both methods will be illustrated by practical examples. Finally we shall look at the complicated phase relationships that can arise in a transition region between a state of stable thermodynamic equilibrium and a state of metastable equilibrium.

The phase diagram and its relation to the (G, x) diagram

First of all we shall recall the basic principles of the phase diagram. *Fig. 1* shows a (G, x) diagram and the phase diagram for two metals A and B that are completely soluble in each other in the liquid and solid states^[7], e.g. silver and gold. (The mole fraction x is the ratio of the number of atoms B to the total number of atoms in the mixture.) The complete mutual solubility is connected with the identical crystal structures of pure gold and silver and the virtually identical atomic radii of both elements. In mixed crystals silver and gold atoms can therefore be substituted for one another in any mixing ratio.

The shape of the curves in the (G, x) diagram for the various phases depends on the temperature T , since

- [1] J. Kramer, Über nichtleitende Metallmodifikationen, *Ann. Phys.* **19**, 37-64, 1934;
R. Hilsch, Amorphous layers and their physical properties, in: V. D. Fréchet (ed.), *Non-crystalline solids*, Wiley, New York 1960, pp. 348-373;
F. E. Luborsky (ed.), *Amorphous metallic alloys*, Butterworths, London 1983;
P. Duwez, Structure and properties of alloys rapidly quenched from the liquid state, *Trans. Metall. Soc.* **60**, 607-633, 1967;
K. H. J. Buschow, Research on amorphous alloys, *Philips Tech. Rev.* **42**, 48-57, 1985.
- [2] J. W. M. Biesterbos, Properties of amorphous rare earth-transition metal thin films relevant to thermomagnetic recording, *J. Physique* **40** (Colloque C5), C5/274-C5/279, 1979;
M. Hartmann, B. A. J. Jacobs and J. J. M. Braat, Erasable magneto-optical recording, *Philips Tech. Rev.* **42**, 37-47, 1985.
- [3] S. Mader, Metastable alloy films, *J. Vac. Sci. & Technol.* **2**, 35-41, 1965.
- [4] R. B. Schwarz and W. L. Johnson, Formation of an amorphous alloy by solid-state reaction of the pure polycrystalline metals, *Phys. Rev. Lett.* **51**, 415-418, 1983;
E. Hellstern and L. Schultz, Amorphization of transition metal Zr alloys by mechanical alloying, *Appl. Phys. Lett.* **48**, 124-126, 1986;
A. W. Weeber, K. van der Meer, H. Bakker, F. R. de Boer, B. J. Thijssse and J. F. Jongste, The preparation of amorphous Ni-Zr powder by mechanical alloying, *J. Phys. F* **16**, 1897-1904, 1986.
- [5] A. G. Dirks and J. J. van den Broek, Metastable solid solutions in vapor deposited Cu-Cr, Cu-Mo, and Cu-W thin films, *J. Vac. Sci. & Technol. A* **3**, 2618-2622, 1985;
J. J. van den Broek, A. G. Dirks and J. L. C. Daams, Phase relationships in binary alloy thin films, *Proc. Int. Symp. Trends and New Applications in Thin Films*, Strasbourg 1987, pp. 421-425.
- [6] J. D. Fast, Entropy in science and technology, *Philips Tech. Rev.* **16**, 258-269 and 298-308, 1955.
- [7] J. L. Meijering, Phase theory, *Philips Tech. Rev.* **26**, 12-26 and 52-60, 1965.
- [8] N. Saunders and A. P. Miodownik, The use of free energy vs composition curves in the prediction of phase formation in codeposited alloy thin films, *CALPHAD* **9**, 283-290, 1985;
N. Saunders and A. P. Miodownik, Thermodynamic aspects of amorphous phase formation, *J. Mater. Res.* **1**, 38-46, 1986.
- [9] The enthalpy H can be calculated with the aid of A. R. Miedema's models. Since the term TS in eq. (1) for the Gibbs free energy is small at room temperature, these models can sometimes be used for determining the Gibbs free energy.
See: A. R. Miedema, The atom as a metallurgical building block, *Philips Tech. Rev.* **38**, 257-268, 1978/79;
A. K. Niessen and A. R. Miedema, The enthalpy effect on forming diluted solid solutions of two 4d and 5d transition metals, *Ber. Bunsenges. Phys. Chem.* **87**, 717-725, 1983;
A. R. Miedema, P. F. de Châtel and F. R. de Boer, Cohesion in alloys — fundamentals of a semi-empirical model, *Physica* **100B**, 1-28, 1980.

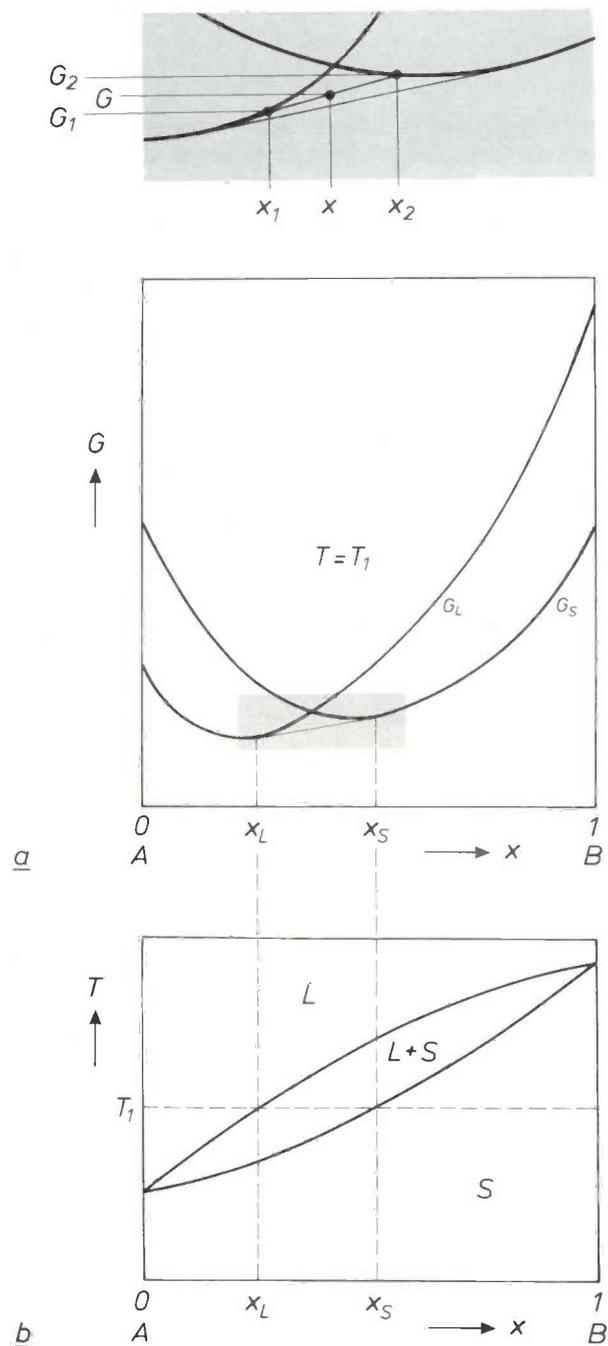


Fig. 1. Principles of the phase diagram. *a*) A (G, x) diagram, the Gibbs free energy G as a function of the mole fraction x of B in the mixture of A and B, at a temperature $T = T_1$, G_L Gibbs free energy of the liquid phase L . G_S the same, for the solid phase S . x_L and x_S boundaries of the composition range where, in stable equilibrium, there are two phases, the liquid phase and the solid phase. In this range the Gibbs free energy for a given mole fraction follows from the common tangent of the (G, x) curves. The enlarged diagram at the top shows that the Gibbs free energy G of a mixture of two phases (G_1, x_1) and (G_2, x_2) at a fraction x follows from the line connecting the corresponding points. In a state of stable equilibrium the appropriate combinations of G and x lie on the common tangent. *b*) The phase diagram for two elements that are completely soluble in each other both in the liquid state and the solid state. The corresponding (G, x) diagram shown in (a) applies at a temperature T_1 . The phase diagram can be thought of as being formed from a number of (G, x) diagrams at different temperatures. The region $L+S$ is two-phase.

the Gibbs free energy G is a function of temperature, of course; see eq. (1). Fig. 1*a* shows the (G, x) diagram for a temperature T_1 higher than the melting point of pure A and lower than that of pure B. At this temperature, A with little B is liquid, B with little A is solid. To the left of the point where the curves intersect, the liquid phase has a lower Gibbs free energy than the solid phase. The liquid phase is therefore more stable here. On the right the same is true for the solid phase. In a state of stable equilibrium there is a composition region near the intersection point of the two (G, x) curves where the solid and the liquid phases form a two-phase mixture. The boundaries x_L and x_S of this range follow from the location of the points of contact with the common tangent. In stable equilibrium the Gibbs free energy of the two-phase mixture is determined by the point on the common tangent that corresponds to the mole fraction x of the alloy under consideration.

It can be seen from the inset in fig. 1*a* that the common tangent of the (G, x) curves must be drawn to find the Gibbs free energy of a mixture in the range $x_L < x < x_S$. Suppose that at a fraction x a liquid phase forms with Gibbs free energy G_1 and fraction x_1 and a solid phase with G_2 and x_2 . The Gibbs free energy of the mixture of liquid and solid is then found for the fraction x by drawing a straight line between the points (x_1, G_1) and (x_2, G_2) ^[7]. The Gibbs free energy of the two-phase mixture is lower than that of homogeneous liquid or homogeneous solid solutions of the same composition. However, of all the possible lines that connect points on the curves, the common tangent to the curves at the fractions x_L and x_S represents the lowest Gibbs free energy that is possible for intermediate fractions.

The stable equilibrium at a temperature T_1 can be described as follows; see fig. 1*a*. Mixtures that are rich in A are liquid when $x < x_L$. Mixtures that are rich in B are solid when $x > x_S$. Mixtures with $x_L < x < x_S$ consist both of a liquid of mole fraction x_L and of single-phase solid solution of mole fraction x_S . According to the 'lever' rule, the amount of solid solution varies from 0% at $x = x_L$ to 100% at $x = x_S$.

The phase diagram corresponding to the system A-B in fig. 1*a* is shown in fig. 1*b*. There is a region, indicated as $L+S$, in which the liquid and the solid phases coexist. The boundaries of this region are given by the contact points of the common tangents of pairs of curves in (G, x) diagrams at different temperatures, as demonstrated for the temperature T_1 . In region S there is a single phase of solid solutions which, unlike pure A or pure B, do not have a melting point but a melting range. The phase diagram describes a state of stable thermodynamic equilibrium.

Fig. 2 gives a (G, x) diagram and the phase diagram for a multiple eutectic system A-B in which compounds A_2B and AB_2 occur. The structures α , β , γ , and δ of pure A, pure B, and the two compounds respectively are different, so that at low temperatures only phases of nearly pure A, nearly pure B or a compound occur. In the (G, x) diagram, which corresponds here to a temperature (T_2) with solid phases only, curves are shown for the Gibbs free energies G_α , G_β , G_γ and G_δ . Here again, the common tangents of the (G, x) curves are shown, resulting in three broad two-phase regions in the phase diagram. Crystal structures of the compounds A_2B and AB_2 only form at the virtually fixed ratios A/B of 2 and 1/2, at which the atoms A and B occupy fixed places in the crystal structures. Examples of systems with compounds are copper-zinc, samarium-cobalt and nickel-zirconium. We shall return to the system nickel-zirconium presently.

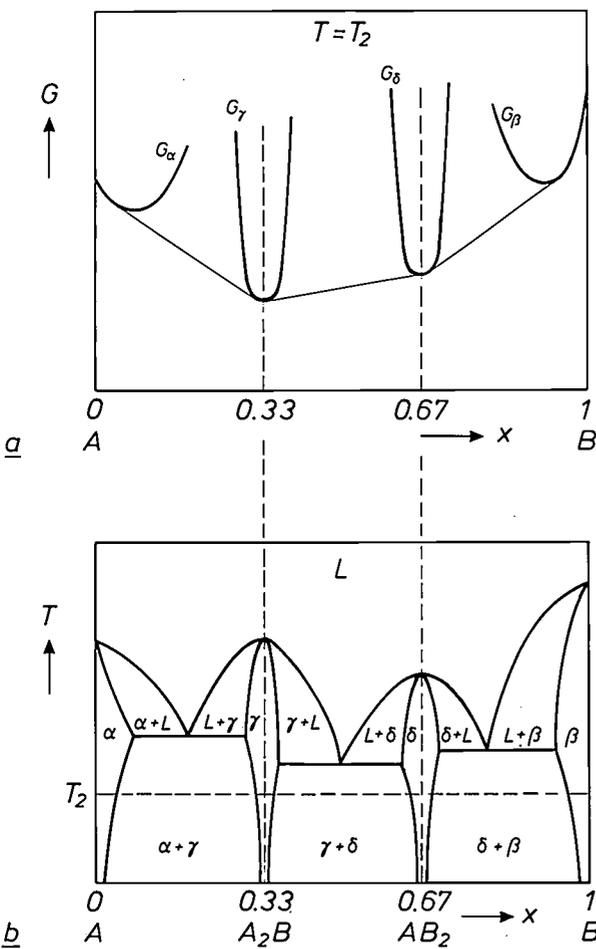


Fig. 2. a) A (G, x) diagram and b) the phase diagram for a multiple eutectic system A-B with compounds A_2B and AB_2 . The crystal structures of pure A and B are indicated as α and β , and those of the compounds as γ and δ . G_α , G_β , G_γ and G_δ are the Gibbs free energies of the various structures. The (G, x) diagram relates to a temperature T_2 . $\alpha + \gamma$, $\gamma + \delta$, $\delta + \beta$, $\alpha + L$, $L + \gamma$, $\gamma + L$, etc. are two-phase regions. L liquid phase.

Predicting metastable phases with the aid of the (G, x) diagram

In determining two-phase regions with the aid of the common tangent of (G, x) curves it is assumed that the atoms possess so much mobility that the two phases with different structures and mixing ratios are in fact formed. After rapid cooling from the liquid phase or deposition from the vapour phase (we shall use the term 'quenching' for both processes) this mobility is sometimes so small, however, that the separate phases cannot form. In such cases, therefore, there is no point in drawing the common tangent. Saunders and Miodownik have pointed out that instead of two separate phases a single phase occurs, which is usually the phase that, according to the (G, x) diagram, has the smallest Gibbs free energy [8].

This is illustrated in fig. 3. For the system A-B the (G, x) diagram in fig. 3a shows curves for the crystal structures α , β and γ . Fig. 3b shows the phases that are formed in a state of stable equilibrium. On either side of the two-phase region $\alpha + \beta$ there are single-phase regions α and β . The structure γ is not formed, because the (G, x) curve lies above the common tangent of the curves for G_α and G_β . Fig. 3c applies to a situation after quenching. The boundary of the regions α and β is defined by the point where the corresponding curves in fig. 3a intersect. It is assumed that no phase γ is formed; this may be the case if it is difficult for crystallization nuclei to form or if the crystallization nuclei grow too slowly, e.g. because the structure is too complicated. If the phase γ does form, we have the situation shown in fig. 3d. There are then three single-phase regions, whose boundaries are given by two intersection points in fig. 3a.

The probability that a third metastable phase, either crystalline or amorphous, will be formed in addition to α and β increases with the area of the grey-shaded region in fig. 3a between the two curves and their common tangent. This area is large, for example, when element B does not 'fit' so well into the structure α of element A, since the difference $G_\alpha^B - G_\beta^B$ will then be greater. This is the thermodynamic background to the 'structural difference rule' [10].

We can illustrate the theory above with data from Saunders and Miodownik [8]. Fig. 4a shows how the composition range for an amorphous phase is determined for the nickel-zirconium system. The continuous curves give the Gibbs free energy G at room temperature as a function of the fraction x for different

[10] B. X. Liu, Ion mixing and metallic alloy phase formation, Phys. Stat. Sol. A 94, 11-34, 1986.

crystalline phases and for the liquid phase. The small squares relate to the amorphous material. It can be seen that there is a relatively large difference between the Gibbs free energy of the amorphous phase and that of the liquid phase. This is a little surprising, since the usual assumption is that they are just about equal. The difference is probably a consequence of some short-range ordering of the atoms in the amorphous phase. The points marked + indicate the minimum Gibbs free energy for different compounds of Ni and Zr; see the phase diagram in fig. 4b^[11]. The atoms do not have sufficient mobility for these compounds to be formed, however, so that we can disregard these points. The range where the amorphous phase would

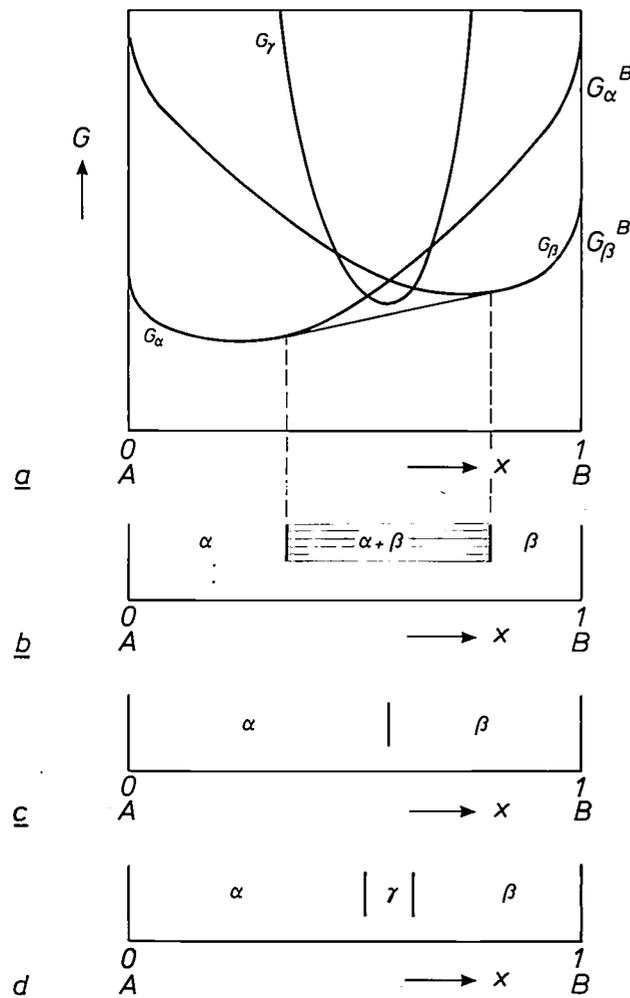


Fig. 3. a) (G,x) diagram at room temperature for the system A-B with phases α , β and γ . G_β^B and G_α^B Gibbs free energy of B with the respective structures β and α . b) The phases at room temperature in a state of stable equilibrium as a function of the different mole fractions x . The single-phase composition ranges are indicated. The two-phase range, whose location follows from the common tangent in (a), is shown by horizontal hatching. c) The same situation, but now for 'quenching' from high temperature. There are now only two single-phase composition ranges, whose separation is given by the intersection point of the curves for α and β in (a). It is assumed that the structure γ does not form, for example because it is too complicated. d) The situation if the structure γ actually did form.

be expected follows from the points where the curve for this phase intersects that for the face-centred cubic structure.

Fig. 4c shows the very broad amorphous range (Sa) found in this way, together with the results of experi-

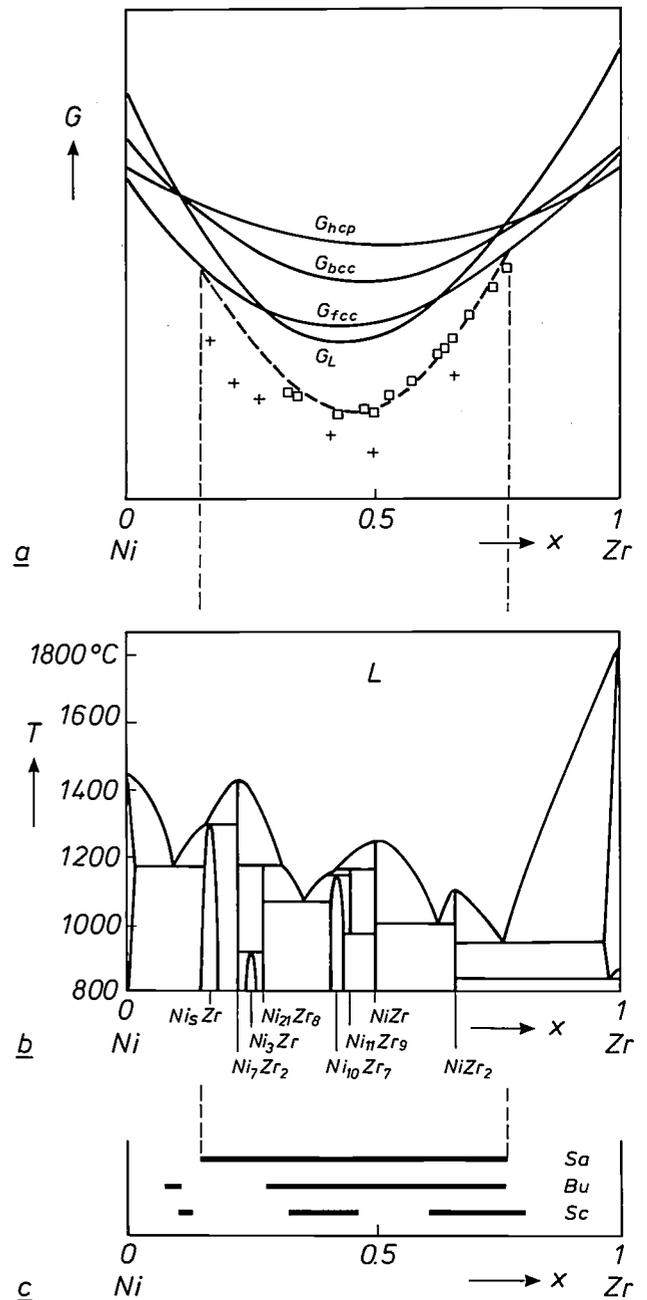


Fig. 4. a) The (G,x) diagram^[8] at 25 °C, and b) the phase diagram^[11] for the system Ni-Zr. hcp hexagonal close-packed structure. bcc body-centred cubic structure. fcc face-centred cubic structure. G_{hcp} , G_{bcc} , G_{fcc} and G_L Gibbs free energy of these structures and of the liquid phase. The points indicated by squares relate to amorphous material. The points marked by + are the minima of (G,x) curves not shown (see fig. 2a) for the various compounds indicated in (b). c) Composition ranges in which amorphous material can be found after quenching. Sa the (theoretical) range that is given by the intersection of the curve for amorphous material intersects the curve for the fcc structure. Bu the range found experimentally by Buschow^[12]. Sc the range found experimentally by Scott^[13].

ments. The ranges *Bu* were found at our Laboratories by K. H. J. Buschow; they are the result of quenching from the liquid phase on a rapidly rotating wheel [12]. The ranges *Sc*, found by M. G. Scott and presented in a review article by H. A. Davies, are also the result of quenching from the liquid phase, but with a different technique [13]. The agreement between the theoretically predicted range and the results of experiments is satisfactory. It should be remembered here that high accuracy cannot be expected in predictions of this kind, since the (*G*, *x*) curves are no more than approximations and because the cooling rate during quenching determines the phases that are actually formed. We shall return to this point later.

The method described can also be used for predicting metastable crystalline phases as well as amorphous phases [8]. We shall illustrate this with the copper-chromium, copper-molybdenum and copper-tungsten systems, which we have investigated. These metals have such poor miscibility in one another that two phases are formed even in the liquid state. We have also investigated the copper-cobalt system; these metals only have poor miscibility in the solid state [5]. On a substrate at room temperature we deposited thin films of these metal alloys from the vapour phase in different proportions and found broad — metastable — single-phase regions that were separated by fairly narrow two-phase regions, as shown in fig. 5. (If the substrate had been cooled well below room temperature, these two-phase regions would have been even narrower, giving the situation illustrated in fig. 3c.) In a state of stable equilibrium, i.e. not after quenching, these alloys always consist of two phases for all compositions, the two phases being the pure metals.

The location of the two-phase regions that form the limits of the single-phase regions should again follow from the intersection points of the (*G*, *x*) curves. Since these curves are not known in this case, we had to determine the location of the intersection point in a complicated and rather roundabout way. We plotted the energy of formation instead of the Gibbs free energy as a function of the mole fraction *x*. The energy of formation is the part of the Gibbs free energy that is necessary to make the metals alloy. The only contribution to the energy of formation we took into ac-

count was ΔG_{latt} , which is a measure of the stability of the metal in the crystal lattice [14]. The energy required for the physical mixing is practically independent of the crystal structure and was not therefore

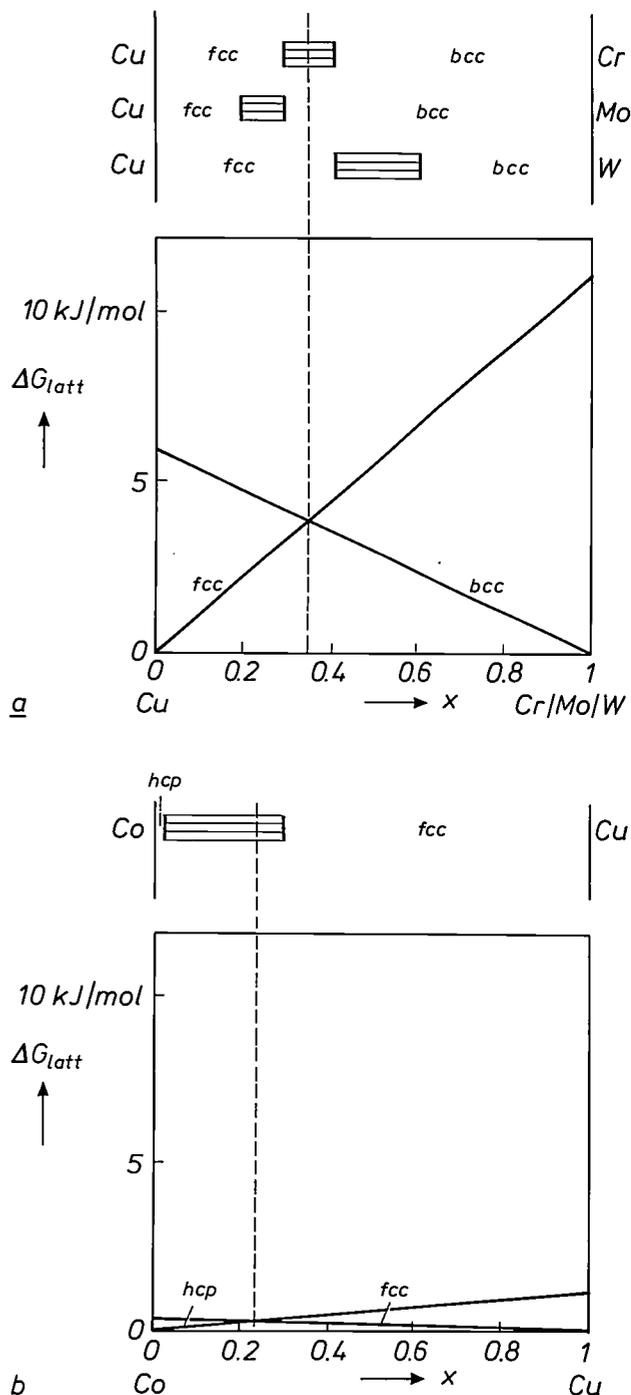


Fig. 5. a) The contribution ΔG_{latt} to the energy of formation as a function of the mole fraction *x* for the systems Cu-Cr, Cu-Mo and Cu-W. A linear relation is assumed between the points at *x*=0 and *x*=1. The horizontally hatched regions at the top of the figure represent the two-phase composition ranges we found experimentally. These ranges separate broad metastable single-phase ranges with the fcc and bcc structures. The single-phase ranges were found after quenching from the vapour phase. The location of the two-phase ranges corresponds reasonably well to the intersection point of the (ΔG_{latt} , *x*) lines. b) The same, for the system Co-Cu, in which cobalt has the hcp structure.

[11] P. Nash and C. S. Jayanth, The Ni-Zr (nickel-zirconium) system, Bull. Alloy Phase Diagrams 5, 144-148, 1984.
 [12] K. H. J. Buschow, Short-range order and thermal stability in amorphous alloys, J. Phys. F 14, 593-607, 1984.
 [13] H. A. Davies, Metallic glass formation, in: F. E. Luborsky (ed.), Amorphous metallic alloys, Butterworths, London 1983, pp. 8-25.
 [14] See: L. Kaufman and H. Bernstein, Computer calculation of phase diagrams, Academic Press, New York 1970. ΔG_{latt} is referred to here as 'lattice stability'; the values of ΔG_{latt} in fig. 5 for pure Cu, Cr, Mo, W and Co in the various crystal lattices are derived from this.

considered. Pure chromium, molybdenum and tungsten always have a body-centred cubic structure in their stable modification; see fig. 5. In these three metals, however, the fcc structure is unstable. Pure copper does have a stable fcc structure, but the bcc structure is unstable. In agreement with this, the ΔG_{latt} of the bcc structure in copper is 6 kJ/mol higher than for the fcc structure. In chromium, molybdenum and tungsten, on the other hand, the ΔG_{latt} of the bcc structure is lower than that of the fcc structure.

We have assumed that the ΔG_{latt} of the metals investigated is a linear function of the mole fraction. If we therefore draw lines connecting the points for the ΔG_{latt} of the pure metals, we find that the intersection points of the $(\Delta G_{\text{latt}}, x)$ lines for Cu with Cr, Mo or W, and for Co with Cu, agree reasonably well with the locations of the two-phase regions that we found.

Predicting amorphous phases from the phase diagram

The phase diagrams of many combinations of two metals, or of metals with metalloids, are known and available in handbooks. This is not so for diagrams in which the Gibbs free energy G is plotted as a function of the mole fraction x , at least not for all the crystalline and amorphous phases that can occur in a particular combination. This means that prediction of metastable phases from the (G, x) diagram usually requires calculation of the Gibbs free energy for different fractions and phases. The method is therefore fairly laborious.

Our method, in which direct use is made of the phase diagram, is much less laborious since it requires no calculations. Fig. 6 shows how the method is used for predicting an amorphous phase for the nickel-zirconium combination in fig. 4. Since at a high quenching rate there is little likelihood of compounds (usually with a complicated crystal structure) being formed, we ignore the central part of the phase diagram. This leaves us mainly with the two-phase regions of the virtually pure elements with liquid, and with the liquid phase L between them. The liquidus curves are then extrapolated. If the extrapolated curves do not intersect at a temperature above room temperature, they form the limits for a composition range in which a metastable equilibrium of the liquid or the amorphous phase can exist. We thus find a composition range for the amorphous phase for Ni-Zr at room temperature that corresponds reasonably well with the ranges in fig. 4c.

For the actual formation of an amorphous phase the rate of cooling from the liquid phase must be higher than a certain critical quenching rate. For a given composition the critical quenching rate increases with

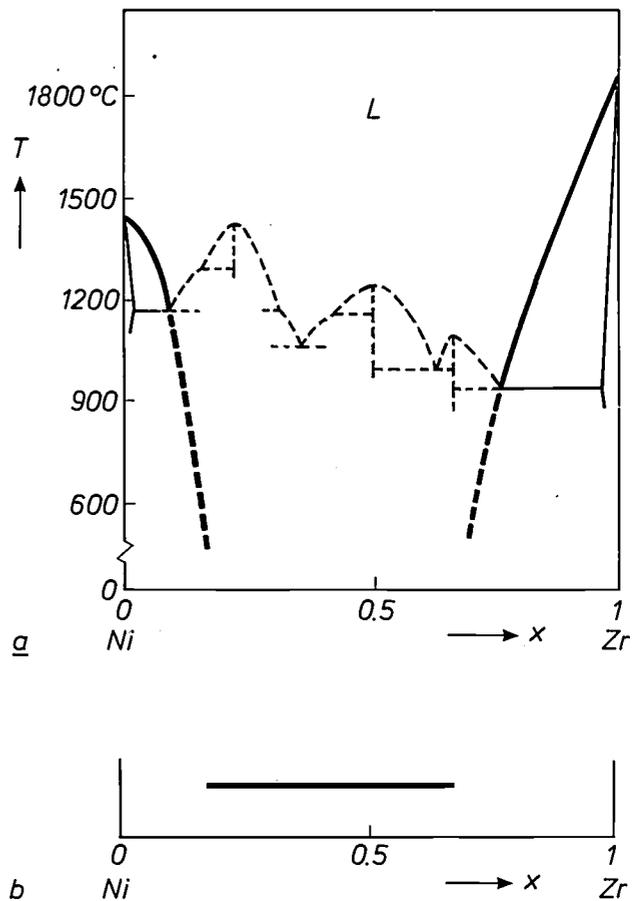


Fig. 6. The phase diagram can be used to determine the mole fractions at which an amorphous phase would be expected after quenching, in this case for the system Ni-Zr (see also fig. 4). *a*) The central part of the phase diagram is not considered because, in general, no compounds will be formed during quenching. The liquidus curves of the almost pure elements are extrapolated to room temperature. *b*) The composition range in which an amorphous phase can be expected; the range is bounded by the extrapolated curves.

the difference between the solidification temperature of the crystalline phase and the glass temperature of the amorphous phase. The glass temperature is the temperature at which the viscosity of the undercooled liquid produced by the quenching has increased to such an extent that it can be said to be an amorphous solid. In quenching from the liquid phase with a rotating wheel a quenching rate is reached that is 10^9 K/s at the most. In deposition from the vapour phase on to a substrate that is kept at room temperature, a quenching rate of about 10^{12} K/s can be calculated from energy considerations. The method used by Scott (see fig. 4c) was probably quenching of molten metal between rollers; the quenching rate that can be reached in this method is not so high, about 10^7 K/s. This explains why Buschow did find amorphous material at the fraction $x = 0.5$, where the solidification temperature of the compound NiZr is high, whereas Scott did not. At $x = 0.22$, where the compound Ni₇Zr₂ has an even higher solidification temperature, both Buschow

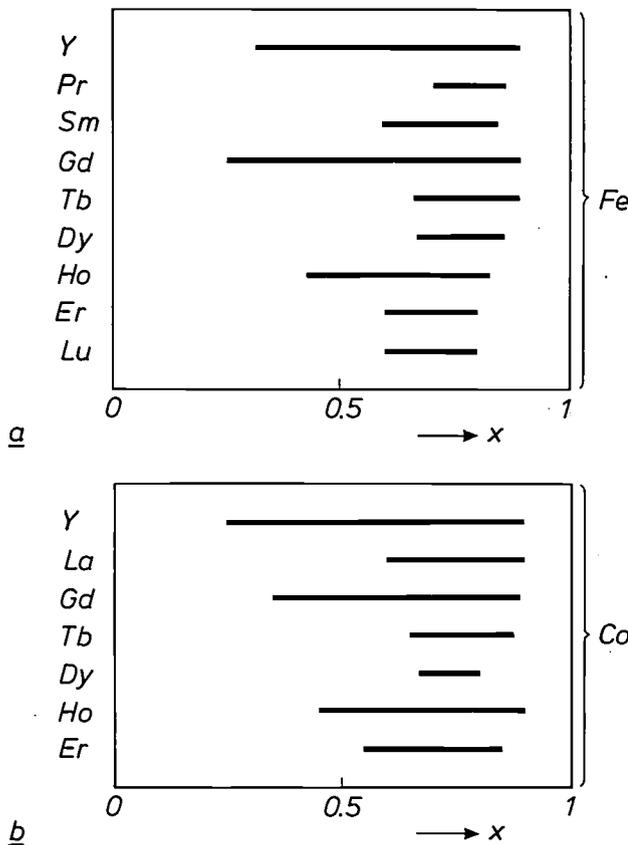


Fig. 7. Composition ranges in which amorphous material was found on quenching from the vapour phase *a*) for the combination of iron with yttrium and rare earths, *b*) for cobalt with yttrium and rare earths^[2].

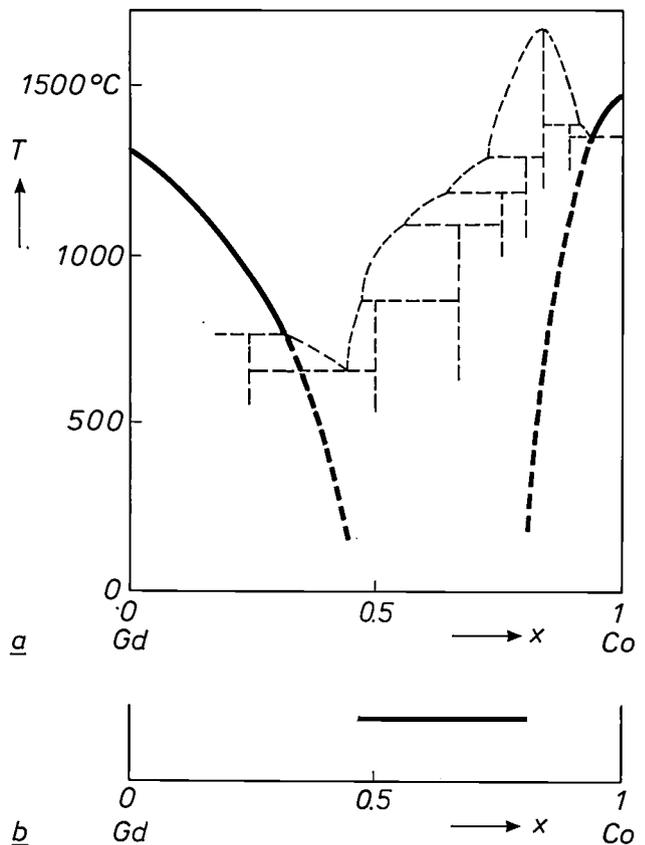


Fig. 8. Prediction of the composition range of an amorphous phase for the system gadolinium-cobalt; see also figs 6 and 7. *a*) Extrapolation of the liquidus curves in the phase diagram. *b*) The range thus found.

and Scott failed to find amorphous material. Presumably they would have done so if they had used the method of deposition from the vapour phase.

In research aimed at finding suitable ferrimagnetic materials for magneto-optical recording, many combinations of elements have been investigated at our Laboratories for the occurrence of amorphous phases^[2]. The results of this work by J. W. M. Biesterbos and A. G. Dirks are presented in *fig. 7a* and *b*. They relate to combinations of iron and cobalt with yttrium and rare earths. The amorphous films were produced from the vapour phase. The advantage of this kind of amorphous material for magneto-optical recording is that the magnetic-compensation temperature can be accurately adjusted by varying the composition in a wide range. If we started from crystalline material of the combinations mentioned above, which can form compounds, we would be confined to fairly narrow composition limits corresponding to the phase diagram.

The usefulness of our method for predicting an amorphous phase is illustrated for gadolinium and cobalt in *fig. 8*. The result agrees reasonably well with the experimentally determined range in *fig. 7b*. The

other ranges in *fig. 7* can also be predicted with phase diagrams, though not so accurately.

Amorphous (metastable) phases are known for most combinations of zirconium with 3d transition metals, e.g. nickel and zirconium; see *fig. 4*. It did not prove possible, however, to produce amorphous chromium-zirconium. *Fig. 9* shows that this can also be explained from the phase diagram. In Cr-Zr, unlike Ni-Zr (see *fig. 6*) there is no question at all of a range with a metastable liquid phase that extends to low temperatures: the extrapolated liquidus curves intersect each other at a point at about 1000 °C — far above room temperature.

The transition from stable thermodynamic equilibrium to a metastable equilibrium

When the substrate is not kept at room temperature during deposition from the vapour phase, but at a higher temperature, the term quenching is not so appropriate. In this situation transitional states can occur between the stable thermodynamic equilibrium and a metastable equilibrium. The effect of the substrate temperature *T_s* has been investigated by

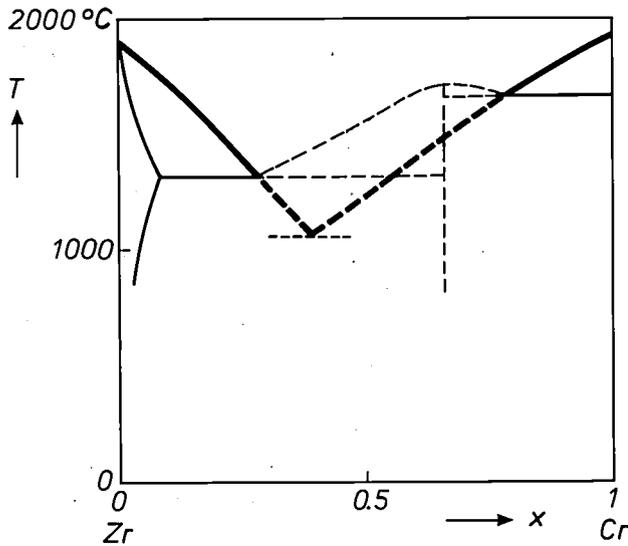


Fig. 9. The extrapolation of the liquidus curves in the phase diagram of the system zirconium-chromium. It can be seen from this that it is almost impossible to obtain amorphous Zr-Cr by quenching, since the intersection point of the extrapolated curves is at about 1000 °C.

H. T. G. Hentzell *et al.* for the system aluminium-nickel in the range $0.4 < x_{Ni} < 1$ [15]. The complicated phase diagram for this system can be seen in fig. 10a. The phases found by Hentzell *et al.* as a function of the substrate temperature are shown in fig. 10b. The intersection point of the curves at 300 K is the result of our own work, involving vapour deposition on an unheated substrate for varying nickel content and investigation of the structures by X-ray diffraction and transmission electron microscopy.

At $T_s > 700$ K the results of Hentzell *et al.* agree reasonably well with the phase diagram. At $300 \text{ K} < T_s < 700$ K the mobility of the atoms is on the one hand too low for a stable equilibrium to be reached, but on the other hand too high for a metastable equilibrium to be 'frozen in'. In this transitional region the phase relations are determined by a combination of atomic mobility and the tendency of the system to seek a minimum Gibbs free energy. This can lead to unexpected situations, because the three-phase region $\beta' + \alpha' + fcc$ in fig. 10b could not exist in a binary system in thermodynamic equilibrium.

At $T_s < 300$ K the α' phase is not present, possibly because the mobility of the atoms is insufficient for this phase to occur. There are then only two broad single-phase regions: that of the fcc structure of pure nickel and that of the β' structure, which corresponds to the (cubic) structure of caesium chloride. The mole fraction that corresponds to the boundary of the two single-phase regions is in agreement with the location of the intersection point of the $(\Delta G, x)$ curves for the fcc and β' structures; see fig. 10c. We do not need to

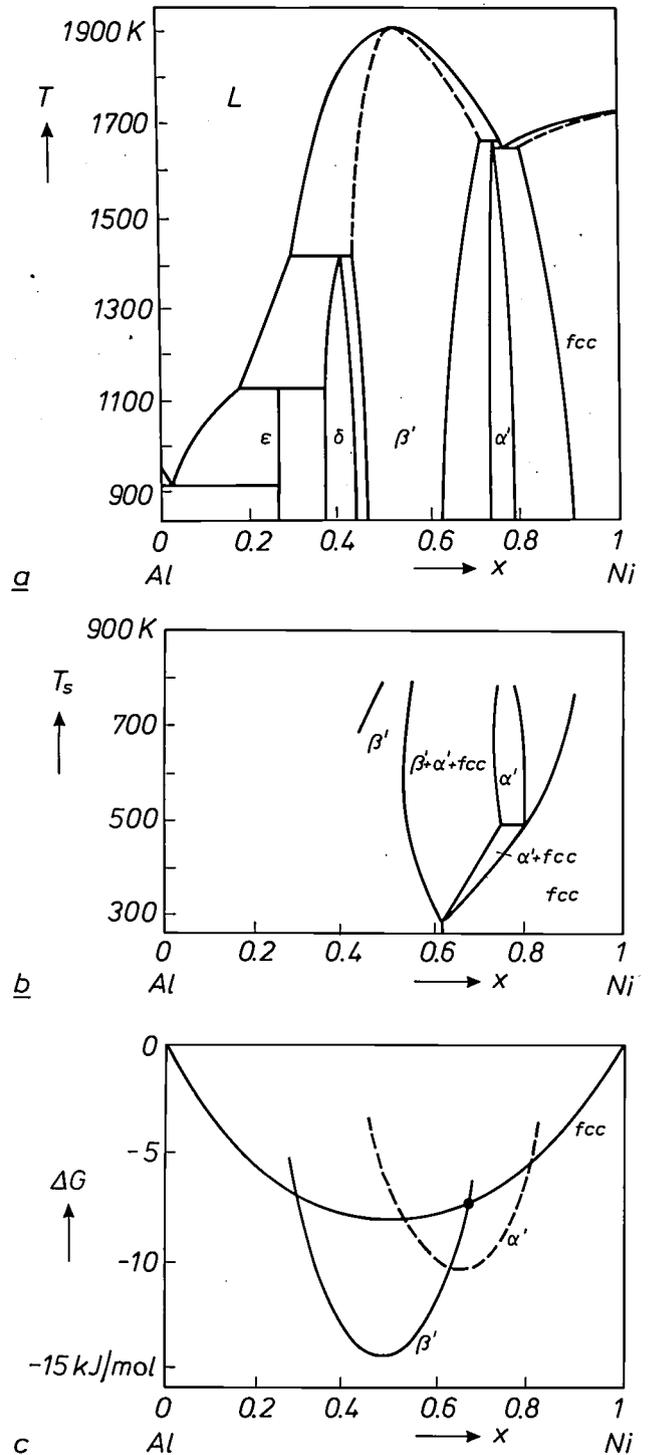
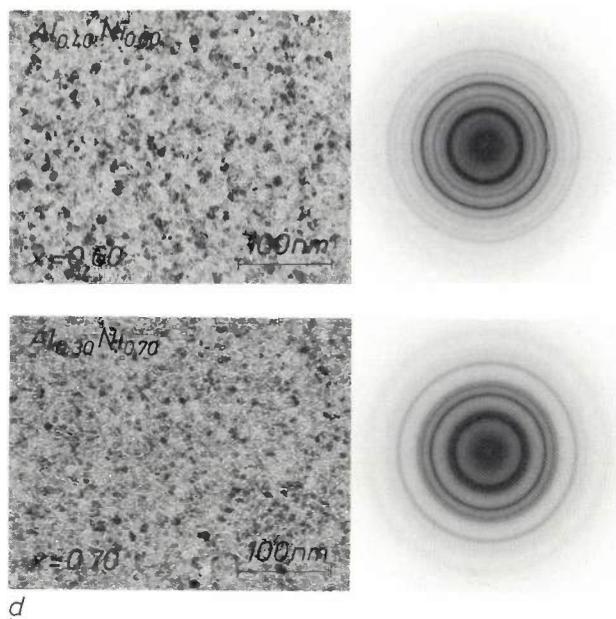


Fig. 10. The effect of the substrate temperature T_s in quenching from the vapour phase for the system aluminium-nickel. a) The phase diagram. ϵ , δ , β' , α' and fcc crystal structures. b) The phase relationships as a function of the substrate temperature [15]. (The intersection point of the lines at room temperature is the result of our own research.) If T_s is equal to room temperature, there are two single-phase regions. c) The energy of formation ΔG as a function of x for the structures β' , α' and fcc at room temperature. The structure α' is not formed at room temperature because it is too complicated; the corresponding curve is therefore shown dashed. The intersection point shown corresponds to the boundary of two single-phase regions on quenching. This agrees well with the results in (b). d) Left: bright-field micrographs, right: electron-diffraction patterns, made with a Philips transmission electron microscope at two different mole fractions x of the $Al_{1-x}Ni_x$ film deposited from the vapour phase. The diffraction patterns are clearly different.



consider the curve for the phase α' , since this phase, as noted, is not formed at room temperature.

The boundary of the single-phase regions also follows from the micrographs in fig. 10d of structures with $x=0.60$ and 0.70 ; these micrographs were made with a Philips transmission electron microscope. The bright-field pictures show an almost identical grain structure, consisting of extremely fine crystallites. The electron-diffraction patterns reveal that there has been a definite change in the crystal structure.

^[16] H. T. G. Hentzell, B. Andersson and S.-E. Karlsson, Grain size and growth of Ni-rich Ni-Al alloy films, *Acta Metall.* **31**, 2103-2111, 1983.

The results in fig. 10 demonstrate that, when the substrate temperature changes during the deposition of thin metal-alloy films from the vapour phase, there are three separate temperature regions:

- a region at high substrate temperature, where the location of the phases corresponds reasonably well to the phase diagram;
- a region at low substrate temperature, where the boundary of wide and largely metastable single-phase regions corresponds to the intersection point of the relevant $(\Delta G, x)$ curves, and
- a region at intermediate temperatures with fairly complicated phase relationships, since on the one hand the system tends towards a stable thermodynamic equilibrium and on the other hand the mobility of the atoms is limited.

Summary. There is a growing interest in amorphous phases, for one reason because of the application of thin films of amorphous metal in magneto-optical recording. The phase diagram, which applies for a state of stable thermodynamic equilibrium, can often be used to predict whether an amorphous phase — invariably metastable — will occur in a binary system. This is so, for example, in combinations of iron or cobalt with rare earths, from which amorphous films have been deposited from the vapour phase to see whether they can be used as a medium for magneto-optical recording. The occurrence of metastable phases — either amorphous or crystalline — can also be predicted from the (G, x) diagram, in which no account need be taken of the (G, x) curves of phases that cannot occur owing to limited atomic mobility. The phase that does form on quenching is the one that has the lowest Gibbs free energy G at a given composition. If the substrate is kept at a higher temperature than room temperature during deposition from the vapour phase, fairly complicated phase relationships arise in a transitional region between a state of stable thermodynamic equilibrium at high temperature and a region of metastable equilibrium at room temperature.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

	Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany		A	
	Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium		B	
	Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands		E	
	Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany		H	
	Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France		L	
	Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.		N	
	Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England		R	
	Philips Research Laboratories, Sunnyvale, P.O. Box 9052, Sunnyvale, CA 94086, U.S.A.		S	
B. Post & J. Ladell	N	The phases of forbidden reflections	Acta Crystallogr. A	43 173-179 1987
O. Boser	N	Electromechanical resonances in ceramic capacitors and evaluation of the piezoelectric materials' properties	Adv. Ceram. Mater.	2 167-172 1987
P. van der Putte, D. K. Sadana, E. K. Broadbent & A. E. Morgan	S	Growth of selective tungsten films on self-aligned CoSi ₂ by low pressure chemical vapor deposition	Appl. Phys. Lett.	49 1723-1725 1986
A. E. Morgan & P. Maillot	S	Depth profiling of As at the SiO ₂ /Si interface using secondary ion mass spectrometry	Appl. Phys. Lett.	50 959-961 1987
J. Petruzzello, D. Olego, S. K. Ghandhi*, N. R. Taskar* & I. Bhat* (*Rensselaer Polytech. Inst., Troy, NY)	N	Transmission electron microscopy of (001) CdTe on (001) GaAs grown by metalorganic chemical vapor deposition	Appl. Phys. Lett.	50 1423-1425 1987
K. Mohammed, D. J. Olego, P. Newbury, D. A. Cammack, R. Dalby & H. Cornelissen	N	Quantum confinement and strain effects in ZnSe-Zn _x Se _{1-x} strained-layer superlattices	Appl. Phys. Lett.	50 1820-1822 1987
A. J. Fox	R	Thermal design for germanium acoustooptic modulators	Appl. Opt.	26 872-884 1987
P. J. Kelly, R. Car (<i>Int. School Adv. Studies, Trieste</i>) & S. T. Pantelides (<i>IBM T. J. Watson Res. Center, Yorktown Heights, NY</i>)	E	Theoretical determination of the vacancy migration energy in silicon	Defects in Semiconductors, H. J. von Bardeleben (ed.), Mater. Sci. Forum, Vol. 10-12, Trans Tech, Aedermannsdorf	115-120 1986
P. Boivin*, J. Rabier*, H. Garem* (*Lab. Métallurgie Phys., C.N.R.S., Poitiers) & M. Duseaux	L	Dislocation substructures and plasticity of GaAs below 400 °C as a function of doping	Defects in Semiconductors, H. J. von Bardeleben (ed.), Mater. Sci. Forum, Vol. 10-12, Trans Tech, Aedermannsdorf	781-786 1986
S. Makram-Ebeid & P. Boher	L	Electron hopping between bombardment induced defects in gallium arsenide	Defects in Semiconductors, H. J. von Bardeleben (ed.), Mater. Sci. Forum, Vol. 10-12, Trans Tech, Aedermannsdorf	1075-1080 1986
E. Dupont-Nivet, J. N. Patillon, J. P. André & G. M. Martin	L	Study of deep levels in Al _y Ga _x In _{1-x} P material grown by MOVPE	Defects in Semiconductors, H. J. von Bardeleben (ed.), Mater. Sci. Forum, Vol. 10-12, Trans Tech, Aedermannsdorf	1207-1212 1986

J. K. Chamberlain*, F. M. Clayton* (*First Res. Centre, Wembley), H. Sari, P. Vandamme (CNET, Lannion)	L	Receiver techniques for microwave digital radio	IEEE Commun. Mag. 24 (No. 11)	43-54	1986
J. R. Brandsma, A. A. M. L. Bruekers & J. L. W. Kessels	E	PHILAN: a fiber-optic ring for voice and data	IEEE Commun. Mag. 24 (No. 12)	16-22	1986
J. Khurgin	N	Theoretical and experimental investigation of amplified spontaneous emission in electron-beam-pumped semiconductor lasers	IEEE J. QE-23	194-204	1987
G. R. Gao	N	Maximum pipelining linear recurrence on static data flow computers	Int. J. Parallel Program. 15	127-149	1986
A. J. M. Houtsma (Inst. Perception Res., Eindhoven) & T. D. Rossing (Northern Illinois Univ., DeKalb, IL)		Effects of signal envelope on the pitch of short complex tones	J. Acoust. Soc. Am. 81	439-444	1987
T. Baller*, D. J. Oostra*, A. E. de Vries* (*FOM, Amsterdam) & G. N. A. van Veen	E	Laser-induced etching of Si with chlorine	J. Appl. Phys. 60	2321-2326	1986
A. H. van Ommen	E	Diffusion of ion-implanted Sb in SiO ₂	J. Appl. Phys. 61	993-997	1987
J. J. Harris, C. T. Foxon, K. W. J. Barnham, D. E. Lacklison, J. Hewett & C. White	R	Two-dimensional electron gas structures with mobilities in excess of $3 \times 10^6 \text{cm}^2 \text{V}^{-1} \text{s}^{-1}$	J. Appl. Phys. 61	1219-1221	1987
J. M. Towner & B. Coulman	S	Effects of vanadium and chromium on aluminum electromigration	J. Appl. Phys. 61	1392-1403	1987
N. D. Young, J. B. Clegg & E. A. Maydell-Ondrusz (Univ. Surrey, Guildford)	R	Low-temperature annealing of shallow arsenic-implanted layers	J. Appl. Phys. 61	2189-2194	1987
R. Bar-Gadda	S	A theoretical prediction of film thickness profiles in a mixed convection-diffusion regime for the chemical vapor deposition of polysilicon in annular tubes	J. Electrochem. Soc. 133	2123-2131	1986
A. E. Morgan, E. K. Broadbent, M. Delfino, B. Coulman & D. K. Sadana	S	Characterization of a self-aligned cobalt silicide process	J. Electrochem. Soc. 134	925-935	1987
H. H. R. Jansen (Philips Lighting Div., Eindhoven)		Aspects of the recrystallization kinetics of doped tungsten	Philips J. Res. 42	3-14	1987
G. Dittmer & U. Niemann	A	Evaluation of thermodynamic data on zirconium and hafnium halides and oxyhalides by means of transport experiments	Philips J. Res. 42	15-40	1987
G. Dittmer & U. Niemann	A	The influence of solid phases on transport cycles in halogen-filled incandescent lamps	Philips J. Res. 42	41-57	1987
E. Fischer	A	Modelling of low-power high-pressure discharge lamps	Philips J. Res. 42	58-86	1987
J. J. de Groot*, J. Schlejen**, J. P. Woerdman** (**Univ. Leiden) & M. F. M. de Kieviet* (*Philips Lighting Div., Eindhoven)		The influence of NaNa, NaHg and NaXe molecules on the spectrum of the high-pressure sodium lamp	Philips J. Res. 42	87-101	1987
W. van Erk* & T. Rietveld* (*Philips Lighting Div., Eindhoven)		The reaction of rare-earth metal iodides with the silica wall in metal halide discharge lamps	Philips J. Res. 42	102-118	1987
G. de With	E	Preparation, microstructure and properties of Y ₃ Al ₅ O ₁₂ ceramics	Philips J. Res. 42	119-130	1987
W. J. A. Goossens	E	Excess free energy due to hard-core repulsions for the cholesteric and smectic phases	Phys. Rev. A 35	1843-1846	1987
G. Duggan & H. I. Ralph	R	Exciton binding energy in type-II GaAs-(Al,Ga)As quantum-well heterostructures	Phys. Rev. B 35	4152-4154	1987
L. Knox, P. Patt & R. Maresca	N	Design of a flight qualified long-life cryocooler	Proc. 3rd Cryocooler Conf. on Refrigeration for cryogenic sensors & electronic systems, Boulder, CO, 1984	1-20	1984

P. van der Putte	S	The reaction kinetics of the H ₂ reduction of WF ₆ in the chemical vapor deposition of tungsten films	Proc. Workshop Tungsten and other Refractory Metals for VLSI Applications, Palo Alto, CA, 1986	77-84	1987
P. van der Putte, D. K. Sadana, E. K. Broadbent & A. E. Morgan	S	Growth of selective tungsten on self-aligned CoSi ₂ by low pressure chemical vapor deposition	Proc. Workshop Tungsten and other Refractory Metals for VLSI Applications, Palo Alto, CA, 1986	101-107	1987
E. K. Broadbent & J. M. Towner	S	Electromigration-induced short circuit failure in aluminum/tungsten (CVD) conductors	Proc. Workshop Tungsten and other Refractory Metals for VLSI Applications, Palo Alto, CA, 1986	247-255	1987
L. Gutai, M. Delfino & J. M. DeBlasi	S	Thermal stability of Al/W/p-n Si contacts during post-metallization annealing	Proc. Workshop Tungsten and other Refractory Metals for VLSI Applications, Palo Alto, CA, 1986	265-272	1987
D. R. Gentner	N	Timing of skilled motor performance: Tests of the proportional duration model	Psychol. Rev. 94	255-276	1987
J. Reimer, A. Bhattacharyya & K. Ritz	S	Electron beam induced current analysis of voltage breakdown sites in thin MOS oxides	Scanning Microsc. 1	23-30	1987
V. Lacroix	B	Pixel labeling in a second-order Markov mesh	Signal Process. 12	59-82	1987
F. J. A. M. Greidanus & J. T. C. van Roosmalen	E	ODMR experiments on the 1.65eV luminescence in GaP:Cu	Solid State Commun. 61	653-657	1987
H. C. de Graaff	E	Polycrystalline silicon in integrated circuits	Solid State Sci. 57	170-184	1984
S. van Lerberghe	E	Digital adaptive multipath equalization of FM signals	Thesis, Brussels	1-226	1987
A. H. M. van Roermund	E	Noise and accuracy in switched capacitor modulation circuits	Thesis, Heverlee	1-250	1987

Contents of Philips Telecommunication and Data Systems Review 45, No. 2, 1987

- J. R. Gray: MAESTRO: science in the art of software production (pp. 1-9)
 J. F. Briend & J. J. Plancke: French PTT MINITEL and LECAM programme (pp. 10-26)
 B. Jurgens: MPX operating system (pp. 27-44)
 L. T. van der Bijl & P. H. A. Timmers: The application environment of the P9X00 (pp. 45-60)

Contents of Electronic Components & Applications 8, No. 2, 1987

- R. Croes & A. de Pagter: ACL . . . Advanced CMOS Logic that lengthens the stride of low-power systems (pp. 66-75)
 S. Baliga, G. Goodhue & J. Jenkins: Microcontroller eases I/O processing burden (pp. 76-80)
 D. Haslam & M. J. Hill: Breakover diodes for the protection of telephone equipment (pp. 81-86)
 F. A. Pieters: Quality — dry reed switches (pp. 87-94)
 G. Conn: The SCC68070 — a monolithic 68000 CPU and peripherals (pp. 95-100)
 T. v.d. Wouw: Cascode-driven SMPS with high-voltage darlington transistors (pp. 101-108)
 N. Siddique & F. Krupecki: Fast controller converts large static RAMs to FIFO buffers (pp. 109-112)
 S. Baliga: Managing graphics displays using a microcontroller (pp. 113-117)
 D. Eckstein & R. W. Stamer: Improved varicaps for tv tuners (pp. 118-124)

CARIN, a car information and navigation system

M. L. G. Thoone

Whatever we think of the car, it is now a part of modern life. At the very least it is useful for getting from A to B. In fact, however, this is not always so easy. We have to be able to read the road map and we need to know where we are. Planning the most economical route can also present problems, particularly if we are unexpectedly held up during the journey. A team at the Philips Project Centre at Geldrop, in the Netherlands, are just rounding off their work on prototypes of an electronic vehicle-guidance system called CARIN (an acronym for Car Information and Navigation System). During the journey CARIN speaks through a speech synthesizer to tell the driver which road to take at crossings and road junctions. Before the journey the system also plans the best route, finding a way round traffic delays. In the design of CARIN the vast storage capacity of the CD-I (Compact Disc Interactive) was of inestimable value for recording the digitized map data.

Introduction

In a car it is the passenger sitting next to the driver who usually has to navigate: plan the route, read the map, relate actual position to the map, and give directions to the driver. The passenger performs these tasks with varying degrees of success, and we all know that any inadequacies of performance can lead to disharmony. Perhaps more serious than this is that mistakes in planning and following a route can mean that greater distances are travelled than is necessary. For professional road users, such as police, fire brigades, delivery and freight services, ambulances and taxis,

economic route planning and good vehicle guidance can give substantial savings. This is certainly so if the guidance helps drivers to avoid traffic delays by an early change of route. All this can become reality if the 'navigator's' tasks are taken over by an electronic computer-controlled system. Electronic navigation and information systems are well established in marine and aviation applications. Car and electronics manufacturers are now trying to develop systems that are suitable for automotive applications.

Almost all of the car navigation and information systems now known to be operational have the disadvantage that they do not have a suitable medium for storing topographical information. Some of the systems also require substantial investments in 'infrastructure', such as beacons for radio range finding.

Ir M. L. G. Thoone, formerly with Philips Research Laboratories, Eindhoven, is now with the Philips Consumer Electronics Division. He is the leader of a combined team (from the Research Laboratories and the Consumer Electronics Division) working on the CARIN project. Various members of this team made valuable contributions to this article.

Nor are most of the systems very user-friendly, since the driver has to enter the starting position and desired destination as a map grid reference [1].

At the Geldrop Project Centre, which is part of Philips Research Laboratories, a car navigation and information system that does not have these disadvantages is now being developed in cooperation with the Philips Consumer Electronics Division [2]. The system, called CARIN (Car Information and Navigation System), plans the best route, guides the driver with the aid of a speech synthesizer, periodically determines the position of the vehicle, selects an alternative route if there are coded digital radio signals reporting traffic obstructions, and can also give tourist information; see *fig. 1*. The special feature of CARIN is that it uses a highly efficient medium for storing digital data, the Compact Disc. The version used for CARIN is the CD-I (Compact Disc Interactive), specially developed for storing audio, video and computer data for consumer applications. This new medium is also ideally

less than that of a magnetic-tape cassette; about a second as compared with about a minute. (The access time is longer, however, than that of a semiconductor memory; the problems connected with the non-negligible access time of the Compact Disc will be dealt with later in this article.)

The units that are carried in the car are:

- an on-board computer for controlling and processing the information, containing a semiconductor main memory with a capacity of 1 Mbyte (8 Mbit);
- a Compact Disc player with amplifier and loudspeakers;
- sensors for determining the position of the vehicle (if wheel sensors are fitted for an antilock braking system (ABS), it may be possible to use them for some aspects of position determination);
- an electronic speech synthesizer;
- a simple keyboard; and
- a small flat liquid-crystal display (LCD); for showing a simplified map of the next intersection.



Fig. 1. The interior of a car fitted with the CARIN vehicular guidance system. The car was designed by the British firm I.A.D. (Industrial Automotive Design), of Worthing, West Sussex. The prototype is the result of ideas for 'the car of the not too distant future'.

suited for the storage of topographical information. The amount of digital data that can be stored on one Compact Disc is vast: 4800 Mbit. The storage capacity of the largest semiconductor memory of the dynamic-RAM type (Random-Access Memory) as yet available is far less: 1 Mbit. The storage capacity of one Compact Disc corresponds to 150 000 typed A4 pages, or — and this is important for navigation systems — to the detailed topographical information for a medium-sized country. The access time is also much

These units complete the basic CARIN system. Other units can also be added:

- a touch-screen graphic display for interactive consultations of map data for a larger area;
- a special car radio that receives messages about traffic delays, road hazards, etc. and inputs them to the system; and
- an interface with sensors present in the car for monitoring the status of oil temperature and pressure, vehicle lighting, and fuel.

The driver can only consult the map data on the touch screen when the car is stationary. The keyboard is used for entering the starting position and destination, and any special requirements. The driver does not have to enter a map reference, but can simply specify the names of countries, towns and streets, since lists of such names are stored with their map references on the Compact Disc. The computer calls the required application programs from the disc, calculates an efficient route from the starting position to the destination, and ensures that the information required for following this route is read from the disc and stored in the main memory. Then the player can be used for playing music again. During the journey the driver is guided by spoken directions from the speech synthesizer, such as 'take first turning left', 'fork right', 'take first exit at roundabout', and 'we have arrived'. The information on the small flat display gives extra support. In CARIN the emphasis is on auditive guidance for the driver, for reasons of traffic safety.

The principle of the vehicle-locating procedure is as follows. Sensors that measure the number of revolutions of the non-driven wheels provide information about the distance travelled and changes in the direction of the vehicle. An electronic compass, which measures the direction of the vehicle in relation to the Earth's magnetic field, gives information about the heading. The computer uses both kinds of information to periodically calculate the location and size of an area in which the car is located with a certain probability. This area is compared with the topographical information in the main memory. Since the car should be on a road whose data is stored in the memory, this comparison determines the part of the road where the car is located. This procedure is continuously repeated. The positioning error with this method is about 20 m. In later generations of CARIN we also intend to make use of a satellite navigation system called GPS (Global Positioning System), due to become operational in the nineties [3].

The European Broadcasting Union has recently reached agreement on RDS (Radio Data System), a system that can be used to transmit digital data at a rate of 1200 bits/s in a free frequency band within the FM broadcast signal. RDS will be introduced in many West European countries within the not too distant future [4]. In future versions of CARIN, if the system receives an RDS message about traffic delays or road obstructions, it will calculate the time to reach the affected position. If the received message says that the obstruction will have disappeared within that time, the route will remain unchanged. If the obstruction will still be there, CARIN will then calculate a route

that bypasses the obstacle and gets the car to its destination in the shortest possible time.

In the rest of this article we shall first consider the method of digitizing map data. Then we shall look at the navigation problems, and see how vehicle location can be improved with the aid of the stored map data. Next we shall consider some aspects of the storage of digital information on a Compact Disc, with special attention to the partitioning and optimum arrangement of the topographical information. Finally, some aspects of the data management will be considered.

Digitizing the map data

In modern cartography some of the topographical information is now processed digitally. Two methods are used, the raster scanning method and the vector method. In the raster scanning method a map is divided into picture elements (pixels), say 0.1 mm by 0.1 mm. The colour of each pixel is represented by a digital code. To digitize the map of Amsterdam (covering 14 km by 12 km), for example, on a scale of 1/15 000, a storage capacity of 375 Mbit is required for the raster scanning method [2]. In the vector method, as we use it, the axes of the roads are approximated by straight-line segments, which each represent a 'vector'. This method requires far less storage capacity. Digitized in this way the map of Amsterdam occupies 1 Mbit, based on the assumptions made in the first article of note [2]. If other information is added to the purely topographical data, such as the names of streets and general information about filling stations and restaurants, etc., this will require a further 1 Mbit of memory. As we noted earlier, the Compact Disc has a capacity of 4800 Mbit, so that more than 1000 maps like the map of Amsterdam can be stored on one disc.

Fig. 2 illustrates the digitization procedure we use in our vector method. The road network is translated

- [1] W. Zimdahl, Guidelines and some developments for a new modular driver information system, Proc. IEEE Vehicular Technol. Conf., Pittsburgh, PA, 1984, pp. 178-182; H. Friedl, Guidance of vehicle traffic flows — ALI, Radio Elektron. Schau 51, 266-269, 1975; D. J. Jeffery, Options for the provision of improved driver information systems: the role of micro-electronics and information technology, Proc. IEE Colloquium on Vehicular Route Guidance, London 1985, pp. 1-3.
- [2] M. L. G. Thoone and R. M. A. M. Breukers, Application of the Compact Disc in car information and navigation systems, SAE Int. Cong. and Exposition, Detroit, Mich., 1984, Tech. Paper 840156; M. L. G. Thoone, L. M. H. E. Driessen, C. A. C. M. Hermus and K. van der Valk, The car information and navigation system CARIN and the use of Compact Disc Interactive (CD-I), SAE Int. Cong. and Exposition, Detroit, Mich., 1987, Tech. Paper 870139;
- [3] T. A. Stansell, Jr, Civil GPS from a future perspective, Proc. IEEE 71, 1187-1192, 1983.
- [4] Main characteristics of a unified European VHF Radio Data System, Doc. Gt R1 290, Eur. Broadcast. Union (EBU), 1982.

into a graphic structure of points connected by straight-line segments. The position of each point is given as the coordinates of a map reference, i.e. as (x,y)-coordinates in a rectangular system. (In a next-generation system these could be geodetic coordinates.) There are two kinds of points, called nodes and intermediate points. A node can be:

- a point where at least three segments meet (1);
- an intersection with the edge of the map or with some other artificial boundary (2);
- an intersection with an administrative boundary (3);
- the end of a cul-de-sac (4); or
- a point where one or more items of data for a street change, for example its name or the type of road surface (5).

All other points are called intermediate points. They are found at a bend in a road or they are used to approximate the curvature of a road by a number of segments, like the points 11 and 12 in fig. 2. The succession of points in a sequence that starts and ends at a node is called a 'chain'. In the figure, 1, 12, 11 and 5 form a chain, and so do 1 and 4, 1 and 3, and 2 and 5.

The graphic structure formed by the various chains can be regarded as the 'skeleton' of the digitized map. 'Attributes' can then be added to the skeleton, such as:

- street names;
- classes of road: motorways, trunk roads, secondary roads and limited-access roads;
- directions of one-way streets.

At the start of a journey (and perhaps during the journey as well) some of the map data on the Compact Disc has to be read and stored in the main memory. It is therefore important to partition the data and

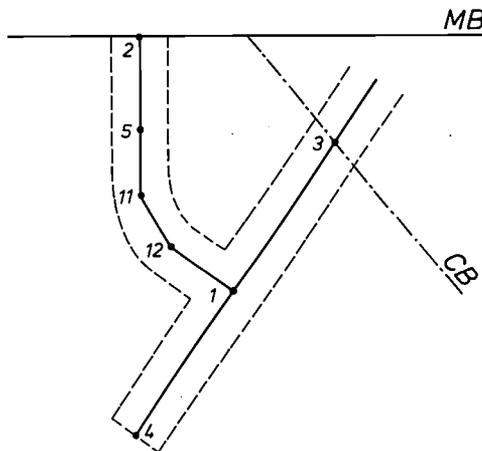


Fig. 2. Vector method for digitizing information from the map. The axes of roads are approximated by line segments that connect points whose position is fixed by the coordinates of map references. 1 to 5 nodes. At node 5 the condition of the road surface changes. 11, 12 intermediate points. MB map boundary. CB local-government boundary. The segment 1-4 is a cul-de-sac. Points 1, 12, 11 and 5 form a 'chain'.

arrange it on the Compact Disc in such way as to give short access times. The CARIN project group have made extensive studies of the problem of the optimum method of data storage, since agreement on proper standards here is of considerable importance to Philips. The theories of the management of topographical information use the topological concepts of the 0-cell, the 1-cell and the 2-cell [6]. The numbers 0, 1 and 2 represent the 'dimension' of a cell. In our case, 0-cells correspond to nodes, 1-cells to chains and 2-cells to areas within joined-up (concatenated) chains; see fig. 3. We shall deal with this in more detail at the end of the article.

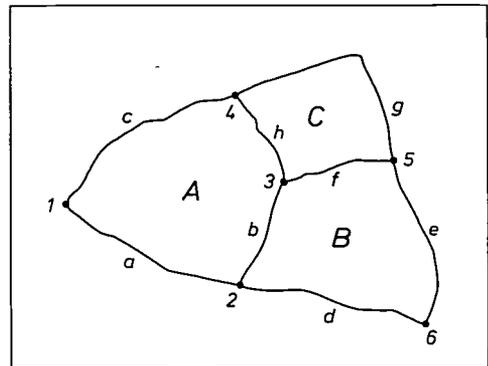


Fig. 3. Some concepts from topology [6]. 1 to 6 0-cells (nodes). a to h 1-cells (chains). A to C 2-cells (areas inside concatenated chains).

Vehicle location

General equations

The electronic navigation system of CARIN has to obtain a 'fix' for the vehicle at regular intervals, say every five metres, by dead reckoning. This amounts to determining the centre of the Vehicle Location Probability Area (VLPA), the area in which the vehicle is most likely to be travelling, with a probability of, let us say, 95%. The coordinates of this centre are x_n and y_n for the n -th dead-reckoned fix. (The y -axis points towards geographic north.) The navigation system takes the coordinates x_n and y_n and calculates new coordinates x_{n+1} and y_{n+1} from the equations

$$x_{n+1} = x_n + \delta_n \sin \frac{1}{2} (\phi_{n+1} + \phi_n) \tag{1a}$$

and

$$y_{n+1} = y_n + \delta_n \cos \frac{1}{2} (\phi_{n+1} + \phi_n), \tag{1b}$$

where δ_n represents the distance travelled between the n th and the $(n + 1)$ th fixes, and ϕ_n and ϕ_{n+1} represent the headings of the vehicle with respect to geographic north at the n th and the $(n + 1)$ th fixes.

The periodic dead-reckoning of a fix requires periodic measurement of the heading and the distance travelled. As noted earlier, the car has sensors for this: an electronic compass and two wheel sensors or 'odometers'. These measure the number of revolutions of non-driven wheels about their axes. We shall now take a closer look at these sensors and at the processing of the signals they produce.

The sensors

Both the electronic compass and the wheel sensors contain transducers that convert a magnetic flux into a voltage. The transducers are made by Philips Valvo, and their type number is KMZ10. They are based on the magnetoresistance effect, an effect in which the resistance of certain materials depends on the angle between the direction of the current and that of the magnetization [6]. The effect is strongest in ferromagnetic material. The transducer is therefore made from a large number of strips of a nickel-iron alloy. The strips are connected in four groups to form a Wheatstone bridge; see fig. 4a.

When the external magnetic field is zero, the preferred direction of magnetization of the strips in fig. 4a will be along the longitudinal axis. This corresponds to the x -axis. The angle between the direction of the current and the magnetization, and hence the value of the resistance, changes when a magnetic field is applied in the y -direction. The change in resistance as a function of this field-strength H_y is approximately parabolic, as shown by the dashed curve in fig. 4b. If $H_y \approx 0$ for such a curve, the output signal of the transducer for a small change in H_y is very weak. This problem is solved by applying a structure of conductive material to the strips so that the angle between the direction of the current and the magnetization at $H_y = 0$ is equal to 45° ; see fig. 4c. The relation established in this way between the strength of the external field and the change in resistance is approximately linear in a fairly wide range, as shown by the continuous curve in fig. 4b. The maximum sensitivity of such a transducer is about 0.1 mV per A/m at a supply voltage of 5 V.

The electronic compass may contain three such transducers for the magnetic field. They measure the components of H_E , the strength of the terrestrial magnetic field, in three orthogonal directions. For the moment only the components $H_E \sin \phi_n$ and $H_E \cos \phi_n$ are of interest. (In future we shall use the vertical component to determine the gradient of the road.) H_E is about 15 A/m. Measurements of such a weak magnetic field are very difficult because of the slow variation of the output signal, due to temperature changes. An answer to this problem is to convert the output

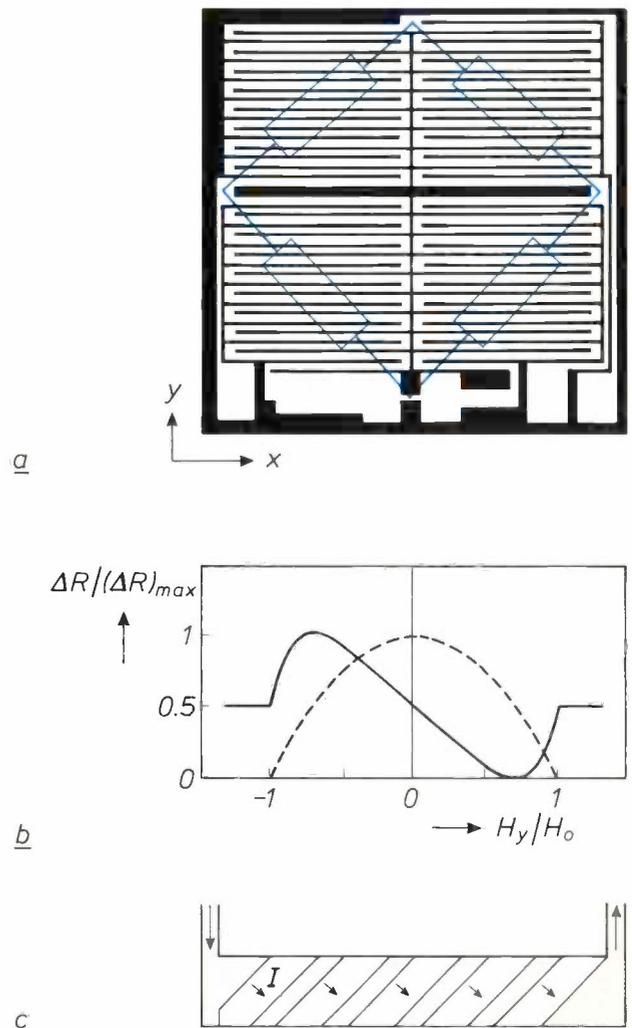


Fig. 4. a) The interior of a KMZ10 magnetic transducer (a Philips Valvo product); three of these can be used in the electronic compass. The transducer measures magnetic fields via the magnetoresistance effect. The transducer is composed of ferromagnetic strips, connected to form a Wheatstone bridge (shown schematically in blue). b) The relative change in the resistance of a strip in (a) as a function of the ratio of the magnetic field-strength in the y -direction to a reference field-strength H_0 [6]. The dashed curve relates to an untreated strip. The continuous curve, which is a straight line over much of its range, is applicable if a structure of conductive material is applied to the strip. c) A strip treated in this way. The grey areas correspond to conductive material. In the regions in between these areas the angle between the current I and the magnetization is 45° for $H_y = 0$.

signal of the transducer into a 'square-wave' signal. This is done by periodically changing the direction of magnetization in the transducer strips by energizing a coil around the transducer with alternate positive and negative pulses of current. Synchronous demodulation of the rectangular voltage then converts it into a direct voltage, with none of the low-frequency components that interfered with the original signal.

[6] S. Lefschetz, Introduction to topology, Princeton Univ. Press, Princeton, N.J., 1949.

[6] W. J. van Gestel, F. W. Gorter and K. E. Kuijk, Read-out of a magnetic tape by the magnetoresistance effect, Philips Tech. Rev. 37, 42-50, 1977.

Transducers of the same type are used for measuring the number of revolutions of the non-driven wheels. Fig. 5 shows how two transducers are mounted in relation to a strip bonded to the inside edge of a wheel rim. The strip is permanently magnetized so that the magnetic field is radial and alternately inwards and outwards, and so that 26 periods of the magnetization, each of length p , correspond to a rota-

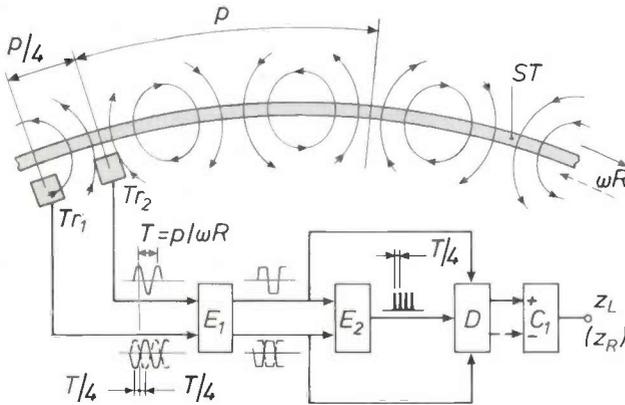


Fig. 5. Principle of a wheel sensor. ST permanently magnetized strip. p length of the period of the radial magnetization. ωR peripheral velocity of the strip; ω angular velocity, R radius. Tr_1 and Tr_2 magnetic transducers; see fig. 4. The output signal is approximately sinusoidal with a period $T = p/\omega R$. The phase difference between the output signals from the two transducers is $T/4$ and is positive or negative, depending on the direction of rotation of the wheel. E_1 circuit that converts the signals into square waves. E_2 circuit that derives pulse trains of period $T/4$ from the two square waves. D circuit that decides whether a positive or negative value should be assigned to the pulses. C_1 counter that supplies a binary number z_L (z_R), which is a measure of the angular displacement of the left-hand (right-hand) wheel. (The circuits are digital except for E_1 .)

tion of the wheel through 360° . The spacing of the transducers is $p/4$. The figure also shows how the transducer signal is electronically processed. When the wheel is turning, the transducers each supply an approximately sinusoidal signal, with a 90° phase difference between the two signals. The direction of rotation of the wheel determines which of the two signals leads in phase. The sinusoidal signals are amplified and converted into square-wave voltages. These are transformed into a pulse train in such a way that one wheel revolution corresponds to $4 \times 26 = 104$ pulses. At the same time the decision is made as to whether each pulse should be positive or negative. Finally, a binary number z_L is produced, which is a measure of the angular rotation of the left-hand wheel with respect to the previous dead-reckoned fix. In the same way a binary number z_R is produced, which is a measure of the angular rotation of the right-hand wheel.

Values for the quantities δ_n and ϕ_n defined in equations (1a) and (1b) are found from the measured results from the sensors. The distance travelled δ_n and the change in heading $\phi_n^W - \phi_{n-1}^W$ are determined with the aid of the wheel sensors by using the equations:

$$\delta_n = \frac{(z_L + z_R)\pi R}{104} \quad (2)$$

and

$$\phi_n^W - \phi_{n-1}^W = \frac{(z_L - z_R)\pi R}{52S} \quad (3)$$

where R is the wheel radius and S the track width. The heading measured by the electronic compass is denoted by ϕ_n^C . A weighted mean of the measured values for the heading is given by:

$$\phi_n = \frac{A\phi_n^W + B\phi_n^C}{A + B} \quad (4)$$

where A and B are weighting factors. Now we shall discuss a relatively simple way of obtaining a value for ϕ_n that is virtually free from interference.

Combining the measured values for the change in heading

We have seen how the signal from the electronic compass and the difference signal from the wheel sensors can both provide a measure for the heading of the vehicle. The compass signal supplies an absolute value of the heading relative to magnetic north; the wheel signal supplies a relative value with respect to the initial value of the heading. The signals must thus be made comparable. Both signals are also subject to interference because of the way in which the signals have been obtained. The difference signal from the wheel sensors has a slow drift due to differences in tyre pressure and wear, and so on. The compass signal contains very fast fluctuations due to the magnetic fields of other vehicles, steel bridges, structural metal in buildings, and tramlines. The compass signal also has a constant error due to the magnetic field of the steel body of the car; see fig. 6.

Fig. 7 shows how the two signals for the heading are corrected and the interfering signals removed. In the computing unit C_2 the successive numbers of pulses z_L and z_R from the left-hand and right-hand wheels are converted into a signal ϕ^W , which still contains an unwanted low-frequency component. In the computing unit C_3 the effect of the car's own magnetic field is eliminated from the signals for $H_E \cos \phi$ and $H_E \sin \phi$. This is done by calculating the equation of a circle of best fit in the diagram for $H_E \cos \phi$ as a function of $H_E \sin \phi$ (see fig. 6) for a large number of measured results. Next, C_3 calculates the signal ϕ^C for the head-

ing. The values of ϕ^W and ϕ^C are each multiplied by a constant factor to make them comparable with each other in magnitude. The signal ϕ^W is then passed through a highpass filter and the signal ϕ^C is passed through a lowpass filter. Both filters have the same cut-off frequency f_c . The interfering components are almost completely removed in the filters.

The special feature of our method is the continuous adjustment of the cut-off frequency f_c of the filters by the control unit R . The interfering components in ϕ^W and ϕ^C are measured in R during a certain period of time by subtracting the mean value from the signals.

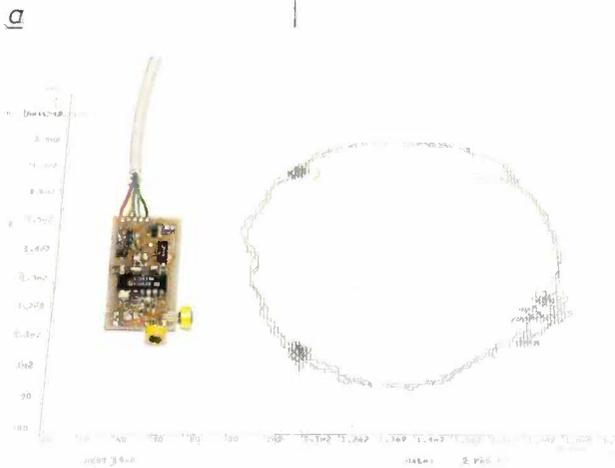
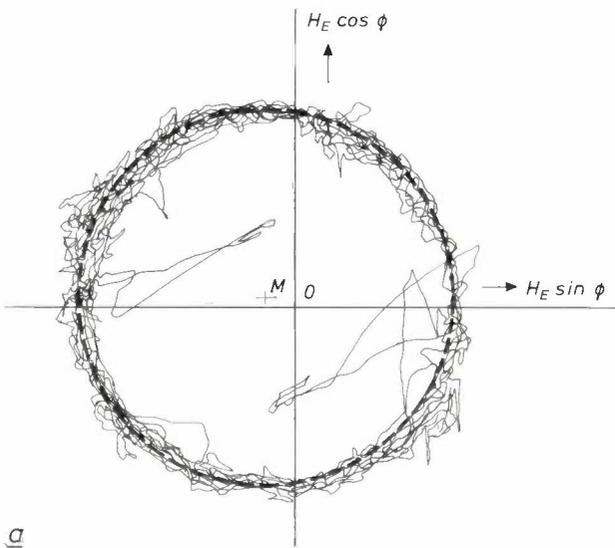
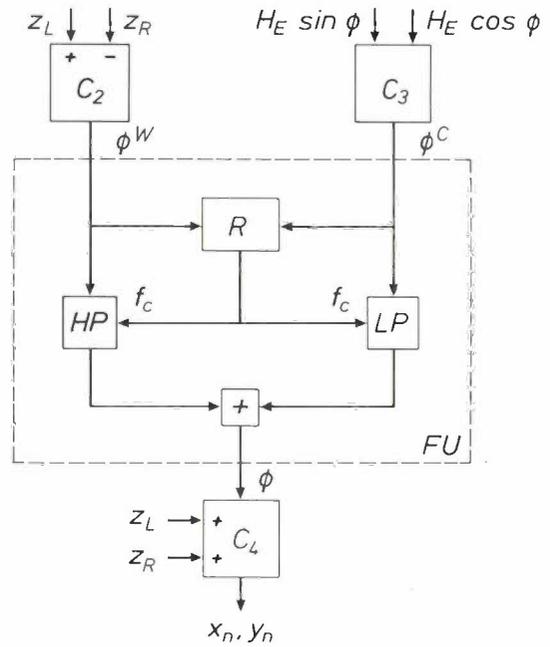


Fig. 6. a) The component $H_E \cos \phi$ of the Earth's magnetic field H_E as a function of the component $H_E \sin \phi$; ϕ heading of the vehicle. The graph is a recording, made during a journey, of measurements made with two magnetic transducers in the electronic compass. Theoretically the measured points should lie on a circle with centre O . The best-fitting circle, however, has M as its centre, because of the effect of the magnetic field of the vehicle. The other deviations are due to external magnetic fields from passing vehicles, steel bridges, electrical cables and so on. b) Photograph of such a recording, combined with the interior of an electronic compass with two transducers (see also fig. 4).

Fig. 7. Processing the signals z_L and z_R from the wheel sensors and $H_E \sin \phi$ and $H_E \cos \phi$ from the electronic compass. C_2, C_3, C_4 computing units. FU filter unit. In C_2 the signals z_L and z_R are converted into a signal ϕ^W for the heading. In C_3 the compass signals are converted into a signal ϕ^C comparable with ϕ^W . HP highpass filter. LP lowpass filter. f_c cut-off frequency of the two filters. R control unit that continuously controls the magnitude of f_c , by measuring the r.m.s. value of the interfering components in ϕ^W and ϕ^C . After addition, the resultant signal ϕ is periodically converted into values x_n and y_n for the coordinates of position, with aid of the sum signal from the wheel sensors.

The root mean square (r.m.s.) value of each of the interfering signals is then determined. If the r.m.s. value of the interference in ϕ^W is large, f_c is increased; if this r.m.s. value is large for ϕ^C , then f_c is reduced. Finally the two filtered signals are added together. Setting the value of f_c also weights ϕ^W and ϕ^C as in equation (4). The computing unit C_4 calculates coordinates of position x_n and y_n from samples of the signals and from the sum signal of z_L and z_R , which is a measure of the distance travelled δ_n . This calculation makes use of equations (1a) and (1b). The corrections described are not executed with samples of analog signals, but are derived from calculations with binary numbers, so that analog/digital converters are used where required. The calculations are made in a special computer for navigation.

Improving the dead-reckoned fix from map data

Errors in the fix arrived at by dead reckoning may be divided into systematic and random errors. The magnitude of systematic errors is essentially predictable, and they can be eliminated from the results of the measurements by calibration. Random errors cannot be eliminated by calibration.

Once the system has been installed in a car, several calibrations are necessary. First the car is driven around a closed loop for the compass to be calibrated. The effect of the car's magnetic field is then eliminated and the transducers are matched to each other. The computer does this by fitting an imaginary circle (a regression circle) to a diagram resembling the one in fig. 6. Then the car is driven along a straight line, so that the wheel sensors can be calibrated by setting the difference signal equal to zero and comparing the sum signal with the distance actually travelled. Such calibrations may also be necessary after new tyres or snow chains have been fitted.

Calibration is also carried out during a journey. The sum signal of the wheel sensors is then compared with the length of map segments and the electronic compass is calibrated as described above. The remaining errors, the ones that cannot be eliminated by calibration, have to be treated as random errors.

The magnitude of random errors can be estimated, since their standard deviation can be calculated with a fair probability from repeated measurements. The standard deviations of the random errors in the measurements of heading and distance can be obtained from a comparison of a large number of dead-reckoned fixes with the actual position of the vehicle. The magnitude of the random errors is expressed by the dimensions of the VLPA, the Vehicle Location Probability Area. In theory this area is bounded by an ellipse whose dimensions are given by rules for the propagation of random errors.

The size of the VLPA increases continuously during the journey. The increase is abruptly interrupted when a position correction is made, where a position correction is a correction of the dead-reckoned fix, associated with a reduction in the uncertainty of position. A position correction can be made when the computer has recognized the pattern of a number of successive fixes by comparing them with topographical information in the main memory. The procedure used, which will now be described, requires a considerable amount of computation. To keep this within reasonable bounds, the VLPA ellipse is approximated by its smallest enclosing rectangle.

Fig. 8 shows how the comparison of map data and dead-reckoned fix is made. Whenever a dead reckoning is made, the computer finds out which segments lie inside the VLPA rectangle. A record of these segments is kept in the form of 'trees' of possible routes, as illustrated in fig. 8b. The lower-case letters in the trees refer to segments whose direction does not really correspond to the heading of the vehicle. If such a segment lies outside the VLPA in the next dead reckoning, it is removed from the tree. A position correc-

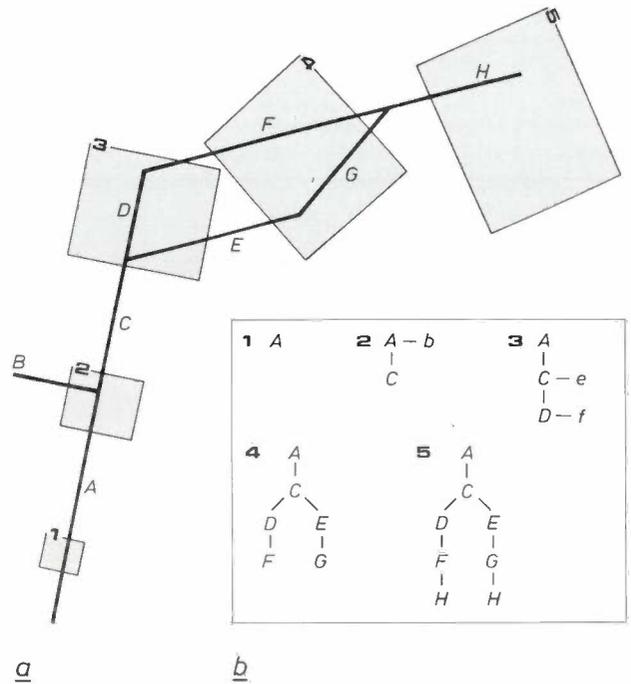


Fig. 8. Comparison of map data with dead-reckoned fix. *a*) The computer finds out which map segments lie inside a 'vehicle location probability area' (VLPA). *b*) 'Trees' of possible routes resulting from this comparison. Segments whose direction does not correspond to the heading of the vehicle are shown in lower-case letters. The size of the VLPA steadily increases because of the increasing uncertainty of position. The tree at VLPA 5 shows that there are two possible routes. (In reality the successive VLPAs overlap each other, since their centres are typically spaced at a distance of 5 m.)

tion (an update) is made if the computer can decide with a good degree of certainty which of the possible routes has been followed.

Fig. 9 shows schematically the position correction made if the computer has recognized the pattern of the centres of successive VLPAs as that of the segments A, C, D, F and H in fig. 8. (For clarity fewer VLPAs are shown in fig. 8 than in fig. 9.) The distances between the segments on the map and the lines of best fit (shown red in the figure) through these centres give two components of a correction vector. After the position correction the VLPA is much smaller and grows again until the next position correction is made.

In addition to position corrections there is also a 'map-matching procedure'. This is a procedure in which, at every dead-reckoned fix, the centre of the VLPA is projected on to a segment inside it. If there are two or more segments inside the VLPA, the choice of the segment is then determined by factors such as whether the segment is on the planned route or not, or the correspondence between the direction of the segment and the heading. Map matching can be used, for

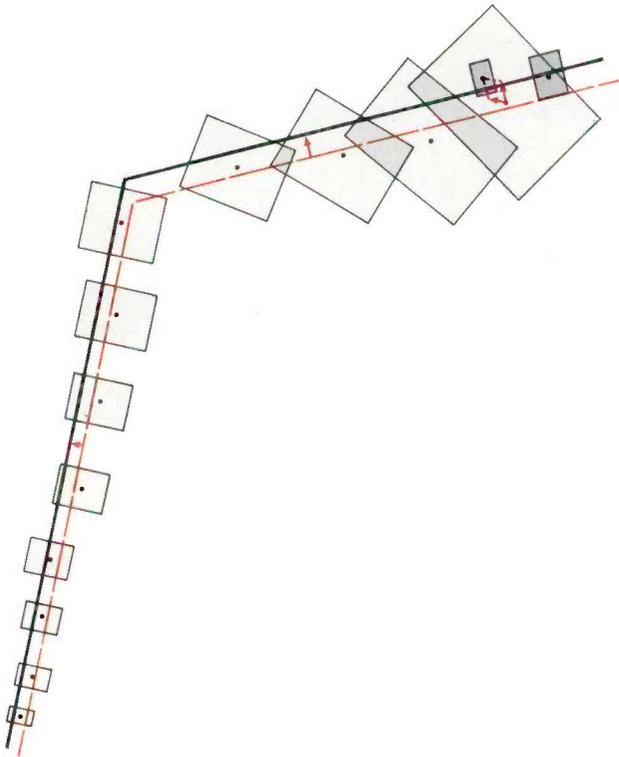


Fig. 9. A position correction after the computer has recognized the road pattern. This pattern is compared with the routes *A-C-D-F-H* and *A-C-E-G-H* on the map, see fig. 8, and the second route is rejected. Best-fitting straight lines (red dashes) through the consecutive centres of VLPAs yield two vectors. The resultant gives the displacement of the centre of the last VLPA, and the positional uncertainty immediately becomes smaller. It grows again as the distance travelled increases.

example, to display the correct segments, nodes and intermediate points on the LCD screen. The position resulting from map matching is always at the centre of the display, which therefore always corresponds to a part of a road.

The Compact Disc as a carrier of digital map data

As mentioned earlier, we decided to use the Compact Disc Interactive, or CD-I, as the storage medium for the CARIN system. An international standard for CD-I is being prepared. The use of a standardized method of storage for the information has the advantage that the data carriers for different navigation systems are interchangeable and that savings can be made in research and development costs. The CD-I storage medium was developed for consumer-electronics applications by the licensors of the Compact Disc Digital Audio system. A CD-I can store audio and video signals of different quality levels, combined with computer data. The proposed CD-I standard expresses detailed agreements relating to the storage of computer data, coding and error-correction sys-

tems, the operating system and the connection of peripherals. This means that discs and systems from different manufacturers can be compatible.

In addition to digital map data, which will occupy most of the disc, a CD-I for geographical applications will shortly contain applications software. Examples of this are parts of the navigation software, the route-planning program and the driver-guidance program. Programs can also be included that will make the CD-I with its player in the home into a kind of interactive atlas. The disc can also be used for teaching geography, if the topographical information is augmented by photographs, drawings and even sound and — with a few restrictions — moving video pictures. Depending on the options available with the hardware, it will soon be possible to plan a route at home, study it on a video display and take a printed copy.

In the development of a geographical CD-I, use is made of a topographical database. The rules to be observed in such a database are also laid down. The database will be regularly updated by publishers of (electronic) geographical data, and it can also be used for other applications, e.g. for town and regional planning, for keeping records of traffic accidents and for processing census data. The part of the database intended for use by CARIN is given the correct structure by a conversion program. The file thus obtained is combined with application software and possibly with other data. The result is then formatted to the CD-I standard. In simple terms, this means that the data is divided into blocks of 2 kilobytes (16 384 bits), each with a 'header', with synchronization bits and with bits for later error correction. After coding and modulation (EFM: eight-to-fourteen modulation) the data is in the correct form for writing to the disc [7]. The blocks are converted into sectors, and the address is included in the header.

The result of these operations is that the successive sectors form a sequence of 'channel bits'. Each '1' in the sequence of channel bits is a land/pit or pit/land transition in the spiral track of pits on the final disc [8]. The pattern of pits and lands is transferred to a photo-sensitive layer on a rotating disc with the aid of a laser. Development of this layer produces a 'master', and the moulds for manufacturing the discs are made from the master in a number of intermediate steps.

As noted earlier, the storage capacity of a Compact Disc is 4800 Mbit. The main memory of the compu-

[7] J. P. J. Heemskerk and K. A. Schouhamer Imminck, Compact Disc: system aspects and modulation, Philips Tech. Rev. 40, 157-164, 1982.

[8] M. G. Carasso, J. B. H. Peek and J. P. Sinjou, The Compact Disc Digital Audio System, Philips Tech. Rev. 40, 151-155, 1982.

ter, on the other hand, has a capacity of only about 8 Mbit. While this is still a very respectable capacity for a semiconductor memory, it is small compared with that of a Compact Disc. The time required for reading all the information from a disc is about an hour. The player used in the CARIN system is an ordinary Compact Disc player for music playback in cars. Moving the optical pick-up in the player from the smallest radius to the largest radius of the useful area on the disc takes one or two seconds. It will take less time in the later generations of players. It is not possible to determine beforehand where exactly the pick-up will arrive after a movement. Since it is not possible to read 'backwards' when the pick-up has to be moved to a particular address on the disc, the data reading must start well before the address is reached. This means that extra information has to be read before the correct address is found, and this takes time. (This extra information is not transferred to the main memory, of course.)

As explained, the access time is not negligible, and so the data required for navigation cannot be read from the Compact Disc at the exact moment when it is required. At the start of a journey, therefore, all the topographical information required for navigation is as far as possible retrieved, read and stored in the main memory. Since this must be done fairly quickly, it makes sense to store the data on the Compact Disc in the most advantageous way.

The topographical information is divided into 'parcels'. The information in each parcel is stored as a 'block' (or 'bucket') of data on the Compact Disc and is distributed over an integer number of successive sectors. Also stored on the Compact Disc is a list of addresses, or 'directory'. This includes the address of the first sector where a block is stored, the number of sectors used for the block, and the coordinates of the lines that define the corresponding parcel on the map. A request for map data from the computer generally means that map data is called for a rectangular area; see *fig. 10*. The directory is then consulted to find out which blocks contain information about the requested area. The information can then be retrieved from their addresses.

The time required for these operations can be short if two requirements are satisfied. The map must be divided into convenient parcels, and the corresponding blocks should be stored on the Compact Disc in a sequence that does not require too much movement of the pick-up. The first requirement means that all the blocks should contain about the same amount of data, which must not be so large that the blocks requested for a particular area take up too much space in the main memory. Nor, on the other hand, should

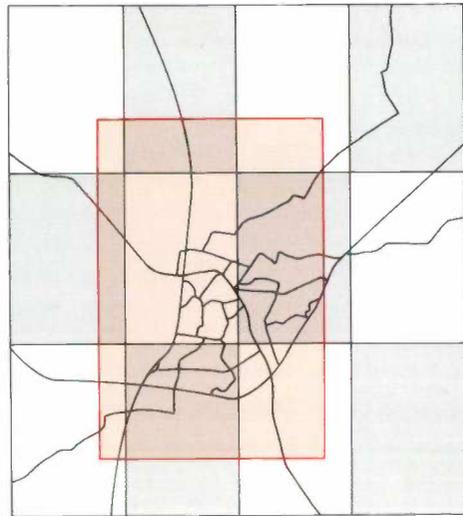


Fig. 10. Partitioning of the map into 'parcels' — the white and grey rectangles. When the computer requests map data from the rectangle shown red, a list of addresses (the directory) supplies the location of the sectors on the Compact Disc where the nine blocks with information about the required nine parcels are stored.

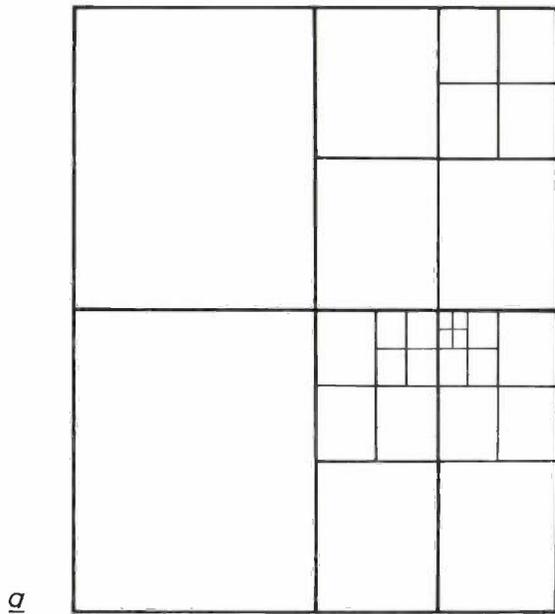
there be so many blocks that too much space is required for the directory or that too many movements of the pick-up are necessary. The second requirement means that parcels adjacent on the map should correspond to blocks close to each other on the disc. These problems will be discussed at greater length in the next section.

Partitioning and arranging the map data

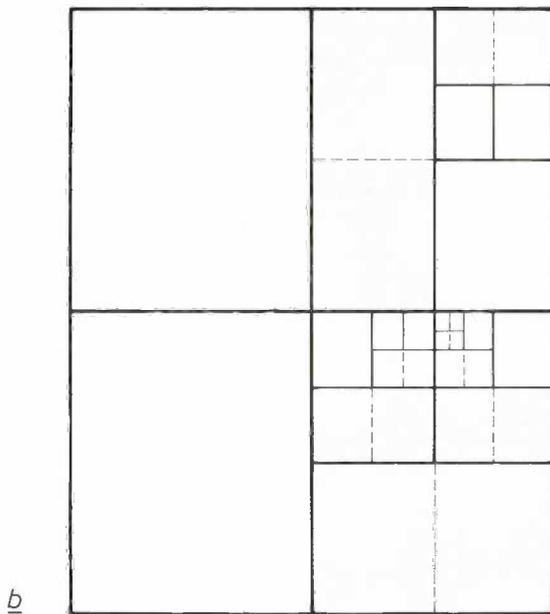
A map should be partitioned into parcels in such a way that the data content in each corresponding block on the disc is about the same. Parcels will therefore be smaller in towns than in rural areas. Several partitioning algorithms are known from the literature. The algorithm we have used is known as the 'region quad tree' [9]. In this algorithm the map is continuously subdivided into four rectangles until the information content per parcel is smaller than a preset value; see *fig. 11a*. We have adapted the algorithm so that, after the repeated partitioning, some pairs of adjacent parcels are joined together again to make the data content of the sum no more than slightly larger than the preset value; see *fig. 11b*.

The problem of arranging the data blocks resulting from the map partitioning on the disc is equivalent to

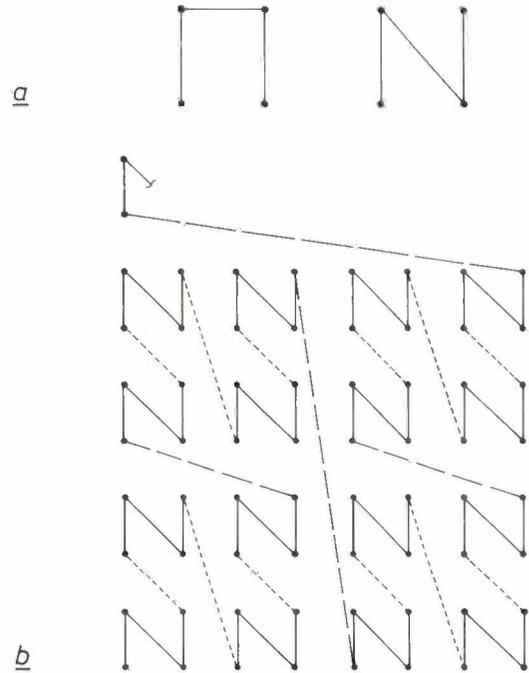
- [9] H. Samet, The quadtree and related hierarchical data structures, *Comput. Surv.* **16**, 187-260, 1984.
 [10] G. Peano, Sur une courbe, qui remplit toute une aire plane, *Math. Annalen* **36**, 157-160, 1890;
 E. A. Patrick, D. R. Anderson and F. K. Bechtel, Mapping multidimensional space to one dimension for computer output display, *IEEE Trans. C-17*, 949-953, 1968.



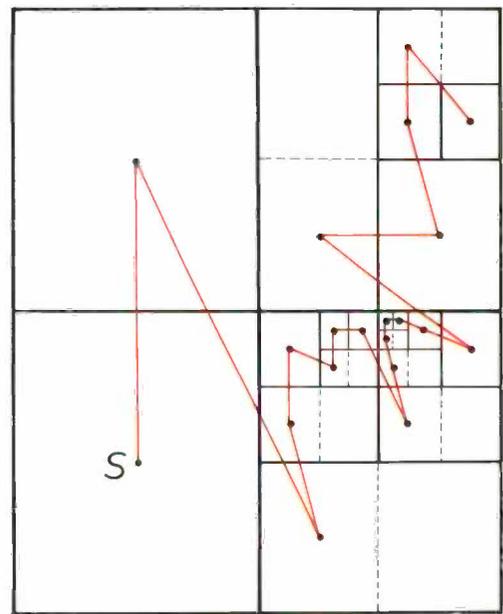
a



b



b



c

Fig. 11. Map partitioning with the 'region quad tree' algorithm^[9]. a) Repeated partitioning into four rectangles until the data content in each block that corresponds to a parcel is smaller than a preset value. b) Some of the pairs of adjacent parcels are then joined so that the block has a data content never more than slightly larger than the preset value. The parcels formed from these mergers are shown grey.

Fig. 12. Generation of a Peano curve^[10]. a) The two possible seed curves for four points of a rectangular network. b) A Peano N-curve, generated from the seed curve on the right. The curve passes once through all 64 + 2 points of the network. c) The Peano N-curve that connects the parcels of fig. 11c and starts at S. In the case of a parcel formed from a summation, the centre of the first original parcel is taken; the other parcel is then ignored.

translating a two-dimensional structure into a one-dimensional structure. To solve this problem we used 'space-filling curves' whose construction was described in 1890 by Giuseppe Peano^[10]. A space-filling curve passes through each point in a space once only. In our case we are working in a plane, and Peano's construction is based on a square network of points; the result is called a Peano curve. Points close to each

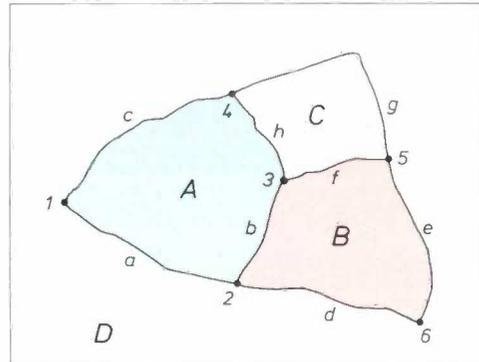
other in sequence along the curve are also generally close together in the plane. The converse is also usually true. A Peano curve is obtained by starting from a 'seed curve', see fig. 12a. The two curves drawn through four points of the network are the only possible seed curves of these dimensions — apart from rotations or reflections. A Peano curve with the left-hand seed curve is sometimes referred to as a Hilbert

curve, and a Peano curve with the right-hand seed curve is known as a Peano-N curve. Fig. 12b shows how a Peano-N curve is generated. Beginning with the four points at the lower left, blocks of 16, 64, 256 points and so on are connected up.

The Peano curve is particularly suitable for arranging the parcels resulting from region quad tree partitioning. Fig. 12c shows how a Peano-N curve is generated for the map partitioning in fig. 11b. The curve connects the centre-points of the different parcels on the map, with *S* as the starting point. First of all an attempt is made to describe an 'N' that contains the parcels whose area is a quarter of that of the map; these are the first-degree parcels. When such a parcel has been partitioned, an 'N' is first described through the centres of the parcels of the second degree, and so on. When two parcels of a particular degree have been joined together the process continues as if they had not been joined together. The centre of the first original parcel is taken and that of the other one is now ignored. The process is applied consecutively to parcels of increasing degree, i.e. with a smaller and smaller area.

It can be seen from fig. 12c that the small parcels in the 'town area' to the right of the centre are close together on the Peano curve. Arranging the corresponding blocks on the Compact Disc in the same sequence as the points on the Peano curve therefore has the intended result that the retrieval of information about adjacent areas requires little or no movement of the pick-up. The algorithm we finally chose for arranging the blocks works in a slightly different way, but the result is the same — a relatively short time for retrieving them.

mediate points on the road connecting the nodes, form a chain. The data management works with the concepts of a 0-cell, which is a node, a 1-cell, which is a chain, and a 2-cell, which is the area inside a number of connected chains. While navigating it is sometimes necessary to retrieve information from the main memory about the chains that surround a 2-cell. This could be necessary if the car is no longer on a road on the map, e.g. on a rough track, in a car park, or on private land. The information about the chains corresponding to a 2-cell on the right or left of the 1-cell where the car was last travelling is then requested by the computer. As soon as the car leaves the area of the 2-cell, the computer can usually determine the 1-cell where the car is located. In the pattern-recognition



a

key	information	pointer to adjacent 1-cell
1	xxxx	c
2	xxxx	b
3	xxxx	h
4	xxxx	h

b

Aspects of data management

In the previous section we saw how the data blocks that correspond to the parcels are stored sequentially on a Compact Disc. For efficient management of the information read from the Compact Disc it is most important that the information is efficiently structured inside each block. As we saw at the beginning of this article, the topographical information is built up from the coordinates of nodes and intermediate points. A sequence of points starting and ending at a node, together with the inter-

key	information	start node		end node		2-cell left	2-cell right
		pointer	thread-pointer	pointer	thread-pointer		
a	xxxx	1	c	2	b	A	D
b	xxxx	2	d	3	h	A	B
c	xxxx	1	a	4	g	D	A
d	xxxx	2	a	6	e	B	D
e	xxxx	6	d	5	f	B	D
f	xxxx	3	b	5	g	C	B
g	xxxx	4	h	5	e	D	C
h	xxxx	4	c	3	f	C	A

c

Fig. 13. a) 0-cells, 1-cells and 2-cells; see the caption to fig. 3. b) Schematic representation of the list of 0-cells forming part of the directory on the Compact Disc. c) List of 1-cells. A separate list for 2-cells is unnecessary, because the thread-pointers in the fourth and sixth columns point unambiguously to the 1-cells that form the boundary of a 2-cell to the left or right of a 1-cell. The thread-pointer at the start or at the end node is the 1-cell calculated clockwise from the 1-cell that contains the node. For 1-cell *b* the figure shows the method of finding the 1-cells of the 2-cell on the right (*B*, shown red) and the 2-cell on the left (*A*, shown blue). For more details, see text.

procedure described above it may also be necessary to request the chains for a 2-cell.

Let us now consider the structuring of the information inside a block. The addresses in the memory where the data for 0-, 1- and 2-cells are stored are linked by means of 'lists'. Fig. 13 shows two such lists, which relate to the parcel shown schematically in fig. 3. (For simplicity, this parcel does not include any roads that cross its boundaries.) The lists use reference addresses called 'pointers' and 'forwarding addresses' called 'thread-pointers', which link the various memory locations. Fig. 13*b* gives a list of the 0-cells inside the parcel. The first column contains the designation or 'key' of the cell, the second the information in the form of coordinates, etc., and the third the 'pointer' to a 1-cell that has the 0-cell as the start or end node. (We shall not consider here which 1-cell is most suitable for this.) Fig. 13*c* gives a list of the 1-cells in the parcel. The columns contain from left to right the key of the 1-cell, the corresponding information, the pointer and thread-pointer for the start node, the same for the end node, and the keys of the 2-cells located to the left and right of the 1-cell.

The thread-pointers in the list given in fig. 13*c* can be used to find out which 1-cells enclose a 2-cell located to the left or right of a 1-cell. The thread-pointer indicates a 1-cell which, as seen from the node, is clockwise from the 1-cell to which the row in the list

relates. The 1-cells that surround a 2-cell on the right of a 1-cell in the list are found by first taking the thread-pointer of the start node of the 1-cell and then always taking the thread-pointer that belongs to the opposite node of the 1-cell indicated by the thread-pointer. The same procedure is applied to the 2-cell located on the left of the 1-cell, but starting with the thread-pointer that belongs to the end node of the 1-cell. This is illustrated in fig. 13*a* and *c* for the 2-cell located on the right (*B*, shown in red), and the 2-cell on the left (*A*, shown in blue) of 1-cell *b*.

The structure of the data on the Compact Disc and in the CARIN main memory is built up from many such carefully designed links. It has enabled us to increase the speed of data retrieval, yet without introducing redundancy.

Summary. The CARIN (CAR Information and Navigation) system plans the best route and guides the driver with spoken directions from a speech synthesizer during the journey. The system does this by making periodic fixes by dead reckoning, with the aid of wheel sensors and an electronic compass. The dead-reckoned fix is continuously updated by comparing it with map data read from a CD-I (Compact Disc Interactive), which can store 4800 Mbit of digital data. Methods have been developed for conveniently partitioning the map data into 'parcels' and arranging the corresponding data blocks in the most efficient sequence on the Compact Disc. An important data-management aspect is the efficient way in which the information is structured in each block: the minimum amounts of information, which correspond to the nodes and 'chains' on the map, are linked via pointer addresses.

Behaviour of the oxide film in MOS devices

D. R. Wolters

The MOS transistor, which contains an insulating SiO₂ film between silicon and gate, is an important component for the semiconductor industry. The oxide film is usually thinner than 0.1 μm, which means that the current in the silicon can be controlled effectively by the gate voltage. The electric field-strength in the oxide film can be very high, a hundred times more than would be possible in high-voltage cables. The charge leakage through the film must be very small, of course and the film must not break down. With the continuing miniaturization of semiconductor components, containing even thinner oxide films, it is becoming increasingly difficult to meet these requirements. It is therefore important to have a detailed understanding of the processes that take place in the film when a voltage is applied to it and to recognize the conditions that lead to breakdown. In the article below the author describes an investigation that he and other colleagues at Philips Research Laboratories have made into the behaviour of the oxide film in MOS devices.

Introduction

Since the fifties the semiconductor industry has made enormous progress in the technology of manufacturing electronic circuits. This owes much to the use of silicon as the basic material. Silicon is a semiconductor possessing a unique combination of properties. It is particularly suitable for the manufacture of large pure single crystals that can be processed by sawing, grinding and polishing into 'slices' or 'wafers'. Electrical conduction in the silicon is easily controlled and can be varied considerably by adding appropriate dopants.

In the sixties the discovery that it was useful to cover a silicon wafer with a thin film of silicon dioxide (SiO₂), e.g. by oxidation at about 1000 °C in an atmosphere of oxygen or water vapour, lent considerable impetus to the development of 'planar' technology. Films made in this way are dense and homogeneous, and they adhere extremely well to the silicon surface. Because they are chemically resistant and impervious, they can be used as doping masks in the manufacture of integrated circuits. Another useful property is that they are good electrical insulators, which is use-

ful for example in the stabilization of bipolar devices. Mastery of the technique of manufacturing very pure oxide films was particularly important in the development of MOS transistors (MOS = metal/oxide/semiconductor). In these field-effect transistors an SiO₂ layer acts as an insulator between doped silicon and the gate^[1]. Electrical conduction in the silicon between the source and drain is controlled by the voltage on the gate. The properties of the SiO₂ films are such that films thinner than 50 nm can be used. This means that the gate in such a device can be extremely close to the silicon to control the electrical conduction in the silicon effectively.

The amount of charge that leaks through an SiO₂ film or is trapped when a voltage is applied must be extremely small, and there must be no premature electrical breakdown, which would make the device useless. The miniaturization of MOS devices to sub-micron dimensions makes it increasingly difficult to meet these requirements: thinner oxide films and stronger electric fields are the requirements of the day. For these and other reasons the manufacture of the films requires extremely accurate process control. It is also necessary to know in detail what happens when a

Dr Ir D. R. Wolters is with Philips Research Laboratories, Eindhoven.

voltage is applied to the film. At Philips Research Laboratories in Eindhoven extensive investigations have been made in recent years into the processes that play a part in the manufacture of the films, in their application in MOS devices and in electrical breakdown [2].

A study of the oxidation kinetics in the manufacture of the film has shown that the mechanism of the oxidation is not the same as has until recently been assumed [3]. In particular, current ideas about the catalytic action of the water vapour present were not really satisfactory. The distribution of radioactive water in the films showed us that water molecules ionize in pairs and that the resulting ions (H_3O^+ and OH^-) travel through the oxide film. The ion migration takes place close to the interface between the silicon and the growing SiO_2 film, while the migration of oxygen molecules takes place closer to the surface [4].

When a high electric field is applied to an SiO_2 film, a certain amount of charge is always injected. Some charge leakage is not in itself so serious, but some of the injected charge becomes trapped and may therefore affect the control of current conduction in the silicon. By accurately measuring the trapping rate and the amount of trapped charge we have been able to model the mechanism of charge trapping and the nature of the trapping centres. We have accounted for the fact that trapped charges have the effect of counteracting any further charge trapping in their immediate vicinity.

Charge leakage is an important factor affecting the life of the SiO_2 film and hence the life of the device. The occurrence of electrical breakdown turns out in fact to be determined by the amount of charge that has leaked through the film. It seems as if the leaking charge gradually 'wears out' the film. This suggests that the 'wear', i.e. a change in properties, starts long before the film breaks down. The understanding thus obtained now enables us to predict fairly accurately when a device will start to deteriorate due to wear of the oxide film and when it will finally break down. We can therefore give better estimates of the useful life of MOS devices.

In this article we shall first consider charge trapping in an SiO_2 film. It will be shown that the repulsion between trapped and injected charges is an important factor. Next we shall look at the phenomenon of electrical breakdown, examining defect-related breakdown, intrinsic breakdown (i.e. breakdown not related to defects) and the effects of field-strength, current and transferred charge. We shall also take a closer look at the conditions under which intrinsic breakdown occurs. Finally, the mechanism responsible for breakdown will be indicated.

Charge trapping

As noted above, when an electric field is applied across an insulating SiO_2 film, some charge will be injected. The current that flows through the film decreases after some time because some of the injected charge is trapped. The space charge thus built up in the film causes a distortion of the electric field at the interface with the silicon. This has the effect of broadening the potential barrier for electrons at the interface between silicon and the SiO_2 film; see *fig. 1*. As a

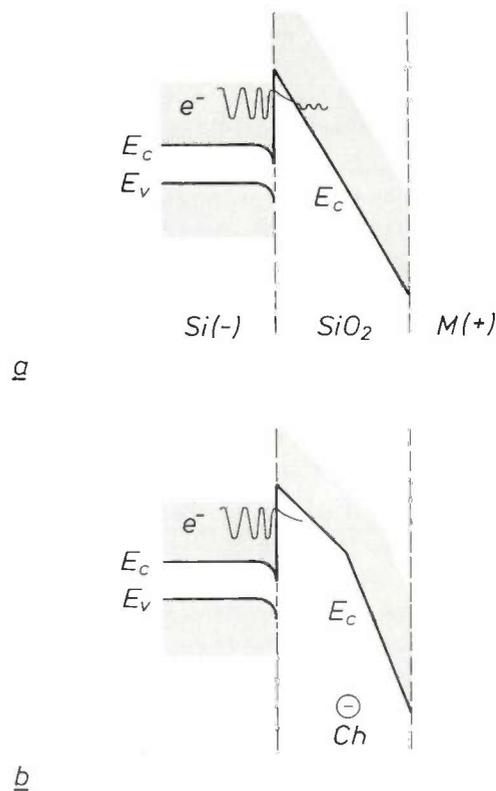


Fig. 1. *a*) Energy diagram for the transition of electrons from silicon to SiO_2 in a MOS device in which the metallic layer M is at a positive potential with respect to the silicon. E_v upper edge of the valence band. E_c lower edge of the conduction band. An electron e^- goes from the conduction band of silicon to the conduction band of SiO_2 . At the interface the electron tunnels through the potential barrier. *b*) In the presence of a negative space charge Ch , caused by the trapping of electrons, the potential barrier becomes wider. An electron tunnelling through the barrier therefore has to cover a greater distance, which reduces the probability of electron injection into the SiO_2 film.

- [1] See the special issue on MOS transistors, Philips Tech. Rev. 31, 205-295, 1970.
- [2] This investigation has been described in detail in the author's doctoral thesis, entitled: Growth, conduction, trapping and breakdown of SiO_2 layers on silicon, Groningen 1985.
- [3] A widely used model for the oxidation is described by B. E. Deal and A. S. Grove, General relationship for the thermal oxidation of silicon, J. Appl. Phys. 36, 3770-3778, 1965.
- [4] D. R. Wolters, On the oxidation kinetics of silicon: the role of water, J. Electrochem. Soc. 127, 2072-2082, 1980.

consequence, there will be less chance of electrons 'tunnelling' through the barrier, and fewer electrons will be injected into the film. The time-dependent injection and trapping of electrons can assume such proportions that the film can no longer perform stably and reproducibly in an electronic device. In MOS transistors, which often operate with a high field-strength (up to 3×10^6 V/cm), the space-charge effects lead to instabilities and drastic changes in the switching or threshold voltage. It is therefore most important to be able to characterize and control these effects.

Charge trapping can usefully be studied in a MOS capacitor, in which the oxide film is situated between an aluminium (or polysilicon) gate and a silicon substrate. It is possible, for example, to measure the increase in voltage necessary to make a constant current flow through the capacitor. From the measured voltage the increase in the number of trapped charges can be derived directly. Information about the number of trapped charges can also be obtained by measuring the capacitance as a function of the applied voltage during interruptions of the injection.

In the past a first-order approximation has generally been used to describe charge trapping as a function of time [6]. It was assumed that the trapping rate was equal to the product of the electron flux, the effective cross-section σ of the trapping centres and the number of vacant trapping centres. The electron flux is Jv_{th}/qv_d , where J is the current density and q the electronic charge, v_{th} is the mean velocity of hot electrons (in all directions) and v_d the effective drift velocity in the direction of the field gradient. The number of vacant trapping centres is equal to the total number of centres (N) less the number of centres already occupied (n). The trapping rate dn/dt can be expressed as:

$$dn/dt = J(N - n)/Q_0, \quad (1)$$

where Q_0 is $qv_d/\sigma v_{th}$. Integration, for a constant value of J , gives:

$$n = N\{1 - \exp(-Jt/Q_0)\}. \quad (2)$$

If this relation applies, a plot of $\log(dn/dt)$ against t should yield a straight line. In practice, however, this is not always so; see for example *fig. 2*. In addition, eq. (2) indicates that n should rapidly saturate, since the exponential term rapidly goes to zero when $Jt > Q_0$. As a rule, however, no such saturation is found.

To obtain better agreement between theory and experiment the presence of centres of many different kinds has been postulated in the literature [6] [7]. It will now be shown, however, that a good description

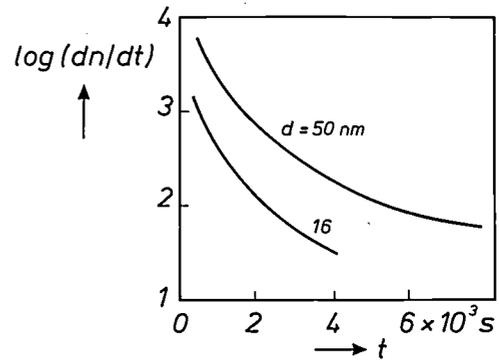


Fig. 2. Logarithm of the trapping rate dn/dt for charge carriers in SiO_2 as a function of the duration t of the applied voltage, for two thicknesses d of the oxide film. The data are derived from results reported by D. R. Young [6]. The first-order approximation, which predicts the straight lines given by eq. (2), is clearly not applicable here.

of charge trapping can be given without any such assumption, provided sufficient account is taken of the repulsion experienced by injected charges from charges already trapped [8] [9].

Effects of Coulomb repulsion

Owing to the Coulomb repulsion between injected charges and charges already trapped, the trapping of new charges becomes less likely in the vicinity of centres already occupied. The trapping probability will decrease as the number of occupied centres increases. For a simple calculation of this diminished probability it may be assumed that an occupied centre makes a region of volume h inactive for further trapping. The volume in which trapping can take place is then no longer V , the total volume of the oxide film, but $V - h$. The trapping probability is thus reduced by a factor of $(1 - h/V)$. If n charges have been trapped, the probability of further trapping is reduced by $(1 - h/V)^n$. The trapping is reduced by this factor and not by $(1 - nh/V)$ because the various inactive regions can overlap. Since $h \ll V$, a good approximation is:

$$(1 - h/V)^n = \exp(-nh/V). \quad (3)$$

This exponential factor must be included in the expression for the trapping rate (eq. (1)), giving:

$$dn/dt = \frac{J(N - n)}{Q_0} \exp(-nh/V). \quad (4)$$

For $n \ll N$ we have, after integration:

$$\exp(nh/V) - 1 = \frac{JNht}{Q_0V}. \quad (5)$$

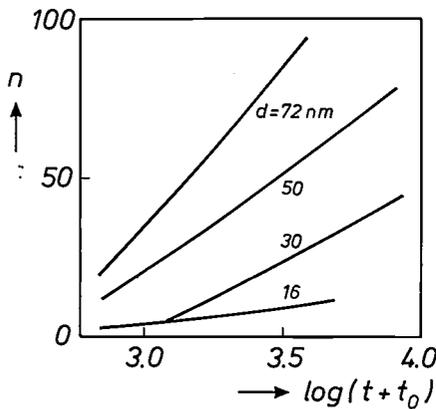


Fig. 3. Number of trapped charge carriers (n) in arbitrary units plotted against $\log(t + t_0)$, see text, for different thicknesses d of the SiO_2 film. These data are also derived from Young's results [6]. In all cases a straight line is found, as predicted by eq. (6).

This leads to

$$n = \frac{V}{h} \ln \frac{t+t_0}{t_0}, \tag{6}$$

where $t_0 = Q_0 V / J N h$. This indicates that the plot of n against $\ln(t + t_0)$ should be a straight line [8]. This is indeed found experimentally, as appears from fig. 3 for oxide films of different thickness.

This model also gives a better description of the variation in the current with the amount of transferred charge. In general the current density J depends on the electrical field-strength E_i at the interface with the gate. To a good approximation the simple exponential relation

$$J = J_0 \exp(E_i/E_0) \tag{7}$$

can be used, where J_0 and E_0 are constants. The value of E_i is determined by the strength E of the applied field and by the field-strength induced by the trapped charges:

$$E_i = E - nq\bar{x}/\epsilon, \tag{8}$$

where \bar{x} is the mean distance of the trapped charges from the electrode and ϵ is the permittivity of the film. The number of trapped charges (n) should now be expressed in terms of the charge transferred per unit area, defined as:

$$Q = \int_0^t J dt. \tag{9}$$

For the first-order approximation we can write, by analogy with eq. (2):

$$n = N \{1 - \exp(-Q/Q_0)\}. \tag{10}$$

For $Q \leq Q_0$, a good approximation is:

$$n = N \ln \left(\frac{Q}{Q_0} + 1 \right). \tag{11}$$

Substituting eq. (8) and eq. (11) in eq. (7) shows that the curve of $\ln J$ against $\ln Q$ is horizontal at first for small Q . Then $\ln J$ decreases linearly with a slope that is proportional to N . Beyond a saturation value (Q_{sat}) the curve should again be horizontal: saturation occurs because n has become equal to N and no more new charges can be trapped. The value of Q_{sat} can be derived from eq. (11) by putting $n = N$. This gives $Q_{\text{sat}} = Q_0 (e - 1) \approx 1.7 Q_0$. The curve of $\ln J$ as a function of $\ln Q$ is shown schematically in fig. 4a.

When the Coulomb repulsion is taken into account, analogy with eq. (6) gives the following expression for n :

$$n = \frac{V}{h} \ln \left(\frac{Q}{Q_1} + 1 \right), \tag{12}$$

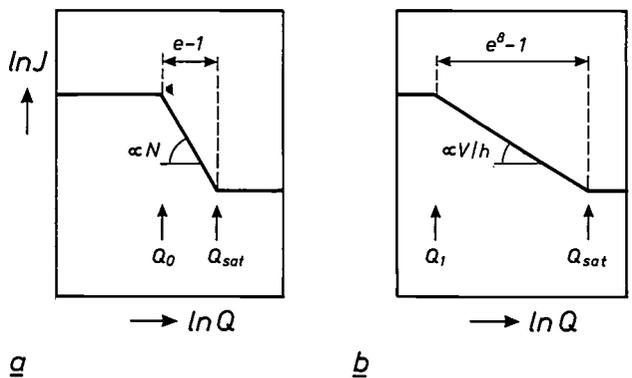


Fig. 4. Schematic curve of $\ln J$ as a function of $\ln Q$ as given by the first-order approximation (a) and after taking into account the Coulomb repulsion between injected electrons and electrons already trapped (b); see text. The main difference is that the sloping section in (a) is much steeper and the saturation point (Q_{sat}) is reached much earlier than in (b).

[6] The applicability of this approximation to earlier experiments is to some extent demonstrated in: E. H. Nicollian, A. Goetzberger and C. N. Berglund, Avalanche injection currents and charging phenomena in thermal SiO_2 , Appl. Phys. Lett. 15, 174-177, 1969; E. H. Nicollian, C. N. Berglund, P. F. Schmidt and J. M. Andrews, Electrochemical charging of thermal SiO_2 films by injected electron currents, J. Appl. Phys. 42, 5654-5664, 1971.
 [6] D. R. Young, Electron trapping in SiO_2 , Inst. Phys. Conf. Ser. 50, 28-39, 1980.
 [7] See also: T. H. Ning and H. N. Yu, Optically induced injection of hot electrons into SiO_2 , J. Appl. Phys. 45, 5373-5378, 1974; R. A. Gdula, The effects of processing on hot electron trapping in SiO_2 , J. Electrochem. Soc. 123, 42-47, 1976; D. R. Young, E. A. Irene, D. J. Di Maria, R. F. De Keersmaecker and H. Z. Massoud, Electron trapping in SiO_2 at 295 and 77 °K, J. Appl. Phys. 50, 6366-6372, 1979.
 [8] D. R. Wolters and J. F. Verwey, Trapping characteristics in SiO_2 , in: Insulating Films on Semiconductors, M. Schulz and G. Pensl (eds), Series in Electrophysics, Vol. 7, Springer, Berlin 1981, pp. 111-117.
 [9] D. R. Wolters and J. J. van der Schoot, Kinetics of charge trapping in dielectrics, J. Appl. Phys. 58, 831-837, 1985.

where $Q_1 = Vqva/hN\sigma v_{th}$. At first sight this expression resembles eq. (11). If, after substituting eqs (8) and (12) in eq. (7), we plot the value of $\ln J$ against the value of $\ln Q$, we again obtain a curve with two horizontal sections, linked by a sloping section; see fig. 4b.

However, two important differences should be noted. The first concerns the slope of the sloping section: this is now much smaller because it is not proportional to N but to the much smaller value of V/h . The second difference concerns the saturation point Q_{sat} . In the model with the Coulomb repulsion, saturation does not occur until the number of trapped charges is equal to the maximum number of spheres of volume h that can be accommodated in the total volume V . If the spheres do not overlap, this number is V/h . If the spheres overlap to the maximum possible extent, this number becomes $8V/h$, which is equal to

the number for the maximum close-packing (with no overlapping) of spheres of half the radius. From eq. (12) this means that $\ln\{(Q_{sat}/Q_1) + 1\} = 8$, so that $Q_{sat} = (e^8 - 1)Q_1 \approx 2980 Q_1$. The value of Q_{sat} is thus more than three orders of magnitude larger than the characteristic value Q_1 , whereas in the first-order approximation Q_{sat} is not much larger than the characteristic value Q_0 for that case. Such a high value of Q_{sat} cannot be reached experimentally because other effects (e.g. breakdown) can occur beforehand. When the Coulomb repulsion is taken into account it can therefore be seen why no saturation is found in the range in which measurements can be made.

Results of measurements of J as a function of Q for different applied field-strengths are shown in fig. 5. The measured curves agree well with the expression that can be derived from equations (7), (8) and (12):

$$\ln(J/J_0) = \frac{E}{E_0} - \frac{Vq\bar{x}}{h\epsilon E_0} \ln\left(\frac{Q}{Q_1} + 1\right). \quad (13)$$

The curves have two asymptotes. At a low value of Q , J is constant; from the values of J/J_0 and the external field-strength E we can then derive E_0 . At a high value of Q , $\ln J$ decreases linearly with $\ln Q$.

Eq. (13) shows that the slope of the asymptote at high Q is $Vq\bar{x}/h\epsilon E_0$. Since q , ϵ and E_0 are known to a first approximation and \bar{x} can be equated to half the thickness of the oxide film, the slope will give the value of V/h , which can be seen as the effective density N_{eff} of the trapping centres. It has been found that N_{eff} does not depend on the applied field; see fig. 6. The experimental curves of $\ln J$ versus $\ln Q$ can also be used for determining Q_1 . The value of the effective trapping cross-section σ_{eff} can then be derived from Q_1 :

$$\sigma_{eff} = q/Q_1. \quad (14)$$

Fig. 6 shows that with increasing field-strength, the value of σ_{eff} decreases by more than an order of magnitude. At high field-strengths, σ_{eff} is virtually constant.

We believe that the results described here will contribute to a much better understanding of charge trapping. The knowledge gained can help us to suppress the undesired effects resulting from charge trapping in the silicon, or to control them better. The results have also put us on the track of the mechanism that plays a major part in the breakdown of the film. We shall now look at some aspects of electrical breakdown.

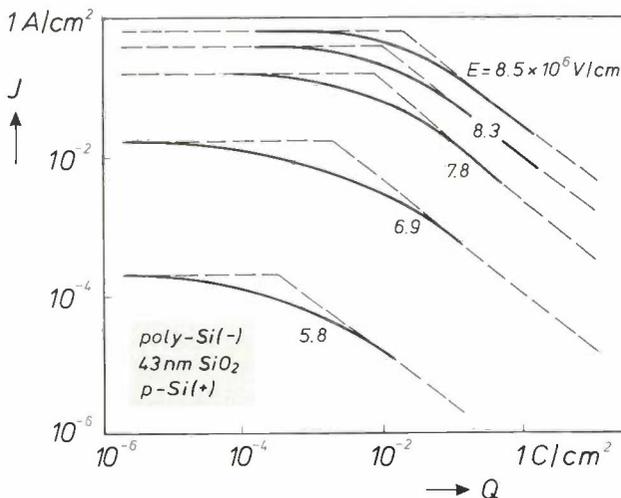


Fig. 5. Log-log plot of the current density J measured in a MOS device as a function of the charge transferred per unit area (Q) for different field-strengths E . As eq. (13) indicates curves with two asymptotes (dashed lines) are obtained: $\ln J$ is constant at low Q and decreases linearly with $\ln Q$ at high Q .

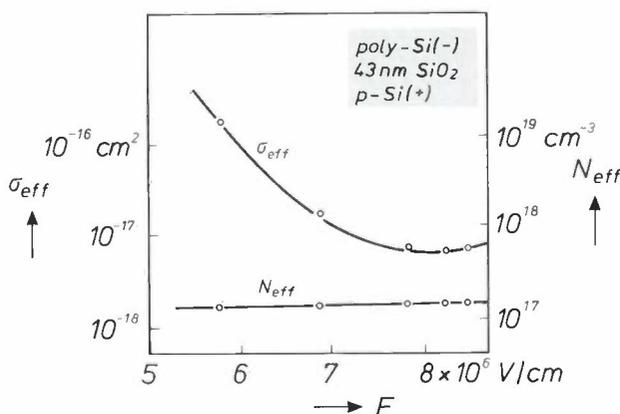


Fig. 6. Effective cross-section σ_{eff} for electron trapping in the SiO_2 film of the MOS device in fig. 5, and the effective density N_{eff} of the trapping centres, as a function of the applied field-strength E . Whereas N_{eff} is independent of E , the value of σ_{eff} first decreases with E until a virtually constant value is reached at high E .

[10] D. R. Wolters and J. J. van der Schoot, Dielectric breakdown in MOS devices, Part I: Defect-related and intrinsic breakdown, Philips J. Res. 40, 115-136, 1985.

Electrical breakdown

When the voltage across the oxide film in a MOS capacitor is increased, a value is eventually reached at which the field-strength in the film becomes so great that it breaks down^[10]. The breakdown takes the form of a sudden current surge accompanied by destruction of the film. This effect is also found at a lower constant voltage if the voltage is applied for long enough. Fig. 7a gives the current density as a function of time measured at various field-strengths for a MOS capacitor with a polysilicon gate. The current density at first gradually decreases, then after a time suddenly increases, indicating that the oxide film has broken down. The breakdown is very rapid; it is complete in less than 100 ns.

The effect of breakdown on the device on a microscopic scale is shown in fig. 7b. At the place where the

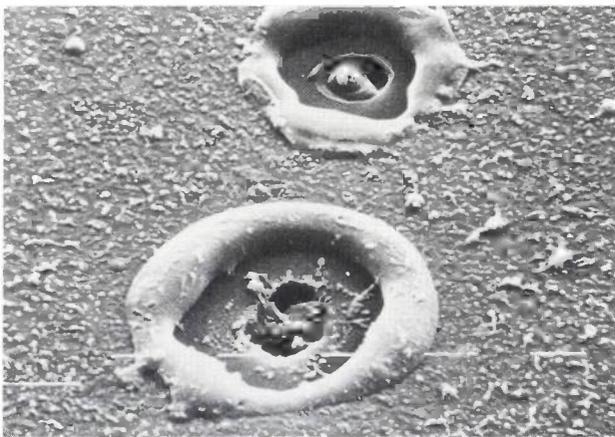
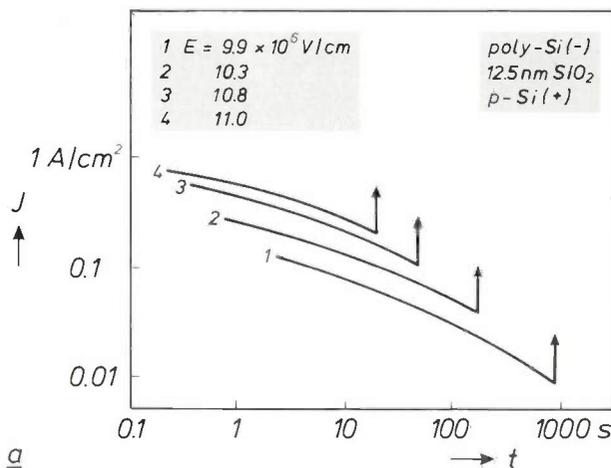


Fig. 7. a) Current density J as a function of time t in a MOS capacitor for four field-strengths E . After a slow decrease, J suddenly rises sharply and the film breaks down. The time to breakdown decreases as the field-strength is increased. b) Scanning electron micrograph (each dash represents $1 \mu\text{m}$) of a failed MOS capacitor consisting of n-doped silicon, a 40-nm thick SiO_2 film and a polysilicon gate. The gate voltage was negative and the area of the capacitor was 0.2 mm^2 . At breakdown a crater was formed that penetrated right through into the silicon. After the first breakdown a voltage was applied to the capacitor again until it broke down once more, and a second crater was formed.

film broke down the polysilicon has melted and a crater has formed with a diameter of 4 to 5 μm . The outer ring of the crater consists of molten polysilicon, and the hole in the centre extends into the silicon substrate. The polysilicon often draws back from the centre of the crater, probably because of the surface tension of the molten gate material. When this happens, the oxide film continues to act as an insulator and a voltage can again be applied to the device until there is a second breakdown. Then a second crater forms, which resembles the first one (fig. 7b). The occurrence of holes in the silicon substrate with a diameter of 1.5 μm indicates that the energy density at breakdown is very high.

The effect of breakdown at constant current is shown in fig. 8 for four very different current values. The dimensions of the craters appear to be virtually independent of the current. It can also be seen that the craters generally form at the corners and edges of the device. This is not so surprising, since these are the places where the highest field-strength would be expected.

Defect-related and intrinsic breakdown

The breakdown of a MOS capacitor often takes place at a weak spot, a defect, in the oxide film. The defect may have been caused during manufacture by the incorporation of impurities or by the formation of hair-line cracks, asperities or pinholes. Defects may also be introduced when the gate is deposited. Breakdown at a defect will occur at a lower field-strength than intrinsic breakdown, i.e. breakdown for a film with no defects. If there are many defects, the breakdown will occur at the weakest spot. In general, therefore, the breakdown field-strength will decrease as the number of defects in the film increases.

Defect-related breakdown can clearly be distinguished from intrinsic breakdown in a statistical analysis of the results of tests on a large number of capacitors^[10]. This analysis uses the same sort of probability calculation as would be used for the breaking of chains at the weakest link. A test capacitor, considered as an imaginary row of capacitors connected in parallel, will fail when the weakest capacitor in the row breaks down. The percentage of failed capacitors (F) is plotted on probability paper. This means that $\ln\{\ln(1 - F)^{-1}\}$ instead of F is plotted as a function of a variable parameter, such as the field-strength E . The advantage of this is that it is easier to make a distinction between defect-related breakdown and intrinsic breakdown. It also offers advantages for predicting the behaviour of electronic devices from the results of measurements on smaller or less complicated devices. For example, if we have a large number of capacitors

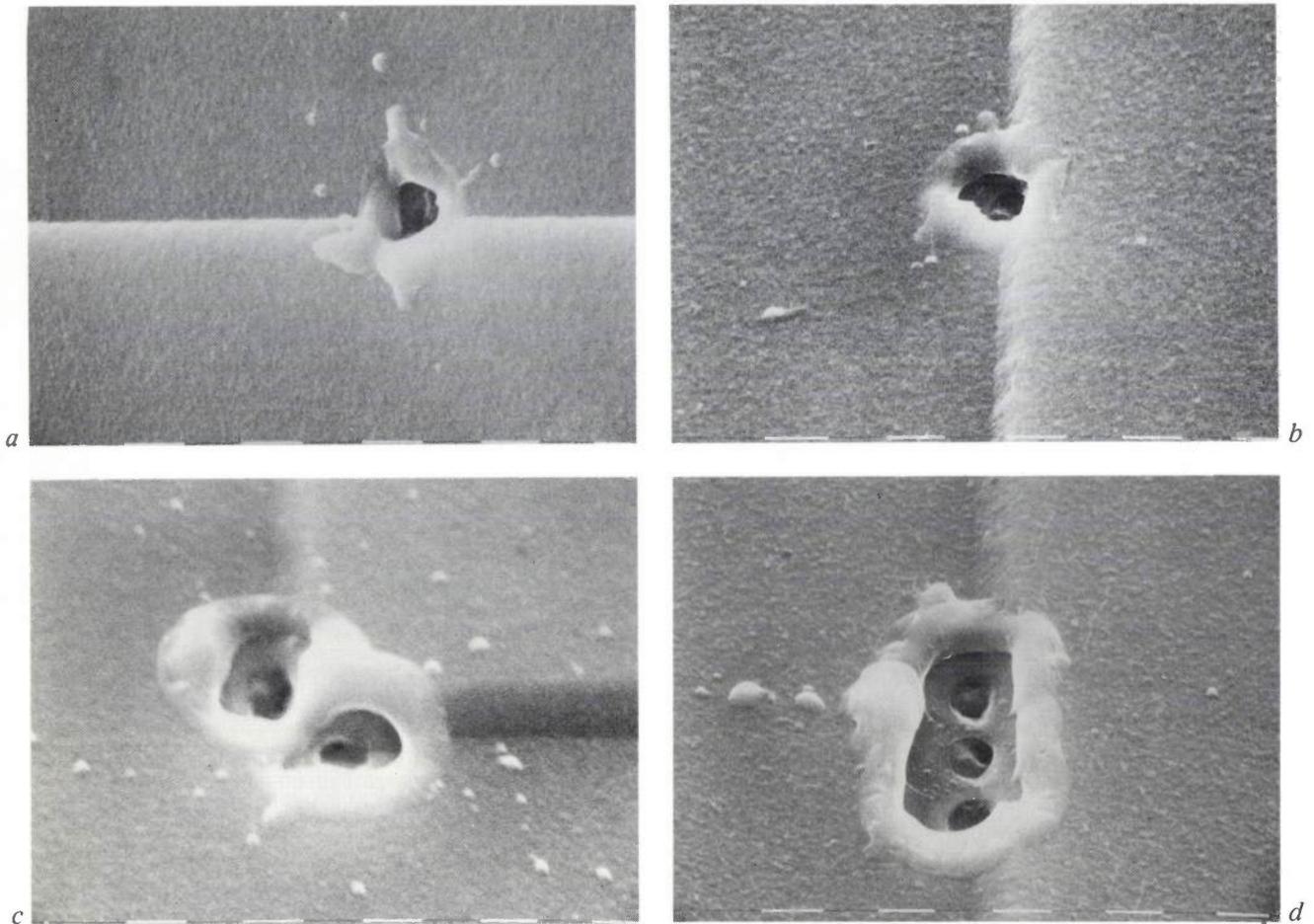


Fig. 8. Scanning electron micrographs (scale division 1 μm) of craters formed during breakdown of MOS capacitors for different current values: a) 10 nA. b) 100 nA. c) 1 μA . d) 10 μA . The dimensions of the craters are virtually independent of the current. The craters usually form at corners and edges of the device.

containing defects and connected in parallel, then to a good approximation [11]:

$$\ln\{\ln(1-F)^{-1}\} = c_1 E + \ln A + c_2, \quad (15)$$

where A is the area of the individual capacitors, and c_1 and c_2 depend on the distribution of the defects and the number of observations.

Fig. 9 gives a plot of $\ln\{\ln(1-F)^{-1}\}$ as a function of the field-strength E for 12 000 capacitors, each with an area of 0.02 mm^2 . On increasing the field-strength a gradual rise is first observed, followed by a very steep increase at $E \approx 9 \times 10^6$ V/cm, where F becomes almost 100%. From the variation of F with E it may be deduced that in about 7% of cases the breakdown is defect-related. The other 93% fail at practically the same field-strength, and we can therefore assume that these breakdowns are intrinsic.

If we now consider the capacitors in groups of say 400, we find that breakdown in a group occurs as soon

as the weakest capacitor fails. The probability of breakdown will increase to the same extent as it would if the area were increased by a factor of 400. The results for the breakdown distribution are also given in fig. 9. As eq. (15) predicts, a straight line is obtained, shifted in the vertical direction by $\ln 400$ with respect to the curve for the single capacitors.

The defect density D can be derived from the point where there is a sudden increase in the slope of the curve for the single capacitors. When this point has been reached all the capacitors with defect-related breakdown have been detected. If the defects are randomly distributed over the surface, a Poisson distribution may be assumed. The probability F_k that one or more defects will be present in an area A is then given by:

$$F_k = \sum_{m=1}^{\infty} \frac{(AD)^m \exp(-AD)}{m!}, \quad (16)$$

The probability of finding capacitors without any defects is then:

$$1 - F_k = \frac{(AD)^0 \exp(-AD)}{0!} = \exp(-AD). \quad (17)$$

Taking $F_k = 7\%$ and $A = 0.02 \text{ mm}^2$, it follows from this that $D = 3.7 \times 10^2 \text{ cm}^{-2}$.

It is also possible to demonstrate indirectly that the failures to the left of the sharp increase in the slope of the curve are due entirely to defects. To demonstrate this, capacitors were made with a much lower defect density than the value given above. Fig. 10 shows that in this case all 16000 capacitors tested failed at vir-

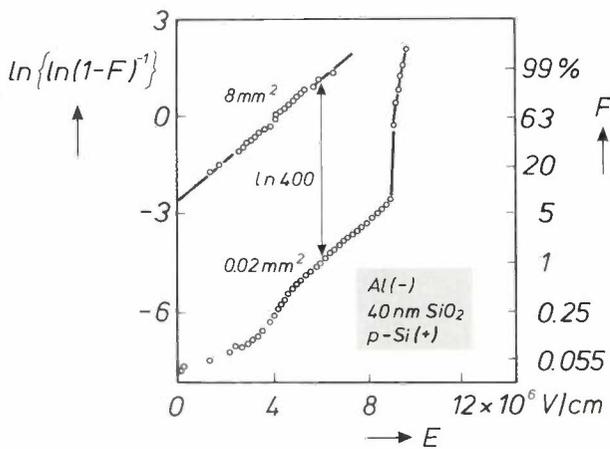


Fig. 9. Lower curve: Breakdown distribution of 12000 MOS capacitors of area 0.02 mm^2 . The value of $\ln\{\ln(1 - F)^{-1}\}$ is plotted as a function of the field-strength E ; F is the percentage of failed capacitors. There is a sharp increase in the curve at about $9 \times 10^6 \text{ V/cm}$, where F is about 7%. Above this value there is a very steep increase in F and hence in $\ln\{\ln(1 - F)^{-1}\}$. Upper curve: Breakdown distribution for groups of 400 capacitors, with a combined area of 8 mm^2 . As eq. (15) predicts, the points lie on a straight line shifted by $\ln 400$ with respect to the lower curve.

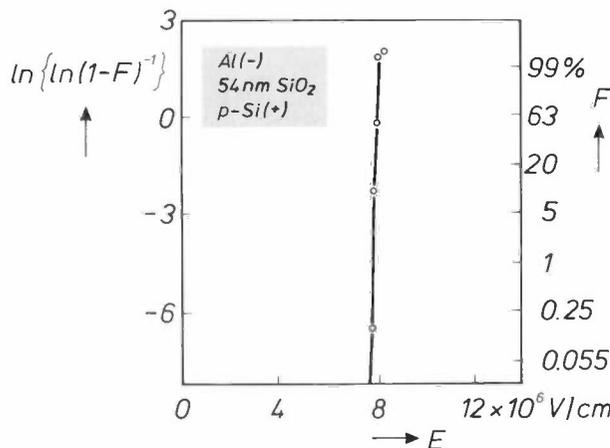


Fig. 10. Breakdown distribution for 16000 capacitors, each with an area of 0.02 mm^2 . The capacitors contain so few defects that they only give intrinsic breakdown, and all of them fail at virtually the same field-strength.

tually the same field-strength. This points to intrinsic breakdown as the sole reason for failure. It can be deduced from the absence of defect-related breakdown that the defect density here must be less than 0.3 cm^{-2} .

Field-dependent and time-dependent breakdown

Instead of testing with increasing field-strength, we can also test the capacitors by applying a constant field and then measuring the time to the moment of breakdown, as shown in fig. 7a. This was done for four groups of 900 capacitors on a slice on which 16000 capacitors had first been tested with a varying electric field. These had a high defect density because of the incorporation of large amounts of impurities. Fig. 11a shows the breakdown distribution with increasing field-strength. The transition to intrinsic breakdown is found at $13 \times 10^6 \text{ V/cm}$, where about half of the capacitors fail. On the basis of this distri-

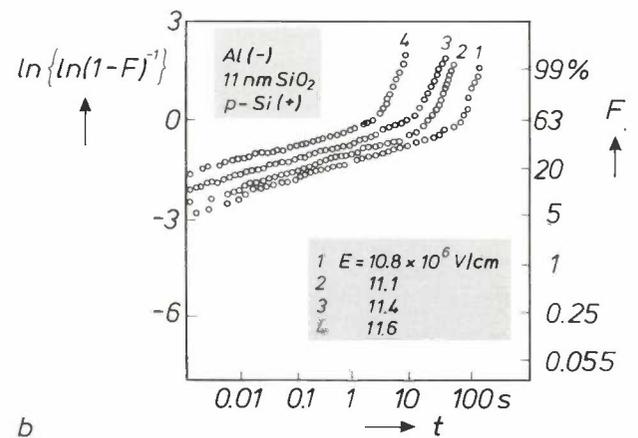
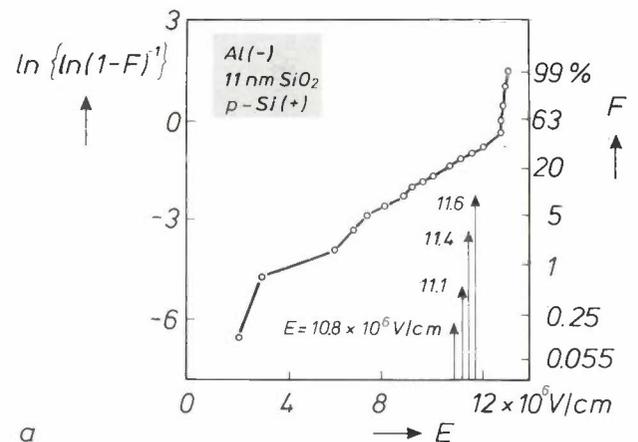


Fig. 11. a) Breakdown distribution for MOS capacitors with a high defect density. The arrows indicate the values of the field E at which time-dependent tests were carried out. b) Breakdown distribution in time-dependent tests on capacitors of (a) for four different field-strengths E . The sudden increases in the curves occur at about the same value of F .

[11] E. J. Gumbel, Statistics of extremes, Columbia Univ. Press, New York 1958.

bution four groups of 'virgin' capacitors were tested at constant field, for four values of E that were significantly lower than the value for intrinsic breakdown. Fig. 11b shows the result in breakdown distributions as a function of time. The curves are not only very similar but they also resemble the distribution in fig. 11a.

The defect densities that can be calculated from the various changes in the slope of the curves are almost identical. A similar correspondence between defect densities determined from field-dependent and time-dependent tests is also found for slices with fewer defects. These results indicate that capacitors that fail because of defect-related breakdown at a relatively weak field will also fail in a constant field after a relatively short time (and vice versa).

Other methods for applying the field to the capacitor can be compared in the same way. Fig. 12 gives the distribution for constant field and for constant current. For convenience the breakdown distribution is shown in both cases as a function of the integrated charge transfer per unit area (Q_{bd}). The curves are virtually coincident above a particular value of Q_{bd} , before the change in slope. In short, the way in which the field is applied for breakdown tests on MOS capacitors is not important.

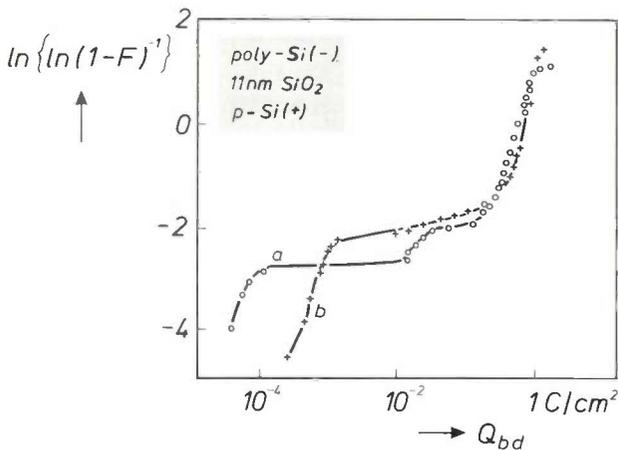


Fig. 12. Breakdown distribution for MOS capacitors as a function of the integrated charge transferred per unit area (Q_{bd}) at a constant current of $40 \mu\text{A}$ (curve *a*) and at a constant voltage of 12.9 V (curve *b*). Beyond Q_{bd} values greater than 0.01 C/cm^2 the curves are virtually coincident.

Principal parameter for intrinsic breakdown

The tests on MOS capacitors indicate that the circumstances that lead to intrinsic breakdown do not depend essentially on the method of application of the field. This can be seen if we consider the time that elapses before breakdown [12]. The measured points in fig. 11a related to measurements in which the field-

strength was increased by $2 \times 10^7 \text{ V/cm}$ per second. An increase in the field-strength from $11 \times 10^6 \text{ V/cm}$ to $13 \times 10^6 \text{ V/cm}$ thus takes about 0.1 s . If the curves in fig. 11b are extrapolated to a curve for $13 \times 10^6 \text{ V/cm}$, the time obtained has about this value.

The simplest method of testing is to use a constant current, as in curve *a* of fig. 12. In this method the field-strength and the time are measured, and the charge transfer at breakdown is proportional to time. We used this method to test a large number of capacitors in which there were differences in gate material, the thickness of the oxide film (between 10 and 50 nm) and the manufacturing process. The current density was varied and only the intrinsic breakdown was studied. Fig. 13 gives a plot of the measured time to breakdown (t_{bd}) as a function of current density J . For not too high values of J , $\log t_{bd}$ is given to a good approximation by:

$$\log t_{bd} = -\log J + c, \quad (18)$$

where the constant c is determined by the characteristics of the capacitor. The total charge transferred per unit area (Q_{bd}) is given by:

$$Q_{bd} = J t_{bd}. \quad (19)$$

Equations (18) and (19) indicate that Q_{bd} is determined entirely by the capacitor and does not depend on the conditions in which the breakdown is produced. This indicates that Q_{bd} is the quantity that most affects intrinsic breakdown. In other words: a capacitor fails as soon as a critical amount of charge has

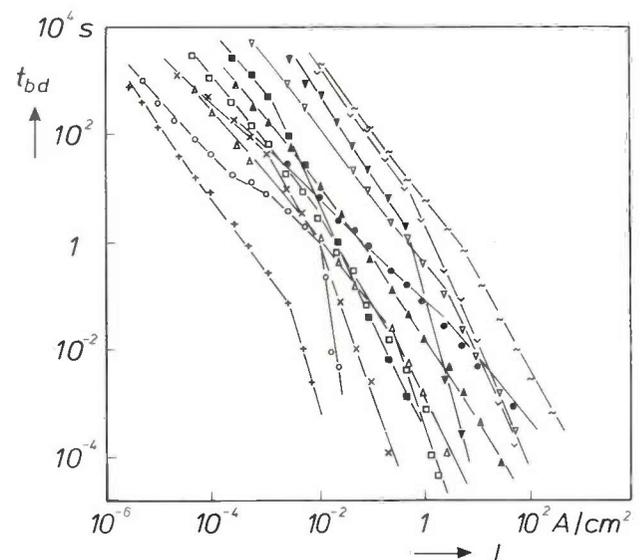


Fig. 13. Log-log plot of the time to breakdown (t_{bd}) against current density J , for MOS capacitors with differences in the gate (aluminium or polysilicon), in the thickness of the SiO_2 film (between 10 and 50 nm) and in the manufacturing process. If the current densities are not too high, straight lines with a slope of about -1 are obtained.

been transferred through the film, however this transfer was produced.

The value of Q_{bd} varies little in a wide range of current densities. At very high current densities, above a critical value J_{cr} , there is often a sharp decrease in the value of Q_{bd} ; this decrease also depends on the direction of the field; see fig. 14a. If the gate voltage is negative, J_{cr} is much lower (by a factor of ten) and the decrease in Q_{bd} is steeper than with a positive gate voltage. The fall in Q_{bd} beyond J_{cr} is closely connected with the change in the gate voltage (V_{bd}) necessary to keep the current constant; see fig. 14b. At moderate current densities, V_{bd} increases linearly with $\log J$. Beyond J_{cr} , however, there is either a steep rise in V_{bd} (at a negative gate voltage) or a steep fall (at a positive gate voltage). These effects are attributed to the accumulation of space charge in the film, which is dependent on the direction of the field. This will be discussed later.

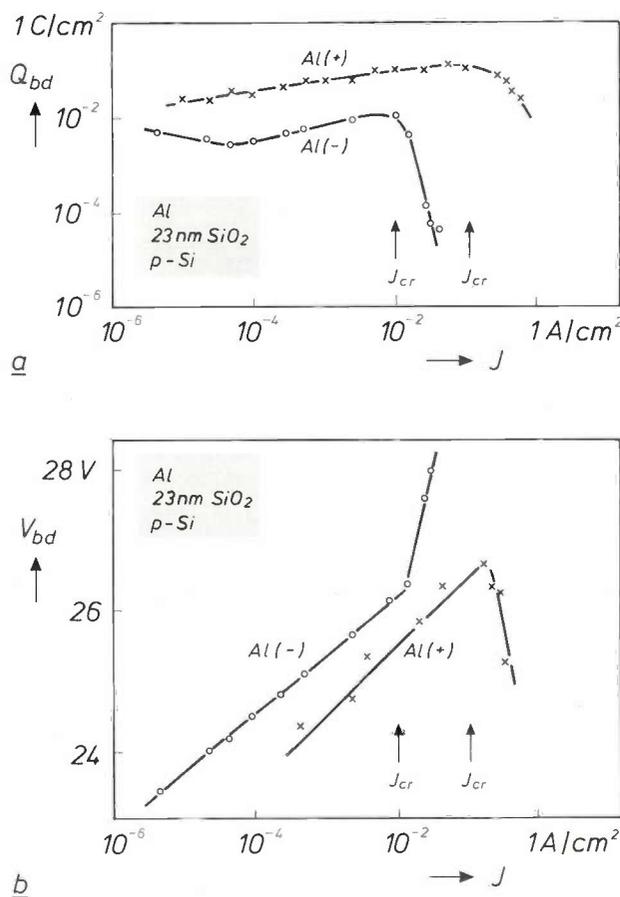


Fig. 14. a) Amount of charge transferred per unit area before breakdown (Q_{bd}) as a function of current density J for MOS capacitors with a positive or a negative voltage on the aluminium gate. Below a critical value J_{cr} the current density has little effect on Q_{bd} , but above it there is a decrease that depends strongly on the direction of the field. b) Breakdown voltage V_{bd} as a function of current density J for these two cases. Beyond $J = J_{cr}$ the value of V_{bd} increases sharply when the gate voltage is negative and decreases sharply when the gate voltage is positive.

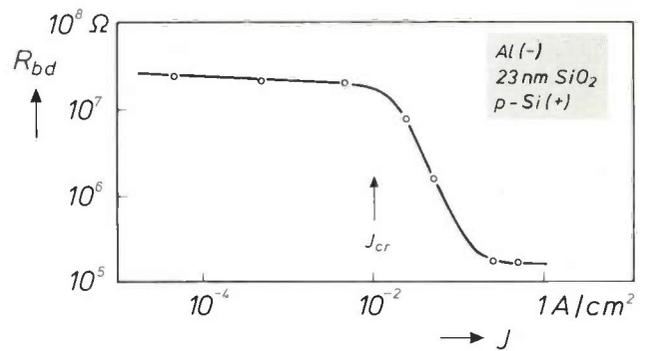


Fig. 15. Resistance of a MOS capacitor after breakdown (R_{bd}) as a function of current density J at breakdown. Beyond $J = J_{cr}$ the value of R_{bd} decreases sharply with J .

Another result that is closely related to the decrease of Q_{bd} at high current densities concerns the resistance of the capacitor after breakdown (R_{bd}). Fig. 15 shows R_{bd} as a function of the current density at breakdown. Beyond J_{cr} the value of R_{bd} falls by two orders of magnitude. Clearly, a breakdown at $J > J_{cr}$ causes a different kind of damage from that at $J < J_{cr}$.

From the results in figs 14 and 15 we can conclude that different mechanisms enter into the breakdown on opposed sides of $J = J_{cr}$. A model will now be presented that gives a good explanation of the behaviour at both low and high current densities.

Mechanism for wear and intrinsic breakdown

The breakdown of a MOS capacitor with no defects may be regarded as a two-stage process. The leakage of charge first causes structural changes in the film that ultimately result in the formation of a low-resistance path between the electrodes. The capacitor can discharge along this path. The energy released causes permanent damage to the oxide film, as can be seen in figs 7b and 8. The leaking charge wears out the film, so that it gradually degenerates and finally breaks down.

This model explains the test results described above fairly well, especially the behaviour at moderate current densities. The greater the leakage in the oxide film, the sooner it will fail. How the leakage arises is immaterial. The capacitor fails as soon as a particular amount of charge has been transferred. We shall now briefly consider why this is so.

Breakdown at constant Q_{bd}

From the observation that Q_{bd} does not vary much at moderate currents and fields it may be deduced that

[12] D. R. Wolters and J. J. van der Schoot, Dielectric breakdown in MOS devices, Part II: Conditions for the intrinsic breakdown, Philips J. Res. 40, 137-163, 1985.

collision ionization^[13] and ion migration^[14] play little part in the breakdown. Mechanisms based on the kinetic energy of the charge carriers predict that Q_{bd} will be closely dependent on the injection rate and the field-strength. The dependence on the previous history of the device excludes mechanisms with an avalanche type of breakdown.

So how can we explain the constancy of Q_{bd} while accepting that the total of dissipated energy is proportional to the applied voltage? This question can only be answered if the kinetic energy does *not* contribute to the breakdown. To understand this, we have to consider the energy of electrons injected into the SiO₂ film; see *fig. 16*. Electrons that tunnel through the po-

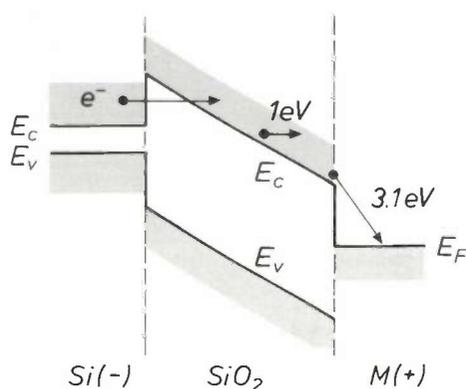


Fig. 16. Energy diagram for electron transfer in a MOS capacitor. E_v upper edge of the valence band. E_c lower edge of the conduction band. E_F Fermi level of the metal. The energy gain for an electron is much greater in the transition to the metal M than in migration in the conduction band of SiO₂.

tential barrier from the silicon into the SiO₂ film have a high potential energy in the conduction band. Provided the values of the electric field-strength are not too high, the kinetic energy of the electrons in the conduction band will be low because of interactions with the vibrating lattice atoms (phonon interactions) in the film. Since the interaction length is no more than 2 to 4 nm, the energy gain in a field of 5×10^6 V/cm will be no greater than 1 to 2 eV. The number of electrons with this gain in energy will also be very small. As the minimum energy required to break a silicon-oxygen bond is about 4 eV, it follows that the kinetic energy of the electrons in the conduction band is insufficient to cause any significant damage.

The energy dissipated by the electrons becomes considerable, however, if they are emitted into the gate from the SiO₂ conduction band or if they are trapped in centres in the SiO₂ film. The high energy, the sum of the potential energy in the conduction band and the acquired kinetic energy, is dissipated in a distance that is short compared with the thickness of the film. Consequently the density of the energy dissi-

pation is much higher than that of the kinetic energy dissipation in the conduction band, which is distributed throughout the entire volume. Since all the injected electrons have a high potential energy that they must dissipate on leaving the SiO₂ conduction band, the dissipation that ultimately leads to breakdown is entirely determined by the number of injected electrons. The damage produced is permanent, so that the effect is cumulative.

Behaviour at high current densities

As the field-strength or the current density increases the injected electrons acquire more and more kinetic energy, while the potential energy remains virtually constant. A field-strength could then be reached for which Q_{bd} was field-dependent. This might explain the decrease in Q_{bd} at $J > J_{cr}$ (*fig. 14a*), except that at a positive gate voltage the injection voltage at breakdown (V_{bd}) decreases instead of increasing with J (*fig. 14b*).

A better explanation can be obtained if we assume that the trapping centres in the film are regularly filled and emptied while a field is applied to the film. At low current densities ($J < J_{cr}$) most of the centres will not on average be occupied. Space-charge effects will then be negligible, and the injection of electrons will be homogeneous over the entire area of the electrode. At high current densities ($J > J_{cr}$) the mean occupancy of the trapping centres can be more than 50%^[12], so that space-charge effects will be significant. This can be seen in the effect that the direction of the field has on the injection voltage at breakdown (*fig. 14b*). If there are trapping centres close to an aluminium gate, then at a negative gate voltage the injection of electrons will be opposed by the formation of a negative space charge. This means that V_{bd} becomes higher. At a positive gate voltage a positive space charge is formed. This stimulates the injection of electrons from the silicon, so that V_{bd} becomes lower.

For both directions of the field the injection will be inhomogeneous at high current densities. When the gate voltage is negative the injection will take place mainly in regions of low space-charge density. The current density will then be very high locally, so that at such places Q_{bd} will be reached while the amount of injected charge is still relatively small. When the gate voltage is positive, the injection will take place mainly in regions where the space-charge density is high. But this also means that J will be higher locally and that the value of Q_{bd} will be reached earlier.

Confirmation of the model outlined here is found from measurements of Q_{bd} for capacitors with different gates and different thicknesses of the oxide film; see *fig. 17*. In the vicinity of an aluminium gate

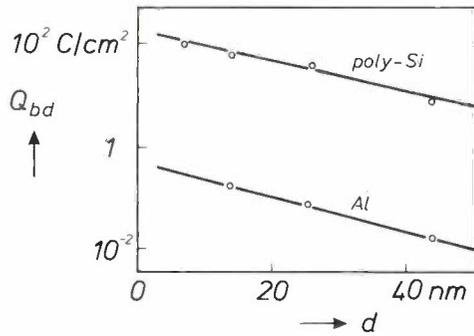


Fig. 17. Amount of charge transferred per unit area up to breakdown (Q_{bd}) plotted as a function of the thickness d of the oxide film for MOS capacitors with a polysilicon and an aluminium gate. In polysilicon, Q_{bd} is a hundred times higher. In both cases, $\log Q_{bd}$ decreases linearly with d .

more centres will be formed than in the vicinity of a polysilicon gate. This means that the space-charge effects will be more pronounced with aluminium, and the film will therefore break down at a lower value of Q . In general, Q_{bd} decreases exponentially with film thickness. In thicker films the trapping probability is higher, so that the space-charge effects will be more marked.

Wear before breakdown

There are now sufficient indications that charge leakage causes wear. For a given amount of transferred charge, for example, breakdown will occur irrespective of the number of current interruptions. Cumulative damage of this nature can be considered to be the result of the formation of a discharge pattern due to the injection of electrons^{[15] [16]}. This pattern has a tree structure, as illustrated in *fig. 18*. As the damage increases, the structure spreads, and the tree acquires more and thicker branches. It appears that the injected electrons by preference lose their potential energy at places that have been damaged before. The growth of the discharge pattern continues until it forms a low-resistance path between the two electrodes, and the capacitor discharges along this path.

This picture is rather like the pattern produced by rivers eroding a landscape. The mechanism for the formation of a discharge pattern can be compared with the erosion caused by the flowing water. That erosion is also cumulative: the wear is not related to the rate of flow but to the total volume of the flowing water. Just as small streams eventually form a broad river and leave deep traces in the landscape, so the leaking charge wears away the insulating film and finally destroys it.

Both the damage and the breakdown require a certain amount of energy. A simple representation of the

energy content of a charged capacitor is given in *fig. 19* in the form of a plot of the applied voltage V as a function of the charge supplied Q . The area under the curve is proportional to the amount of energy supplied. As long as no charge leaks through the dielectric, the stored charge is proportional to V . The proportionality factor is the capacitance C , and the stored energy is $\frac{1}{2}CV^2$. Above a particular voltage (V_1) the curve departs from a straight line as a result of charge leaking through the SiO_2 film. At a voltage V_2 the stored energy is given by $\frac{1}{2}CV_2^2$; the remainder of the area under the curve up to $Q = Q_2$ gives the energy that has been dissipated in the dielectric and has been

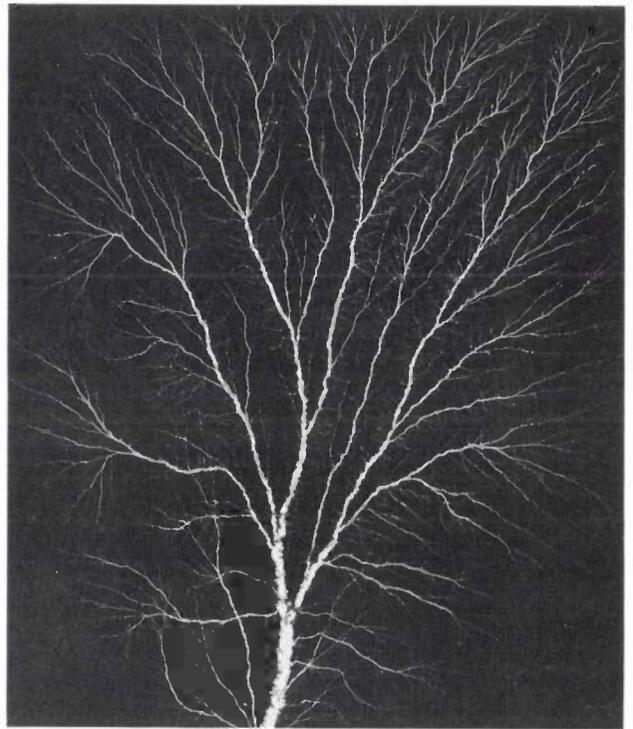


Fig. 18. Tree-shaped discharge pattern after injection of 3-MeV electrons into an insulator that was then earthed at a central point^[16].

[13] See the book by E. H. Nicollian and J. R. Brews, MOS (Metal Oxide Semiconductor) physics and technology, Wiley, New York 1982, and the articles by N. Klein:

Electrical breakdown in thin dielectric films, J. Electrochem. Soc. **116**, 963-972, 1969;

Switching and breakdown in films, Thin Solid Films **7**, 149-177, 1971;

Electrical breakdown of insulators by one-carrier impact ionization, J. Appl. Phys. **53**, 5828-5839, 1982.

[14] See C. M. Osburn and D. W. Ormond, Dielectric breakdown in silicon dioxide films on silicon, J. Electrochem. Soc. **119**, 591-603, 1972;

H. J. de Wit, C. Wijenberg and C. Crevecoeur, The dielectric breakdown of anodic aluminium oxide, J. Electrochem. Soc. **123**, 1479-1486, 1976.

[15] D. R. Wolters and J. J. van der Schoot, Dielectric breakdown in MOS devices, Part III: The damage leading to breakdown, Philips J. Res. **40**, 164-192, 1985.

[16] J. G. Trump and K. A. Wright, Injection of megavolt electrons into solid dielectrics, Mater. Res. Bull. **6**, 1075-1084, 1971.

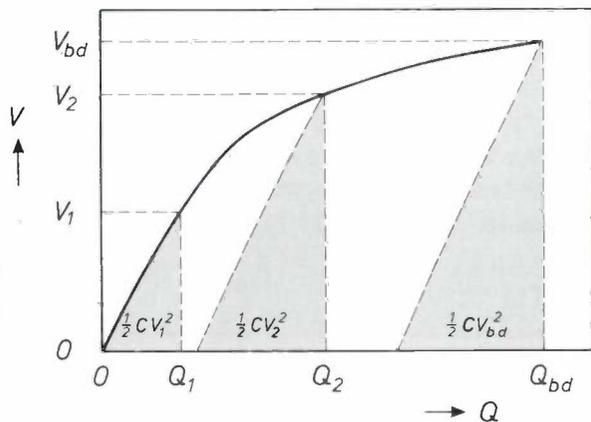


Fig. 19. Schematic representation of the energy content of a charged capacitor. The voltage V is plotted as a function of the charge supplied Q . The area of a shaded region corresponds to the conservatively stored energy and is equal to $\frac{1}{2}CV^2$, where C is the capacitance of the device. Below $V = V_1$ the value of V increases linearly with Q : all the energy supplied is conservatively stored. Above $V = V_1$ there is a less marked increase with Q because of charge leakage. At $V = V_2$ the stored energy is $\frac{1}{2}CV_2^2$; the remainder of the area under the curve up to $Q = Q_2$ gives the energy dissipated by the capacitor. At $V = V_{bd}$, where $Q = Q_{bd}$, the energy dissipation is so great that the film breaks down. The stored energy ($\frac{1}{2}CV_{bd}^2$) is then released and causes permanent damage to the capacitor.

lost, most of it as heat. If the voltage is increased still further, then at V_{bd} the total dissipated energy will damage the capacitor to such an extent that breakdown occurs. The stored energy is then released and causes the damage observed after breakdown (see for example figs 7b and 8). The variation in the crater dimensions with the area of the capacitors confirms that the damage is entirely due to the release of the stored energy ($\frac{1}{2}CV_{bd}^2$).

Analogy with mechanical wear

The behaviour of the oxide film in an operating MOS capacitor can be compared with the behaviour of solid materials under mechanical stress. In both cases there is a force or a field that causes a displacement of material in the mechanical system or a displacement of charge in the capacitor. The energy supplied is to some extent stored conservatively, and to some extent lost as heat. In the conservative storage of energy the mechanical energy is proportional to the

extension (elastic behaviour). If mechanical energy is lost as heat, then the energy is no longer proportional to the extension (plastic behaviour). Something similar takes place in a capacitor: the stored energy is proportional to the charge, until the film starts to leak and energy is converted into heat. Beyond a certain extension a material will often start to oppose further extension ('stress hardening'). The analogy here is that charged centres oppose further injection of charge because of the Coulomb repulsion.

The collapse or fracture of material follows the application of a short-term high stress or a lower stress applied over a longer period (fracture due to creep or fatigue). A dielectric film with a high voltage across it breaks down in a very short time. At a lower voltage it takes longer for breakdown to occur. And just as cracks and dislocations accelerate mechanical wear, so dielectric wear is accelerated by the presence of defects and charges.

Finally, there is also an analogy with the behaviour of materials subject to rapidly changing stress. A solid material that gives considerable extension when subjected to a slowly varying stress may break without any significant extension if the rate at which the stress is applied exceeds a critical value. Below this value the material is considered to be ductile, above it brittle. The deformation after brittle fracture is relatively great. Something similar is found in the electrical behaviour of an oxide film. Above a certain injection rate there is a marked decrease in Q_{bd} (fig. 14a) and the structure of the film is altered to such an extent that the resistance also decreases markedly after the breakdown (fig. 15).

Summary. Charge trapping by the oxide film in a MOS capacitor and the current flowing through it are greatly affected by the interaction between injected electrons and trapped electrons. If an electric field is applied to the film it eventually breaks down. This can happen very quickly if the field-strength is high, if a large amount of charge has leaked through the film or if there are many defects. An important factor in breakdown is the amount of charge transferred. The mechanism for wear and breakdown is reminiscent of the erosion of a landscape by rivers, and there is an analogy with the occurrence of mechanical wear followed by fracture.

1937

THEN AND NOW

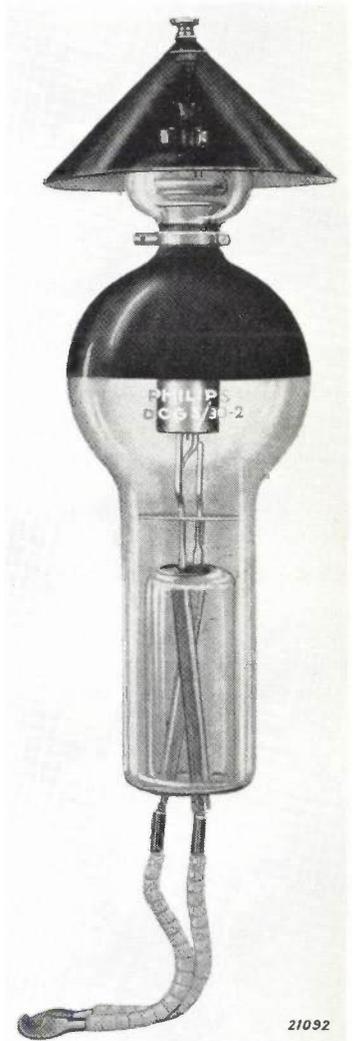
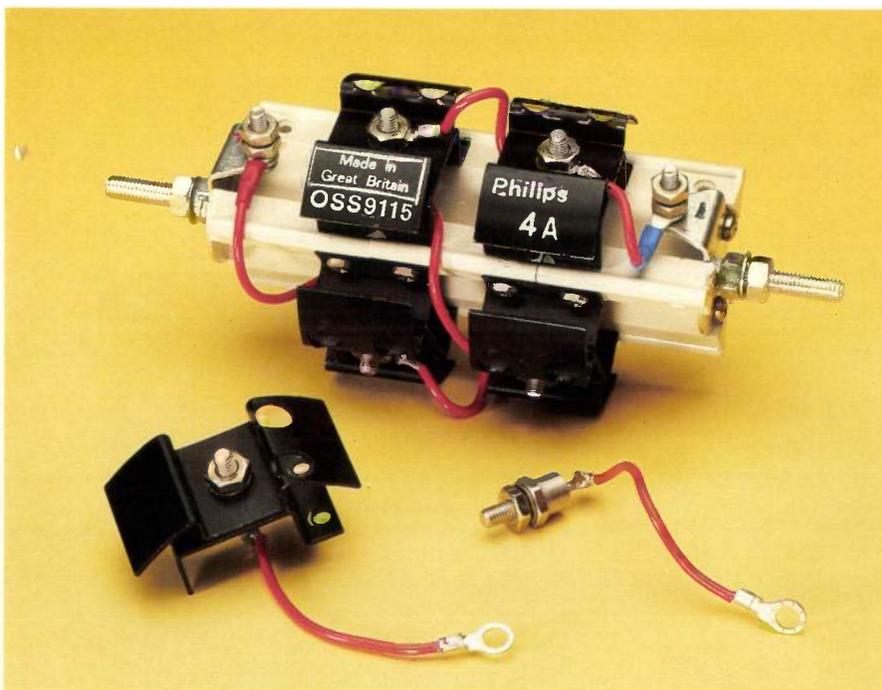
1987

High-voltage rectifiers

The term 'solid-state' probably first makes us think of integrated circuits for low-voltage applications like digital memories or microprocessors. But there are other fields in which solid-state devices have been responsible for radical changes.

The figure^[*] at the right shows the single-phase rectifier DCG5/30 II of 1937, for a current of 6 A and a voltage of 6 kV. This 'rectifying valve' had an incandescent cathode and a mercury-vapour filling. The overall length was approximately 56 cm. A mica cone at the top trapped the rising hot air to prevent condensation of mercury in the top bulb. The ignition voltage of this 'relay valve' could be controlled by varying the voltage on the metal ring just below the upper sphere.

In 1987 high-voltage rectifiers look rather different. The figure below shows a modern solid-state rectifier of the type OSS9115-4A, also suitable for a voltage of 6 kV and a current of 6 A (with oil cooling) or 3.5 A (with air cooling). The rectifier consists of a series circuit of four specially mounted diodes ('diode cells'), whose operation is based on the avalanche effect. A separate diode cell and a separate diode are also shown in the photograph. Up to 36 diode cells can be connected in series, giving rectifiers that will work up to 54 kV. The diode cells are mounted on a non-flammable triangular plastic core. The core is assembled as desired from modules about 4 cm long, which can



each contain up to three diode cells. These basic units are also convenient for use in slightly different diode configurations such as voltage doublers or full-wave rectifiers.

[*] From Philips Technical Review, April 1937.

An optical communication system with wavelength-division multiplexing and minimized insertion losses

II. Multiplexing and demultiplexing

A. J. A. Nicia, C. J. T. Potters and A. H. L. Tholen

MULTIPLEXING

Two kinds of multiplexer

The symbol used in part I of this article^[1] for an optical multiplexing unit shows many of the features of the 'butt-joint' multiplexer (*fig. 1a*), henceforth labelled 'CM', an abbreviation stemming from direct-contact multiplexing. This is the simpler of the two types of multiplexer that have been developed in the course of this system study^[2].

The first, rather striking feature is the direct contact, that is to say the absence of lenses in the coupling between the pigtails (P_1 , P_2) at the input and the fibre (F) for long-distance transmission at the output (all fibres are of the multimode type, with square-law core; see also *fig. 3a* in part I).

The second main feature to be observed is the spatial separation between the incident signals. The directional separation of the pigtails that this requires implies that each input signal (λ_1 , λ_2) has its own optical axis, with three fixed directional coordinates, at the point where the signals combine. The output signal ($\lambda_1 + \lambda_2$) also has a fixed optical axis at this point, which in fact is a 'spatially averaged' continuation of the optical axes of the incident rays. The 'fixed directions' of the various optical axes are of course not fixed in an absolute sense; they refer only to fixed orientations in relation to the multiplexer unit. The angle (2α) that represents the spatial separation has a value of a few degrees in practice. In this respect the longitudinal sec-

tion shown in the figure (*fig. 1b*) is therefore evidently more realistic than the perspective drawing.

The block diagram in *fig. 1c* gives a number of factors that can be used for systematically investigating the extent to which a ray present in P_1 or P_2 is acceptable for guidance by F . As illustrated, the ray will only be accepted for guidance when it is captured within an input aperture cone (NAC) of F . This will be dealt with in more detail in the theoretical section.

The number of pigtails at the input of the multiplexer depicted, two, is of course the smallest number. Four-channel CMs have also been made and tested; an example is shown in the title photograph of part I.

In the pigtails at the input the characteristic diameter for signal guidance (\emptyset in *fig. 3b*, part I) is smaller than in the fibre at the output, which is advantageous for the transmission of the radiant energy. As will appear later, the difference is insufficient, however, for placing the collective cross-section of the pigtail cores completely on the cross-section of the core of the output fibre (compare the situation in the cross-section TP , *fig. 1a, c*). This is therefore the case even with a CM with only two pigtails, which makes insertion

[1] A. J. A. Nicia, An optical communication system with wavelength-division multiplexing and minimized insertion losses, I. System and coupling efficiency, *Philips Tech. Rev.* **42**, 245-261, 1986.

[2] A. J. A. Nicia, Wavelength multiplexing and demultiplexing systems for singlemode and multimode fibers, *Proc. 7th Eur. Conf. on Optical communication*, Copenhagen 1981, pp. 8.1-1 - 8.1-7.

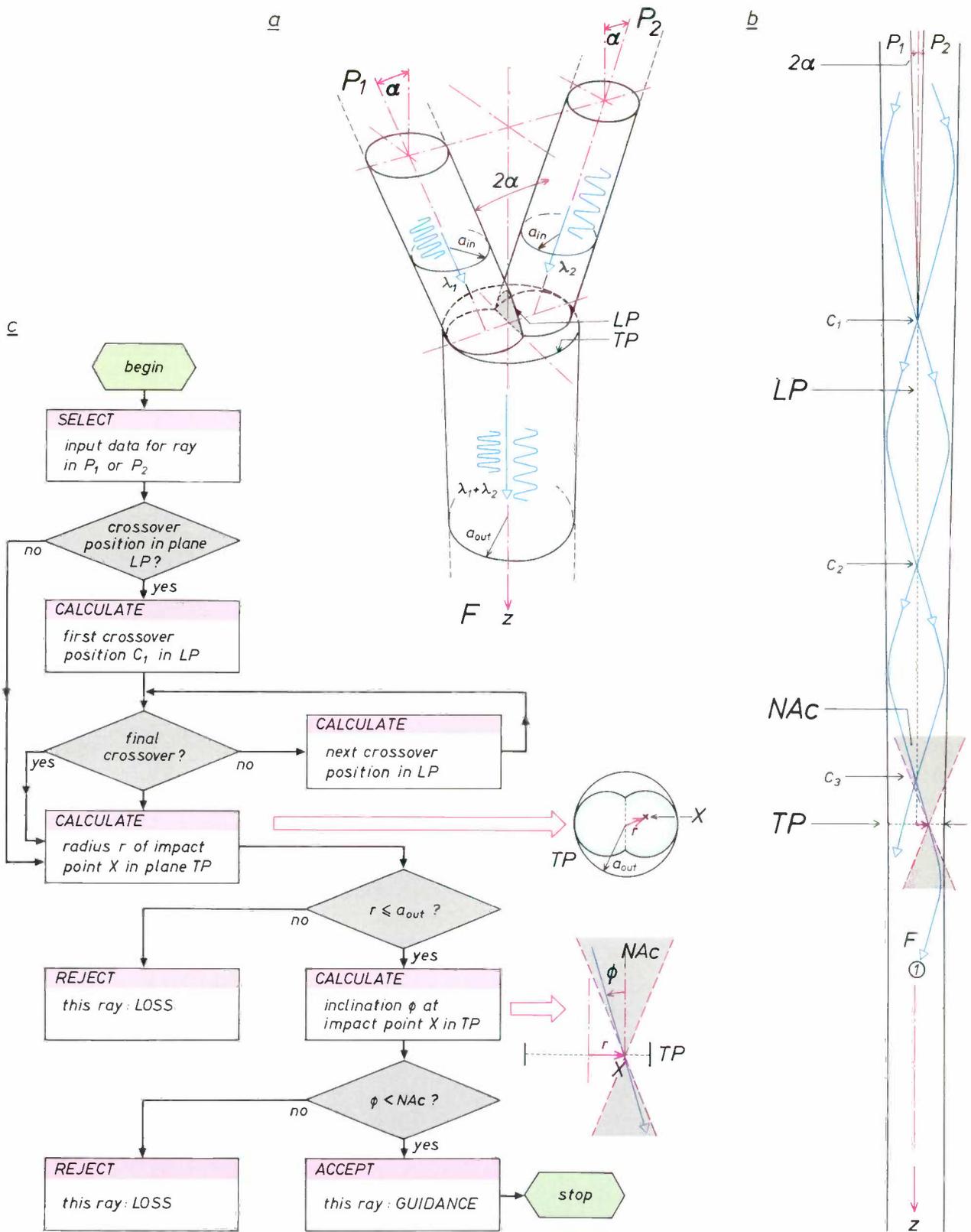
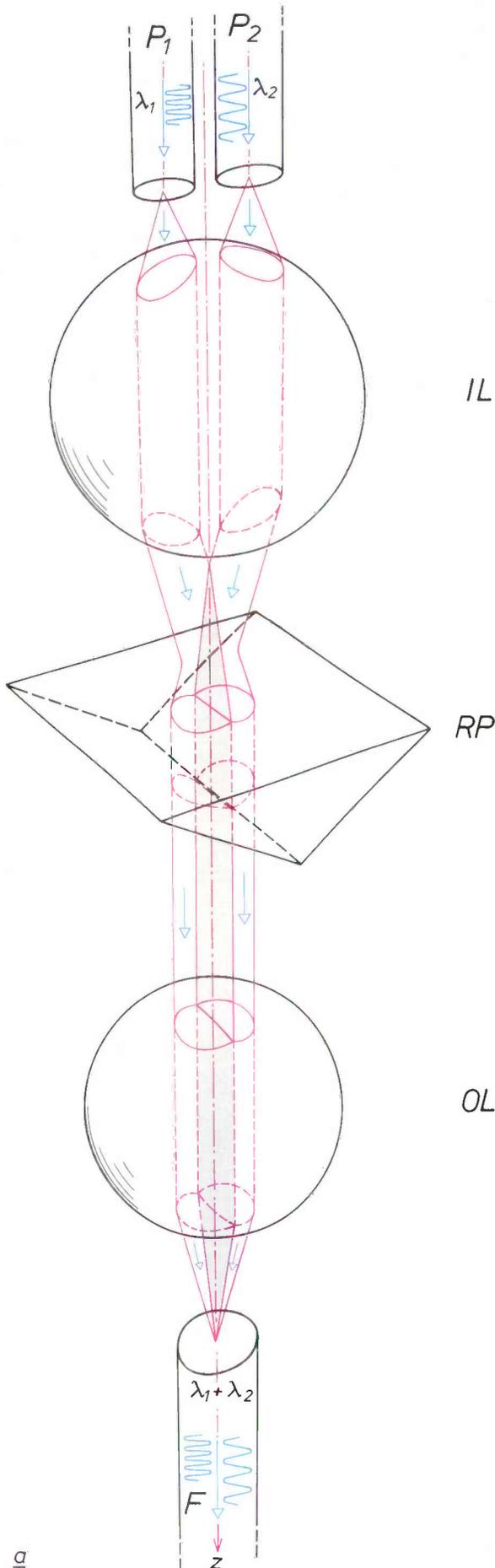


Fig. 1. a) Perspective drawing of a butt-joint multiplexer with two input pigtailed (P_1, P_2) input signals. α inclination angle, giving the slope of the longitudinal polished surface of the cores of P_1 and P_2 . LP longitudinal plane produced by polishing; P_1 and P_2 are bonded here with optical cement. F output fibre. $\lambda_1 + \lambda_2$ output signal. TP contact plane, also core cross-section (πa_{out}^2) of F . More than 80% of the total core cross-section ($2\pi a_{in}^2$) of P_1 and P_2 contributes to the radiation transfer in F . z optical axis of the output signal. b) Longitudinal section of the multiplexer. $c_{1,2,3}$ crossover positions where individual rays pass through the flattened area (see the section on ray tracing). c) Calculations for ray guidance in the multiplexer. NAC input aperture cone of F . ϕ inclination angle of (acceptable) ray, which meets TP at X .



a

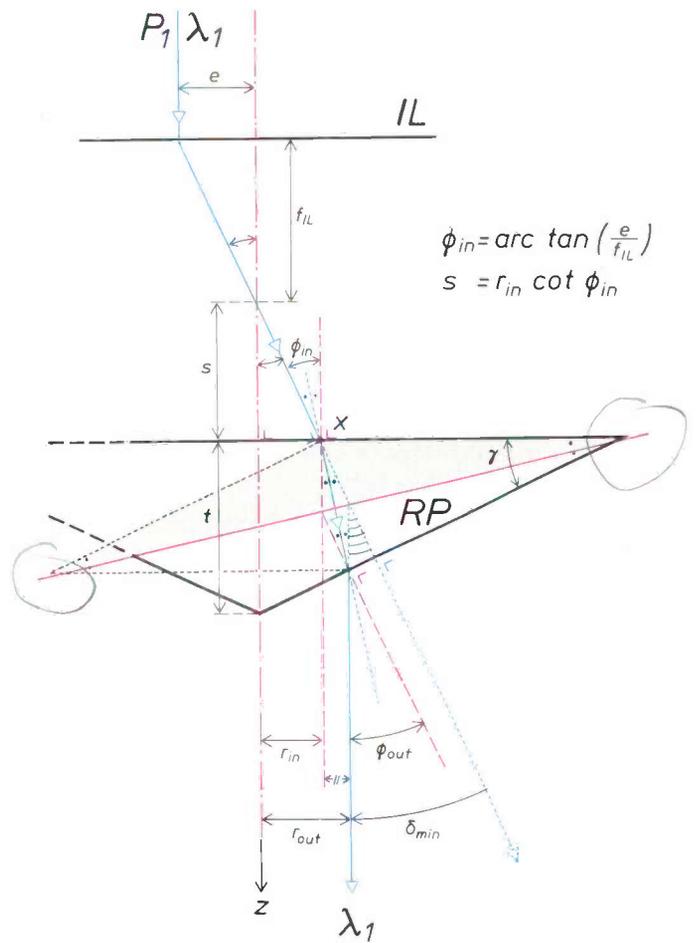
losses at the coupling interface unavoidable. Before comparing some theoretical estimates and measurements of these insertion losses, we shall first say something about the other type of multiplexer and about the method of fabricating both types.

In the introduction of part I the second type of multiplexing unit, the prism multiplexer (or PM), was briefly mentioned as an example. The version shown in *fig. 2a* has two pigtailed (P_1, P_2) at the input. In this case there is no direct contact with the fibre (F) at the output; situated in between are two ancillary lenses (IL, OL) and the prism (RP), known as a roof prism, with an obtuse apex angle.

IL

RP

OL



b

$$r_{in} = r_{out} - (t - r_{out} \tan \gamma) \tan \left(\frac{1}{2} \gamma \right)$$

Fig. 2. a) Perspective drawing of a prism multiplexer with two input pigtailed (P_1, P_2). λ_1, λ_2 input signals. F output fibre. $\lambda_1 + \lambda_2$ output signal. IL, OL ball lenses for collimation and decollimation. RP roof prism. b) Reconstruction of the ray diagram for the signal λ_1 (blue line). z optical axis, also direction of output ray, and of the normal to the base of RP . δ_{min} deviation angle. γ refraction angle. f_{IL} focal length of the lens IL . t height of RP . The chosen 'total' symmetry leads to $\phi_{in} = \phi_{out} = \gamma = \delta_{min}$. The grey triangle, whose apex is the point X where the ray meets RP , forms half of the rhombus shown (partly dashed). The equations stated and the rhombus determine the coordinates of position for the ray paths between the two ball lenses.

The first lens collimates the beams coming from P_1 and P_2 , so that the rays incident on RP , via the base, are all virtually parallel. (Without a collimating lens it would be impossible to combine the different beams into a single relatively narrow beam, which is the hallmark of a successful multiplexing operation.) After passing through the prism all the rays of the two beams take the direction of the optical axis (z).

The second lens (OL) then decollimates the composite beam, so that all the rays 'fit' inside the core and the input aperture of F , allowing their capture and further propagation. The PM, although seeming to be more complicated than the CM, offers more freedom in design and dimensions, since the focal lengths of the two lenses constitute a pair of extra parameters. Depending on the space available, they could be made larger or smaller, rather arbitrarily.

As illustrated in fig. 2*b*, the prism is used in a kind of symmetrical situation (the base being as perpendicular as possible to the z -axis) with minimum deviation. This means that any remaining divergence in the input rays will cause the least perturbation of the output rays, which improves the coupling of the PM to the long-distance fibre (F in fig. 2*a*).

A general rule for a prism is that the sum of the angle of refraction and the angle of deviation is equal to the sum of the angle of incidence (before the first refraction) and the exit angle (after the second refraction). This gives:

$$\gamma + \delta_{\min} = \phi_{\text{in}} + \phi_{\text{out}}. \quad (1)$$

In the case of minimum deviation the angles ϕ_{in} and ϕ_{out} are equal, hence:

$$\cos\left(\frac{1}{2}\delta_{\min}\right) + \cot\left(\frac{1}{2}\gamma\right) \sin\left(\frac{1}{2}\delta_{\min}\right) = n_{\text{RP}}, \quad (2a)$$

where

$$\phi_{\text{in}} = \phi_{\text{out}} = \frac{1}{2}(\gamma + \delta_{\min}). \quad (2b)$$

The quantity n_{RP} is the refractive index of the glass of which the roof prism is made. The value of the refractive index depends on the wavelength of the radiation, of course; the dependence is relatively weak, however, and will therefore not be taken into account here.

As a result of the two refractions in the path of the rays (fig. 2*b*) the directions of ray and normal are interchanged. This special symmetry implies that, in addition to the minimization condition ($\phi_{\text{in}} = \phi_{\text{out}}$), eq. (1) has to meet the condition:

$$\gamma = \delta_{\min}. \quad (2c)$$

Because of the extra condition, eq. (2a) in this case becomes:

$$n_{\text{RP}} = 2 \cos\left(\frac{1}{2}\gamma\right), \quad (2d)$$

so that for a given refractive index the refraction angle is also fixed for the required path of the rays in fig. 2*b*.

Of course it is also possible to use a larger number of input channels with the PM than the minimum number considered here. In that case, however, a modified prism is needed; for example, if the number is doubled to four channels, a roof prism with four facets instead of two is required. Such a prism is shaped like a shallow pyramid with a square base, which should not be too difficult to make.

Fabrication

The butt-joint multiplexer

The photograph (fig. 3) shows a CM of the type with four input channels, i.e. two more than the version used to illustrate the principle in fig. 1. Each of the four channels is provided with a connector. The pigtails in these channels have a core of parabolic material, with a radius of about $16\mu\text{m}$.

To make the multiplexer an easily interchangeable part of the system, the output channel also contains a connector. The pigtail to this connector and the fibre to be connected to it for long-distance transmission have the same parabolic profile and the same core radius ($25\mu\text{m}$) as well [3].

The difference between the core radii is so considerable that in both cases, i.e. with four input pigtails and two, about 52 and 81 per cent respectively of their

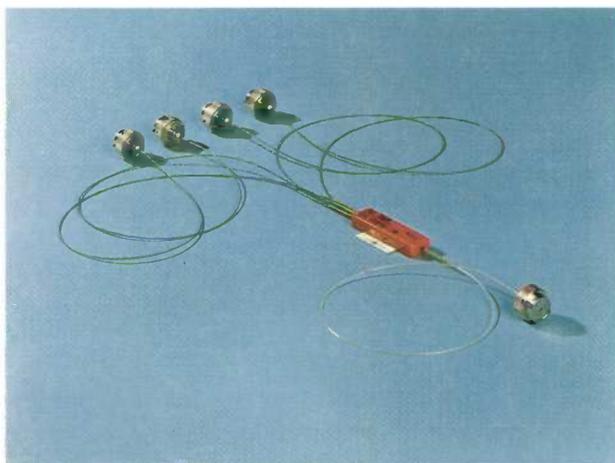


Fig. 3. One of the butt-joint multiplexers. The version shown here, which is about 58 mm long, is suitable for four input channels. The output channel is in the foreground on the right. All the signal channels are terminated in a connector[4].

[3] In part I of the article the symbol F refers solely to the fibre for long-distance transmission. Here, however, F also refers to the output pigtail; both have identical characteristics.

[4] A. J. A. Nicia and C. J. T. Potters, Components for glass-fibre circuits, Philips Tech. Rev. 40, 46-47, 1982; see also A. J. A. Nicia, Lens coupling in fiber-optic devices: efficiency limits, Appl. Opt. 20, 3136-3145, 1981. 'Dry' means that the optical medium in the coupling device between the fibre end faces and the nearest lenses is not an immersion liquid but simply air.

collective cross-section nevertheless remains active, i.e. can continue to contribute to the transfer of radiation through the contact plane — referred to as *TP* in the situation shown in fig. 1*a*. Fig. 4 shows just how the four core cross-sections (1, 2, 3, 4) of the input pigtailed fill the contact plane (the butt joint), which of course is also the end face of the output pigtail. They are just tangential at the periphery of the core. Because of the two-sided (*LP*₁, *LP*₂) longitudi-

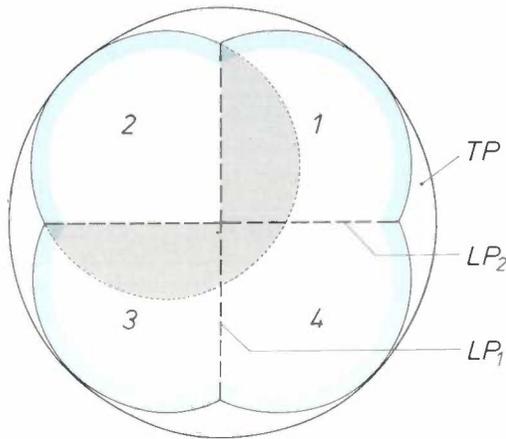


Fig. 4. Transverse section, in the contact plane *TP* of the butt-joint multiplexer in fig. 3. The blue margin forms the boundary of the core cross-sections of the four input channels (1, 2, 3, 4). *TP* is also the end face of the output pigtail core. The input pigtailed are each polished with a slope towards the end on two sides (*LP*₁, *LP*₂) and bonded together with optical cement. The grey hatching at one of the cross-sections (2) indicates the part that no longer contributes to radiation transfer through *TP*, because of the polishing on two sides.

nal flattening of their cylindrical shape, the four input pigtailed can be positioned closely and centrally together. The geometry of this longitudinal modification of the shape of the pigtailed, produced by polishing to an angle of 1.5–2 degrees and optical finishing, is shown in fig. 1*a, b*, for the case with *two* input pigtailed.

A cross-sectional deficit of 48% may perhaps not seem small, but the deterioration it causes in the transfer of radiation through the multiplexer is considerably less than 48%. The parabolic profile of the refractive index, which concentrates the radiant power closely at the axis of the fibre, is the essential background to this relatively smaller effect of the loss of cross-section. The core radius of over 16 μm used for the four input fibres is about the smallest possible value. If pigtailed with an even smaller core radius were to be used — these could in principle be made — it would be necessary to accept substantially higher radiation losses in coupling to available lasers.

The material can be polished with the least chance of fracture by first bonding a few centimetres of the fibre, starting at the butt end, to the inside of a glass capillary. The embedded fibre is then polished ob-

liquely from a point on the outer circumference of the capillary towards the butt end. The butt end itself is finally polished slightly to ensure that it is accurately at right angles to all longitudinal planes.

After being machined in this way, the ends of the input fibres are carefully aligned under a microscope and brought together with a micromanipulator. The polished faces are all coated with a thin film of optical cement to provide a number of firm and completely transparent connections, so that the four butt ends, at least in the region of the flat parts (*LP*₁ and *LP*₂, fig. 4) form a reasonably homogeneous medium for the propagation of radiation. The fourfold assembly is finally bonded to the optically polished end face of the output fibre, again with the thinnest possible film of optical cement, of course. In this way the CM is given a rugged and stable structure; the examples made were no bigger than about 6 cm.

Loss measurements. The extent to which the method of fabrication described can be expected to give satisfactory and reproducible results can be evaluated most reliably from measurements of the radiation loss in a number of CMs made for test purposes. The test radiation, from a halogen lamp in combination with an interference filter, is almost monochromatic ($\lambda \approx 850$ nm, line width $\Delta\lambda \leq 10$ nm). With one photodiode for each input channel, successive measurements are made of the radiant power ($\Phi_{\text{out}}^{(i)}$) available at the multiplexer output and of the radiant power ($\Phi_{\text{in}}^{(i)}$) available in each input channel ($i = 1, 2, \dots$, see fig. 4). The latter measurement is made just in front of the multiplexer, so the input fibre has to be cut^[5].

The insertion loss of an input channel (i) is expressed in decibels by means of a transmission-loss coefficient:

$${}^N\alpha_i = 10 \log_{10} (\Phi_{\text{in}}^{(i)} / \Phi_{\text{out}}^{(i)}), \quad (3)$$

where the superscript N is the number of input channels. The measured value of ${}^N\alpha_i$ also depends on the way in which the power is distributed among the various modes in the relevant input channel. By suitably injecting the test radiation, all modes can be made to contain the same amount of power, producing a well-defined situation, referred to as 'uniform launching'. The configuration is also such that the launched radiation completely 'fills' the input aperture of the test sample. In this situation the measurements will give rather conservative values for ${}^N\alpha_i$.

For some ten samples of a two-channel multiplexer it was found that ${}^2\alpha_i$ was between 0.6 and 0.8 dB, values representing coupling efficiencies of 87 and 83%, respectively. Owing to non-uniform launching, the distribution of power among the modes is in fact

found to be concentrated at the optical axis of the fibres. The coupling efficiency of the CM may therefore be a few per cent better in practice, and for each input channel it can be put at about 90% [6]. It appears from this that the fabrication method described is indeed adequate, and gives a two-channel CM whose characteristics are entirely satisfactory and are also reproducible.

Table I. The spread in the transmission-loss coefficient ${}^4\alpha_i$ (in decibels), measured for each input channel i for four butt-joint multiplexers (A, B, C, D), all with four input channels (fig. 3). In total eight input channels have a transmission loss between 1.5 and 2.0 dB (the grey area); four work better and four not so well (see the sum column Σ_i).

${}^4\alpha_i$ (dB)	input channel i				Σ_i	
	1	2	3	4		
1.2			C		4	
1.3	D			A		
1.4		B				
1.5			A		8	
1.6	B,C			C		
1.7		D				
1.8				D		
1.9		A				
2.0		C				
2.1						4
2.2			D			
2.3	A		B			
2.4				B		

The reproducibility of the CM-type versions made with four input channels (fig. 3) is less satisfactory. The experimentally determined transmission-loss coefficients are listed in *Table I*. Doubling the number of input channels makes it much more difficult to polish and align the individual fibres sufficiently accurately, and to make the extremely thin bonded joints. This leads to a greater spread in the values (1.2 to 2.4 dB, corresponding to coupling efficiencies of 76 and 57.5% respectively, again to be taken for each input channel). The coupling efficiency of these multiplexers, too will turn out to be more favourable in practice.

The prism multiplexer

The structure of a two-channel PM is shown in *fig. 5*. The number of channels is the same as in the illustration of the principle (*fig. 2a*), but it could be increased to four. Most test specimens made thus far, however, are two-channel structures, with an extra 'dummy'

pigtail, for input-channel alignment. The radii of the fibre cores are 16 μm and 25 μm , as for the CM, and the material is again parabolic.

Our design for this version of the PM [7] is derived from the fibre connector described earlier, a 'dry' coupling device with two ball lenses and bayonet catches, which is relatively easy to make [4]. To create space for the roof prism (*RP*, *fig. 5*) the two ball lenses (*IL* and *OL*), with diameters of 5 and 7 mm, are situated somewhat further apart. This is no problem because the rays between two lenses of such a coupling device are virtually parallel.

The relatively large transverse dimension of the beams between the lenses has the advantage that unavoidable imperfections, such as marks on the lens surfaces (due to the penetration of dust, for example) or radial adjustment errors (in the positioning of one lens relative to the other) will cause fewer radiation losses. The alignment of the optical axes of the two lenses is therefore less sensitive to lateral displacements than to directional deviations. Since the lenses themselves are spherical, there are no problems of angular orientation, and hence fewer awkward adjustments.

This design clearly offers a great deal in the way of simplified fabrication and assembly. There are of course considerable benefits in designing a PM in such a manner that only few parts require the ultimate in mechanical finish (errors of form less than 1 to 2 μm). In the design shown in *fig. 5* only one part, the fibre carrier rod *vr*, requires an ultrafine finish. This part is a tiny brass rod with a groove milled in it. The fibres of the input channels are grouped symmetrically and very accurately around the central dummy fibre in this groove, with the three optical axes all in one plane. The optical axis of the dummy fibre is by definition the reference line for aligning the other optical axes in the signal region inside the multiplexer.

The number of positioning and aligning operations that have to be carried out with the highest accuracy on assembly should also be kept to the minimum, of course, so as to cut down on expensive production time. This must not be done at the expense of the coupling efficiency, however. It must be borne in

[6] C. J. G. Verwer *et al.*, Precision fracture of optical glass fibres, *Philips Tech. Rev.* 39, 245, 1980. Subsequent splicing introduces no more than about 0.06 dB of additional loss; see A. J. J. Franken, G. D. Khoe, J. Renkens and C. J. G. Verwer, Experimental semi-automatic machine for hot splicing glass fibres for optical communication, *Philips Tech. Rev.* 38, 158-159, 1978/79.

[6] The insertion losses depend to some extent on the inclination angle of the longitudinal polished plane sloping towards the ends of the input fibres. The chosen value, 1.5 to 2°, gives a minimum. See D. Opielka and D. Rittich, Low-loss optical Y-branch, *Electron. Lett.* 15, 757-759, 1979.

[7] A. J. A. Nicia, U.S. patent No. 4 472 797.

mind that a PM will in any case have higher losses than a CM, which has fewer interface zones and fewer alignment problems.

For fabricating the ball lenses two approaches were successfully adopted. First, a number of ball lenses were made within our company from good optical

glass (LASF9, from Schott, with refractive-index values above 1.80) by spherical polishing followed by optical finishing and antireflection coating. Second, certain commercially available sapphire beads, with a refractive index greater than 1.7, also proved highly suitable. They also have to be coated, or unwanted re-

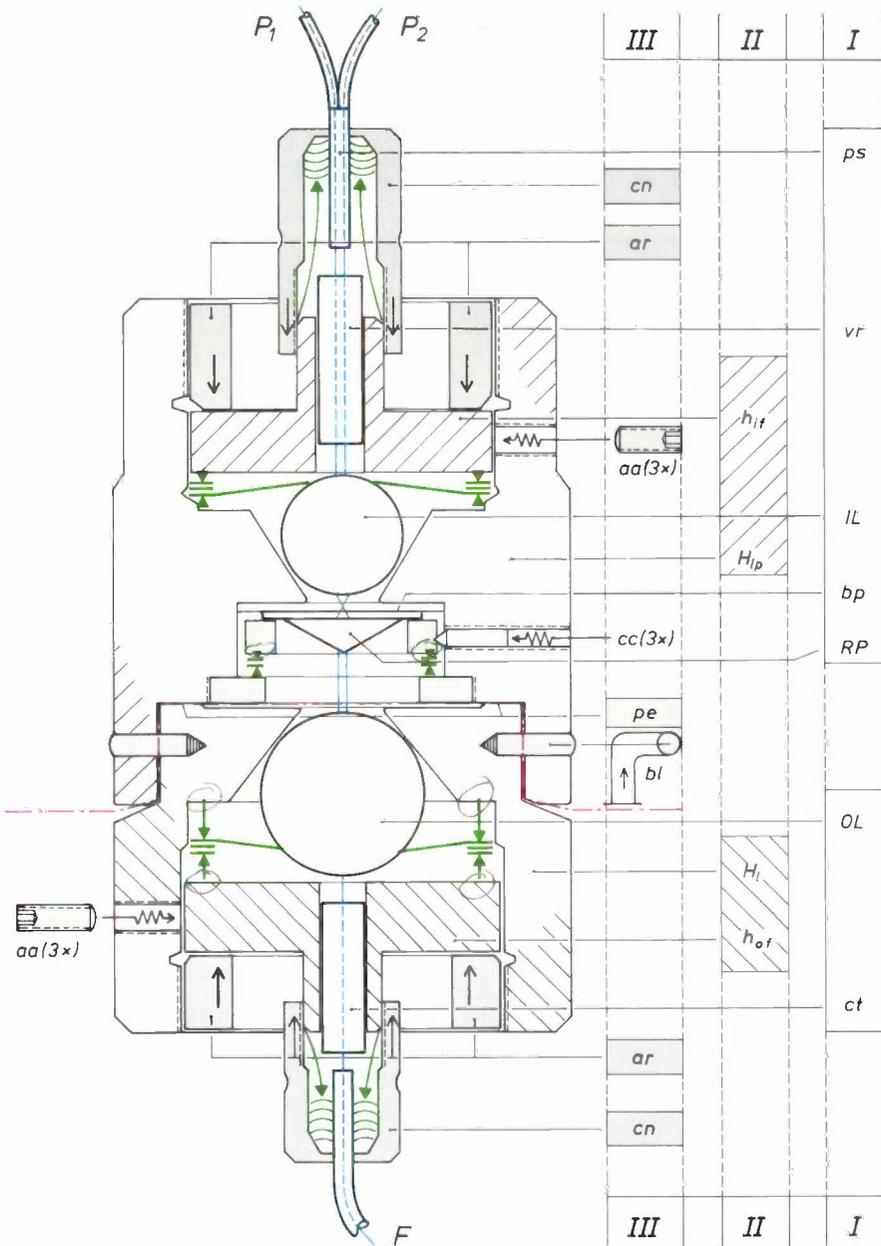


Fig. 5. Simplified diagram of the construction of a prism multiplexer with two input pigtailed (P_1 , P_2 ; see also fig. 2a). F output fibre. The three groups of symbols refer to the signal path (I), the housing (II) and the localization mechanisms and components (III).

I: ps plastic sleeve around P_1 and P_2 . vr fibre carrier rod with V-groove, which acts as a critical reference in the alignment of the optical axes of the fibre cores. IL ball lens. bp base plate fixed to the roof prism (RP). OL ball lens. ct capillary tube, acting as a sleeve for F .

II: h_{if} input fibre holder. H_{ip} metal housing for IL and RP . H_1 metal housing for OL . h_{of} output fibre holder.

III: cn clamping nut, which clamps the plastic sleeve ps via the pressure spring (green). ar pressure-adjusting ring, which fixes the axial location of h_{if} without significant play, by means of counter-pressure from a cupped spring (green). aa , cc three adjusting screws. pe locating edge, which establishes the alignment. bl bayonet catches.

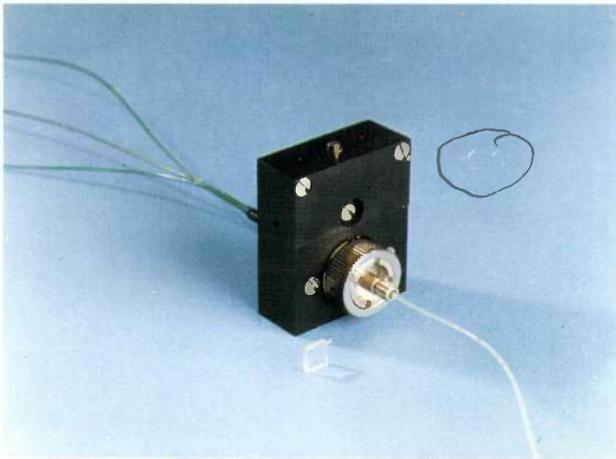


Fig. 6. This prototype prism multiplexer, whose design corresponds to the design shown in fig. 5 (though with a slightly different housing), has three input pigtails. Two of them are signal channels, and a third one (in the centre) is a dummy, used only for optical alignment. In the foreground the roof prism can be seen; after one of the bayonet catches has been opened the prism can be mounted in the central element (the black box). The pressure-adjusting ring and the clamping nut (*ar* and *cn*, fig. 5) can both be seen close to the bayonet coupling on the output side.

reflections would introduce an extra insertion loss of at least a few tenths of a decibel.

The photograph (*fig. 6*) shows a PM fabricated here, a two-channel prototype with the dummy fibre mentioned earlier. This version, with the 'black box' for the prism as the central element, is a rather more 'modular' version of the design in *fig. 5*. The black box contains the mechanisms for the precise adjustment and correct positioning of the prism, which are functionally identical with those shown in *fig. 5*. The black box has a 'half-bayonet' coupling on opposite sides, each containing a ball lens. The structure of such a half-coupling is almost identical with the part of the arrangement in *fig. 5* containing the larger ball lens (the lower part). The large knurled nut ($\varnothing 2$ cm) in the photograph is in fact the adjusting ring *ar* in *fig. 5*, and the small one corresponds to *cn*, the clamping nut for relieving tensile stresses (which could all too easily cause fracture of the fibres).

The largest dimension of this PM is about 6 cm, which is therefore little different from the CM described before. The focal lengths of the ball lenses are only a few millimetres. The focal planes themselves, at least the upper and the lower ones (see *fig. 2a*), ought of course to coincide with the fibre end faces, because of the need for the rays to be parallel in the region between the two ball lenses. In practice it turns out to be better to position the fibre end faces slightly *outside* the focal plane, a few microns away from the lens surface. The explanation for this 'erroneous' positioning — which fortunately has very little effect on

the coupling efficiency — lies in the procedure used for aligning the various optical axes inside the PM^[8]. Mounting and aligning such a multiplexer, consisting of two or three modules, is best done on a special optical bench.

The loss measurements. In the same way as described for the CM, the radiation loss in a number of test samples of the PM was measured under uniform launching conditions. Versions with two input channels (plus dummy fibre) were investigated and also one with four input channels (with no dummy fibre). The

Table II. The transmission loss and its four causes (*r*, *o*, *sa*, *m*) in a two-channel prism multiplexer. The coefficients giving the magnitude of the transmission (insertion) loss are $\delta\alpha_i$ for partial losses and ${}^2\alpha_i$ for the total loss (both in decibels and taken per input channel). The third column shows how the value of the coefficient is found. The loss contribution due to overlap of the input signals at the roofline of the prism (*RP* in *fig. 2a*) is estimated by a computing method (*ca*) based on the 'cross-sectional deficit' (see the section 'Theoretical approaches'). Overlap and spherical aberrations are almost unavoidable loss factors (the grey area). It appears that antireflection coating can reduce the reflection losses to less than 0.2dB, and it seems that the losses due to mechanical tolerances are still well above the fundamental limit by far. The ray-tracing method mentioned in note [b] is essentially derived from the acceptance test in *fig. 1c*. *r* reflection, *o* overlap, *sa* spherical aberrations, *m* mechanical tolerances, *me* measurement, *ca* calculation, Σ summation.

Cause of transmission loss	Magnitude/input channel (dB)	Method of determination
Coefficient $\delta\alpha_i$		
reflection	0.35 ^[c]	me
overlap	0.2	ca ^[b]
spherical aberrations	0.3	me ^[b]
mechanical tolerances	≈ 0.6	${}^2\alpha_i - \sum_{r, o, sa, m} \delta\alpha_i$
Coefficient ${}^2\alpha_i$		
$\sum_{r, o, sa, m}$	1.45 ^[a]	me

[a] best value

[b] ray tracing gives $\sum_{o, sa} \delta\alpha_i = 0.7$ dB

[c] uncoated roof prism

best result, obtained from one of the two-channel versions, was a transmission-loss coefficient of 1.45 dB, which corresponds to a coupling efficiency of nearly 72%. In the total transmission loss four contributions can be distinguished (*Table II*). The reflection losses (0.35 dB) are due to the absence of an antireflection coating on the prism under test. The estimated 0.2 dB

[8] See also A. J. A. Nicia, Practical low-loss lens connector for optical fibres, *Electron. Lett.* 14, 511-512, 1978.

is a loss connected with the fact that the two incident beams, entering via the base of the prism, overlap slightly in the centre after the first refraction (compare fig. 2). The fraction of the radiation arriving at the wrong side of the roofline of the prism will be lost. The partial loss of 0.3 dB is due to spherical aberration in the two ball lenses. This can be measured, via the central dummy fibre, just before the prism is mounted. The fourth contribution to the total transmission loss arises from a number of small inaccuracies in the various mechanical parts of the PM. With our method of fabrication and assembly that loss should therefore not amount to more than about 0.6 dB.

A final point to be noted about the favourable effect of ball lenses on the coupling efficiency is that the high values of the refractive index, close to 2, help to minimize the spherical aberration (for a given lens strength the lens curvature may be kept correspondingly small). Since spherical aberration is the main lens imperfection^[9], considerable importance is attached to making it smaller.

Theoretical approaches

From the description above it is clear that both types of multiplexers work without requiring an extra component with a strong wavelength-selective effect. Demultiplexers, on the other hand, which are the next main topic of this system study, cannot do without such a component, a good example of which is an interference filter^[10].

Efficiency and radiant power

The use of fibres for the input channels with an undersized core radius (compared with the output fibre) is an approach that makes it possible, in the multiplexers at least, to omit an additional wavelength-selective component without incurring losses. Even with uniform launching conditions, the most critical situation, the input loss per channel can then be kept within acceptable bounds. In the case of fibre cores with a parabolic refractive-index profile this can be seen from a well-known expression for the radiant power guided by such a fibre^[9]:

$$\Phi = \frac{1}{2}\pi^2 L (a \times NA)^2, \quad (4)$$

where a and NA are the radius and the (maximum) numerical aperture of the fibre core, and L is the radiance, assumed to be constant across the fibre core (radiance expressed in watts per square metre and per steradian is a kind of surface brightness). In the case of a multiplexer with N input channels, which ideally all have the same efficiency for energy transfer to the

output channel, the partial radiant powers available at the output of the multiplexer are given by:

$$\Phi_{\text{out}}^{(i)} = \frac{1}{2}\pi^2 \frac{L_{\text{out}}}{N} (a_{\text{out}} \times NA_{\text{out}})^2, \quad (5a)$$

where i is the number 1, 2, ... of the input channel. The radiant power in the corresponding input channel is similarly equal to:

$$\Phi_{\text{in}}^{(i)} = \frac{1}{2}\pi^2 L_{\text{in}} (a_{\text{in}} \times NA_{\text{in}})^2. \quad (5b)$$

The theoretical efficiency ^{$N\eta_i$} , expressed in decibels, of the input channel can be calculated from:

$$N\eta_i = 10 \log_{10} (\Phi_{\text{out}}^{(i)} / \Phi_{\text{in}}^{(i)}). \quad (6a)$$

The quantity ^{$N\eta_i$} corresponds to the transmission-loss coefficient introduced earlier (^{$N\alpha_i$} , eq. 3), except for the sign. Assuming that a multiplexer is a passive element and, again ideally, involves no extra radiation loss due, for example, to reflection or absorption and scatter, the radiance will be maintained constant during the transfer of the radiant power (i.e. $L_{\text{in}} = L_{\text{out}}$)^[11].

Using this we can derive from eq. (6a) the expression:

$$N\eta_i = -10 \log_{10} N + 20 \log_{10} \left(\frac{a_{\text{out}} \times NA_{\text{out}}}{a_{\text{in}} \times NA_{\text{in}}} \right). \quad (6b)$$

It is easily seen from this expression for the theoretical coupling efficiency how unsatisfactory the choice of equal radii would be ($a_{\text{out}} = a_{\text{in}}$, and hence $NA_{\text{out}} = NA_{\text{in}}$). In the case of a multiplexer with four input channels ($N = 4$), for example, the coupling efficiency, taken for each channel, would be only about -6 dB, which is certainly not acceptable.

It can also be seen from eq. (6b) how the coupling efficiency can be made equal to 0 dB by a better choice of $a_{\text{out}}/a_{\text{in}}$. It is therefore necessary to ensure that:

$$\frac{a_{\text{out}}}{a_{\text{in}}} = N^{\frac{1}{2}} \frac{NA_{\text{in}}}{NA_{\text{out}}}. \quad (7a)$$

In commonly used fibres there is little variation in the (maximum) numerical aperture, so that the condition formulated in eq. (7a) implies the approximate equality of the total core cross-sections in accordance with:

$$N a_{\text{in}}^2 \approx a_{\text{out}}^2. \quad (7b)$$

The choices mentioned earlier for a_{in} and a_{out} satisfy eq. (7b) reasonably well, at least when N is equal to

two. They also, of course, match the 'trumpet geometry' used in the system, as could be seen in fig. 3b of part I.

Efficiency and the cross-sectional deficit (CM)

Eq. (6b) for the coupling efficiency gives a result that must differ from the reality for the CM, because of the assumption that the output-fibre core is completely 'filled' with radiation.

Looking again at fig. 4, we can estimate the radiation loss by determining which part of the total circular cross-section of an input channel, calculated in the contact plane (*TP*), is no longer available for energy transfer. (This part is shown grey for channel 2.) This means that a certain fraction of the total radiant power is lost in this channel. Calculation of that fraction gives the required estimate of the radiation loss. It is no more than a rough estimate, because the relative loss of cross-section is taken as the starting point. This represents the simplest 'geometrical translation' of the effect of the bevelling of the end of a fibre core on the transfer of radiation through it.

The loss fraction is relatively easy to calculate by using a well-known expression for the radiant exitance (M_e), the guided power per unit area in the transverse section of a parabolic fibre core^[12]:

$$M_e(r) = M_{e,co} \left\{ 1 - \left(\frac{r}{a_{in}} \right)^2 \right\}, \quad (8)$$

where r is the radial coordinate inside the fibre core; the optical axis of the core is characterized by equating r to zero. It can be seen that M_e has cylindrical symmetry and that M_e does not remain constant as a function of the radius r but decreases continuously from the maximum value $M_{e,co}$ (at $r=0$). At the edge of the core, where of course the guidance of radiation ceases, M_e will in fact exactly vanish.

The derivation of eq. (8) is based on the fact that L , the radiance, is constant over the cross-section of the fibre core (because of the postulated uniformity of the original injection of radiation into the input channel). Applying the general definition of radiance to optical fibre cores we find:

$$M_e(r) = \{ \pi \times NA^2(r) \} L, \quad (9)$$

where $NA(r)$ is a completely radiation-filled local numerical aperture. (The product $\pi \times NA^2(r)$ gives the corresponding solid angle in steradians.) For $NA(r)$ it can be shown that^[12]:

$$NA(r) = \{ n^2(r) - n^2(a_{in}) \}^{1/2}, \quad (10)$$

where $n(a_{in})$ is the refractive index of the cladding material.

Eq. (10) may be seen as a generalization of the analogous expression for NA_{co} , the maximum value that the local numerical aperture can have (and which, having been introduced in part I as NA , the maximum numerical aperture, has already been used here in eqs (4)-(7)).

The material of the fibre core is parabolic, which means that (see eq. 3 in part I):

$$n(r) = n_{co}(1 - g^2 r^2)^{1/2}, \quad (11)$$

where n_{co} is the value of the refractive index on the optical axis and g is a fibre constant. Substitution (twice) of eq. (11) in eq. (10) gives:

$$NA(r) = n_{co} g a_{in} \left\{ 1 - \left(\frac{r}{a_{in}} \right)^2 \right\}^{1/2}. \quad (12)$$

Bearing in mind that NA_{co} must be equal to $n_{co} g a_{in}$, we readily arrive at eq. (8) from eq. (9) and eq. (12).

The calculation of the loss fraction is further illustrated here with the example of four input channels. This requires summation of the radiant exitance (eq. 8) both over the full cross-section, a circle of diameter $2a_{in}$, of one of the input channels, (e.g. No. 2, see fig. 7) and over the part of this cross-section to be polished away. In fig. 7 this part is bounded by two dashed straight lines (LP_1, LP_2) and the grey arc. The point H is the centre of the grey arc and thus lies exactly on its symmetry axis (the line C_2M_0). It can be seen that the 'marginal' centre-point (M_0 , relating to the situation just before the occurrence of insertion loss) of the output channel and the centre-point (C_2) of the input channel used here establish an interval in the plane *TP* over which the centre-point (M) of the output channel can be varied by the designer of the multiplexer. The amount of material removed by polishing determines just where the point M will come; corresponding to this choice as a kind of independent variable is the product $a_{in}^{-1} \times \frac{1}{2} R_4 (= \cos \psi_{max})$. The radius a_{in} is taken as fixed, the radius a_{out} is variable, however. The tangent point T remains where it is, of course.

Summation of the radiant exitance is in all cases simplest when eq. (8) is integrated over the radial coordinate (r) and over the azimuthal coordinate (ψ).

[9] A. J. A. Nicia, Micro-optical devices for fiber communication, thesis, Eindhoven 1983.

[10] See for example: Iridescence in technology and nature, Philips Tech. Rev. 41, 149, 1983/84.

[11] This equality is sometimes called the radiance law. A more accurate formulation states that L/n^2 (and hence not L itself; n is the local refractive index) in a passive and lossless optical element remains constant along every possible ray path in the element. See also G. Cancellieri and U. Ravaoli, Measurements of optical fibres and devices: theory and experiments, Artech House, Dedham, Mass., 1984.

[12] D. Gloge and E. A. J. Marcatili, Multimode theory of graded-core fibers, Bell Syst. Tech. J. 52, 1563-1578, 1973.

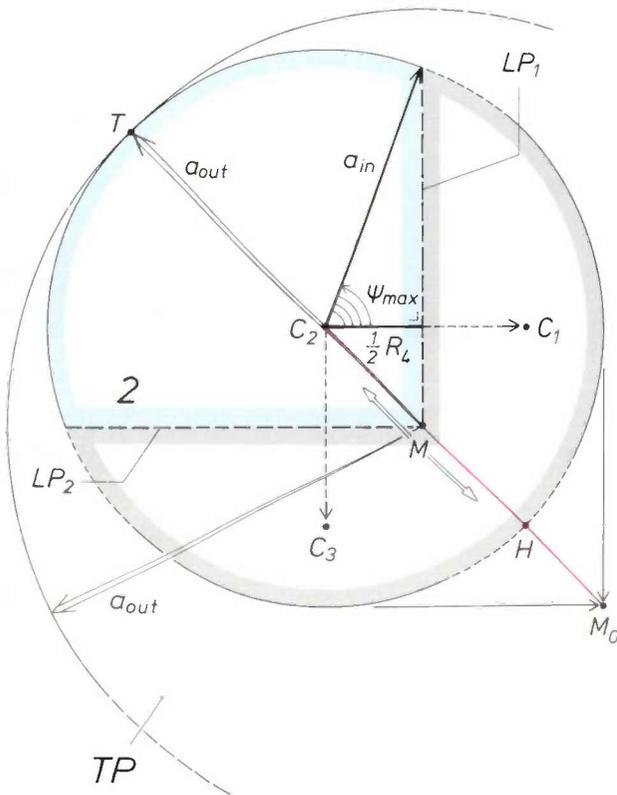


Fig. 7. Configuration details for the location of an input channel (No. 2, fig. 4) on the contact plane *TP*, i.e. the end face of the output channel (core radius a_{out}), for a butt-joint multiplexer with four input channels (core radius a_{in}). The part polished away (framed in grey) from the core cross-section of 2, as seen in the plane *TP*, is bounded by LP_1 , LP_2 and the arc with the centre-point H . The part with the blue edge determines the total radiant power available from 2 for transmission and hence for the multiplex operation. The significance of the other symbols is explained in the text.

The result of the integrations is given in *fig. 8* as a plot of the coupling efficiency, the desired quantity.

The total radiant power in the selected input channel is:

$$\Phi_{in} = \int_{\psi=0}^{2\pi} \int_{r=0}^{a_{in}} M_e(r) dr d\psi = \frac{1}{2} \pi a_{in}^2 M_{e,co}. \quad (13)$$

There are also two expressions for the insertion loss ($\Delta\Phi_{in}$), corresponding to the choice of the centre-point M either in the interval C_2H , or in the interval HM_0 . The equations for ${}^4\eta_i$ are:

$$\frac{1}{100} {}^4\eta_i = 0.75 - \frac{\psi_{max}}{\pi} + \frac{\sin(2\psi_{max})}{3\pi} \times [1 + \{2 - \cos(2\psi_{max})\} \{ \frac{1}{2} + \cot \psi_{max} \}], \quad (14a)$$

for $0 \ll \cos\psi_{max} \ll \frac{1}{2}\sqrt{2}$, which corresponds to the interval C_2H , and

$$\frac{1}{100} {}^4\eta_i = 1 - 2 \left[\frac{\psi_{max}}{\pi} - \frac{\sin(2\psi_{max})}{6\pi} \{4 - \cos(2\psi_{max})\} \right], \quad (14b)$$

for $\frac{1}{2}\sqrt{2} \ll \cos\psi_{max} \ll 1$, which corresponds to the interval HM_0 .

An independent variable that can also be used in the efficiency calculations is the ratio a_{out}/a_{in} , given by:

$$\frac{a_{out}}{a_{in}} = \frac{\frac{1}{2}R_4}{a_{in}} \sqrt{2} + 1. \quad (15)$$

The values of a_{in} and a_{out} used in practice thus provide a calculated coupling efficiency of over 61.5%. This calculated result, although well within the total spread of 57.5 to 76% (mentioned in the section on the fabrication of the CMs) as established from the coupling efficiencies measured for the prototypes, is on the conservative side. Such four-channel CMs will in general have a somewhat smaller insertion loss than may be predicted from *fig. 8* (see *Table I*). Also plotted in *fig. 8* is the quantity ${}^4\sigma_i$, the purely geometric effect of removing material by polishing, expressed in the active cross-section remaining for radiation transfer (the dashed curve C_{eff}), again as a function of $a_{in}^{-1} \times \frac{1}{2}R_4$ calculated as a percentage of the complete cross-section. The fact that ${}^4\eta_i$ is well above ${}^4\sigma_i$ is related to the parabolic profile of the refractive index, a choice which therefore helps to give a good coupling efficiency.

Matching of the local numerical aperture. When radiation is transferred through the contact plane, each ray should of course emerge from the input channel at an inclination angle ϕ small enough to enable it to remain within the local numerical aperture of the output channel (in other words within the cone NA_c , as depicted in *fig. 1c*). It should be noted here that the assumption made earlier that L remains constant during radiation transfer already implies such matching of local numerical apertures.

The trumpet geometry discussed in part I also helps in meeting this constraint on the inclination angle ϕ in the plane *TP*, certainly if the dimensioning of the fibres satisfies the relation:

$$\frac{a_{in}}{a_{out}} \geq \left(\frac{NA_{in,co}}{NA_{out,co}} \right)^2. \quad (16a)$$

This inequality relation stems from the requirement that for each point of the cross-section of the input channel (i.e. each point inside the blue boundaries in *fig. 7*) the local numerical aperture of the channel should be smaller than — or at worst equal to — the corresponding local numerical aperture of the output channel. This local numerical aperture has circular symmetry, of course, and its relation to the radius r , the distance from the point considered in the contact plane to the centre-point M , is given by:

$$\frac{NA_{out}^2(r)}{NA_{out,co}^2} + \frac{r^2}{a_{out}^2} = 1. \quad (16b)$$

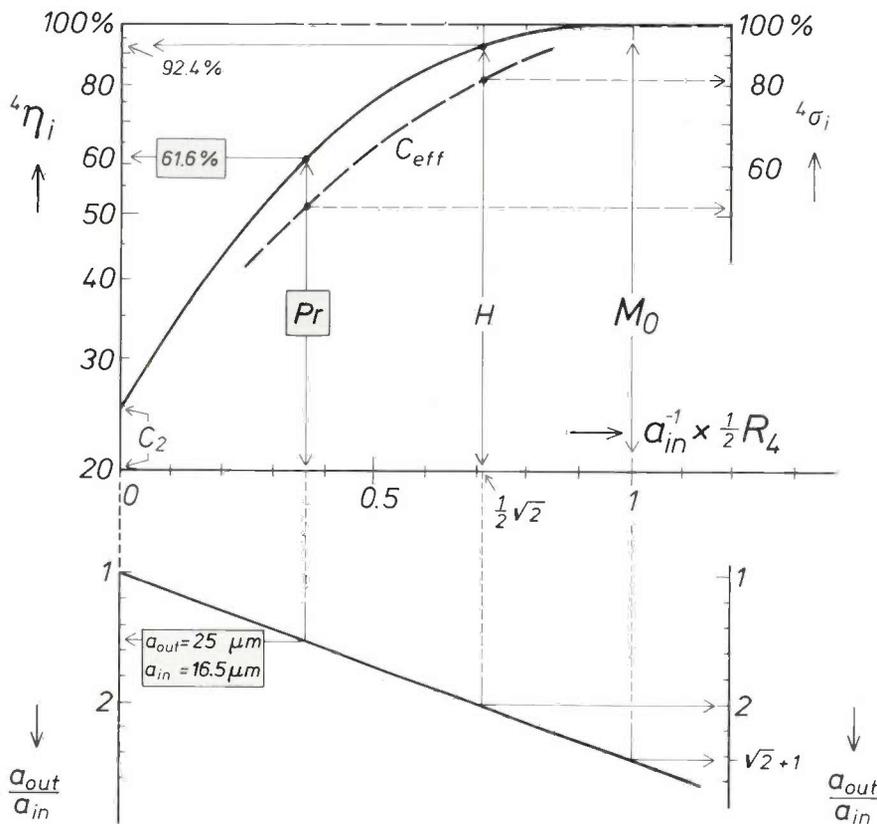


Fig. 8. The coupling efficiency ${}^4\eta_i$ calculated for each channel i for a butt-joint multiplexer with four input pigtails. The degree of polishing (on two sides) of the fibre core material is considered to be variable, so that the dimensionless quantity $\frac{1}{2}R_4/a_{in}$ can be treated as an independent variable (see fig. 7 for the associated contact-plane configuration). The core radius a_{in} remains constant, while R_4 and a_{out} vary. These variations also obey the linear relation shown between $\frac{1}{2}R_4/a_{in}$ and a_{out}/a_{in} . In the limiting case where polishing of the input channels is only just superfluous ($\frac{1}{2}R_4/a_{in} = 1$, M coinciding with M_0), the value of ${}^4\eta_i$ is at its maximum (100%). With the maximum polishing ($\frac{1}{2}R_4/a_{in} = 0$, M coinciding with C_2) ${}^4\eta_i$ drops to its lowest value (25%). Pr practical case calculated with $a_{in} = 16.5\mu m$ and $a_{out} = 25\mu m$. If M_0 is the initial situation, then in the situation H polishing reaches the point where the sides LP_1 and LP_2 'really' intersect for the first time. C_{eff} curve for the calculated coupling efficiency (${}^4\sigma_i$) assuming constant radiant exitance (see eq. (8) in the text); with this assumption only the loss in effective cross-section of the input fibres is taken into account. $i = 1, 2, 3, 4$.

This expression for $NA_{out}(r)$, which can easily be derived from eq. (12), is a well-known equation for an ellipse. The semi-axes of the ellipse are $NA_{out,co}$ and a_{out} (fig. 9). The local numerical aperture for the input channel follows from a similar equation, not given here. Fig. 9 gives a plot of the pairs of values of the calculated local numerical apertures for all the points of the path MC_2T in the contact plane TP (see fig. 7). It can be seen that to the right of the point of intersection of the two branches of the ellipses the local numerical apertures of the input channel, $NA_{in}(r)$, are too large (the grey segment between 'suff' and the top point T). This is due to a minor violation of eq. (16a). (The inequality relation is simply a reformulation of the requirement that at the common top point T the ellipse for the input channel should have greater curvature than the ellipse for the output channel (fig. 9).)

Efficiency and ray tracing (CM)

We now return to the acceptance test for the CM, shown earlier as a block diagram in fig. 1c. This test is based on ray tracing, and also helps to give a better understanding of the behaviour of the coupling efficiency. This approach gives a good approximation to the physical reality of radiation transfer through the contact plane (TP) and permits comparison with the results of the method of calculation based on the cross-sectional deficit in TP .

A computer program for the test has been written. It follows the general pattern of fig. 1c and also includes a number of refinements discussed elsewhere^[9]. Some efficiency curves found with the program are shown in fig. 10, relating to CMs with two input channels. Two values (1.5 and 0.1°) have been chosen for the inclination angle (α) that determines

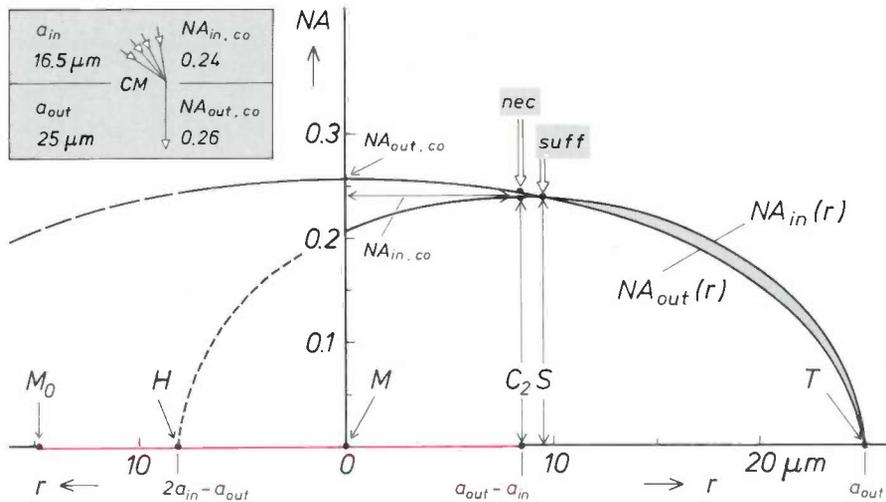


Fig. 9. Calculated behaviour of $NA_{in}(r)$ and $NA_{out}(r)$, the two local numerical apertures in the contact plane (TP , see fig. 7) in a butt-joint multiplexer with four input pigtails. The two curves (branches of an ellipse) relate to the practical case, whose four characteristic quantities are quoted in the grey block. The variable plotted horizontally is the distance r between the contact point in TP and the centre-point M . In the grey segment, to the right of 'suff' the solid angle of the input channel is too large, which prevents a (fairly small) fraction of the radiation from being transmitted^[13]. $NA_{in,co}$, $NA_{out,co}$ maximum numerical apertures of the input and output channels, respectively. 'nec' refers to C_2 , the point where $NA_{in}(r)$ reaches its maximum (which is smaller than $NA_{out}(r)$ at the same position). The significance of all the other symbols is as in fig. 7.

the slope at which material was polished away towards the end face of each input fibre. The values available in practice for the core radii, $a_{in} = 16.5 \mu m$ and $a_{out} = 25 \mu m$, form the basis of all the calculated points. In the contact plane, therefore, only the distance $\frac{1}{2}R_2$ was varied (from 1.5 to 10.5 μm). When $\frac{1}{2}R_2$ is equal to 8.5 μm , we have the preferred coupling geometry (Pr), with maximum 'fill' of the contact plane. In the case of the two points with $\frac{1}{2}R_2$ greater than 8.5 μm , the input channels have a certain 'overhang' (indicated in red), which implies additional insertion losses.

The variation of R_2 (in a path indicated by the colours green and red) simulates the polishing away of core material, always in the direction towards the end face of the fibre (P_1), where only the polished length (l_{gr}) is varied, in accordance with the linear equation:

$$\frac{\alpha \times l_{gr}}{a_{in}} = 1 - \frac{R_2}{2a_{in}}, \quad (17)$$

which is also plotted in fig. 10.

Measurements versus calculations. In the case of the preferred coupling geometry (Pr) the calculated coupling efficiencies are 89.5% and 87.8%, as can be seen from the figure. Measurements on prototype CMs with two input channels gave values between 87 and 83%, as discussed in the section on fabrication. For use in the system a measurement will indicate a slight-

ly better value; a value of about 90% would be expected. The calculated results are thus reasonably comparable with this better value. It may be concluded from this that the test calculations in fig. 1c of the acceptance of radiation in P_1 or P_2 for further guidance (through the output channel of the CM) are in every way satisfactory.

The third curve (SF) in fig. 10 was calculated from:

$$\frac{1}{100} {}^2\eta_i = 1 - \frac{1}{\pi} \psi_{max} + \frac{1}{6\pi} \sin(2\psi_{max}) \{4 - \cos(2\psi_{max})\}, \quad (18)$$

an equation derived by taking the cross-sectional deficit in TP into account, and analogous to eq. (14). This method of calculation is a relatively rough approximation — although the agreement for the case Pr is surprisingly good — which was considered earlier for the case of four-channel CMs (fig. 8); the equation then used (14b) contains a correction term that is twice as large:

$$1 - \frac{1}{100} {}^4\eta_i = 2 \times (1 - \frac{1}{100} {}^2\eta_i). \quad (19)$$

Simulating the oblique polishing. When the value of the quantity $\frac{1}{2}R_2$, in the geometry Pr , is gradually reduced (this corresponds to polishing away more material) the behaviour of the calculated points, in the 'green

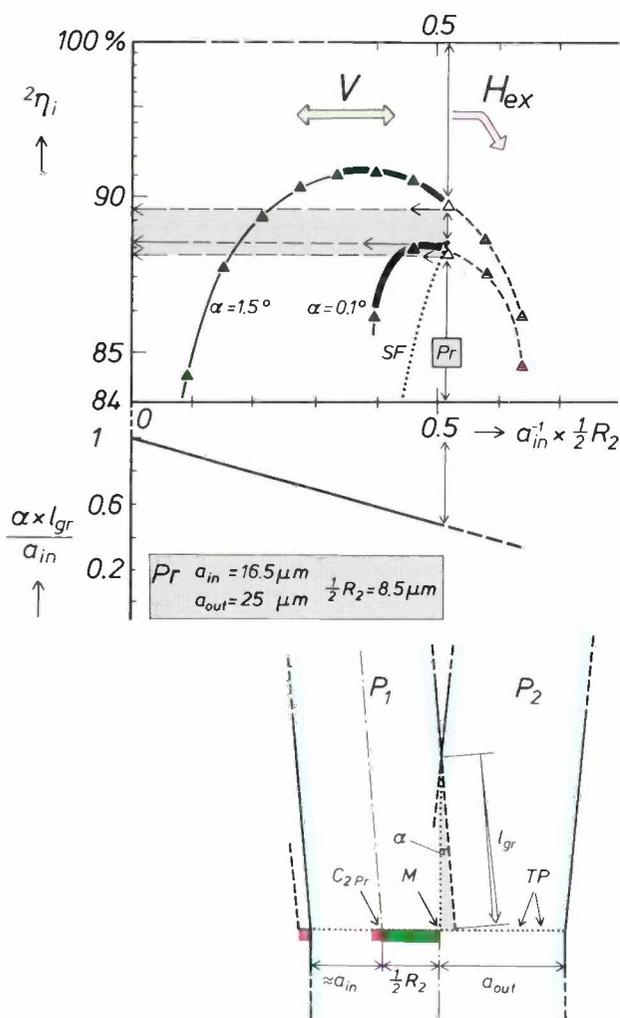


Fig. 10. Efficiency calculations for a butt-joint multiplexer with two input pigtailed (P_1, P_2). The coupling efficiency ${}^2\eta_i$ per input channel i ($= 1, 2$) is found with a computer program for ray tracing that corresponds in general terms to fig. 1c. The longitudinal section shows the fibre cores (with blue edges) in the region of the contact plane (TP). α inclination angle at which core material is polished away from P_1 (and P_2). l_{gr} polished length. a_{in}, a_{out} core radii, assumed to be constant. Pr preferred coupling geometry with characteristic data (grey rectangle); the cores of P_1 and P_2 touch the periphery of the core of the output fibre channel, as shown by the cross-section. Along the horizontal axis $\frac{1}{2}R_2$ varies, which simulates the polishing. The curve in the lower diagram illustrates the linear relation between l_{gr} and R_2 which then change as well. V the region, shown green, where the end faces of the cores of P_1 and P_2 are completely inside TP . H_{ex} the region, shown red, where P_1 and P_2 have some 'overhang', with extra insertion loss. SF calculated coupling efficiency based on cross-sectional deficit (determined by a method corresponding to the method in fig. 8).

zone' V , is found to be subject to two opposing influences. On the one hand the insertion loss decreases, because of better matching of the local numerical apertures; on the other hand the insertion loss increases, since more rays intersect the longitudinal interface (LP in fig. 1a) and are lost. (As R_2 decreases, the polished length l_{gr} increases.) It can be seen from the curves in fig. 10 that when the decrease in R_2 is small the positive effect predominates, but the increase in the insertion loss soon gains the upper hand. The reversal can be

understood from the changing relation between the fixed (spatial) period of the envelope of the radiation and the polished length. If this spatial period is still relatively large (and the perturbing interface is thus relatively small in length) then the positive effect is of greater importance.

The efficiency curve with the smallest inclination angle ($\alpha = 0.1^\circ$) gives the fastest reversal, because the polished-length values are larger by a factor of 15 ($\approx 1.5/0.1$). In this case the difference between the efficiency predicted with SF (88.2%) and the efficiency for the coupling geometry (87.8%) is only 0.4%. (In the limiting case, $\alpha \rightarrow 0$, SF gives the exact result.) For a CM with such a small difference between the input and the output signal directions there is therefore very little difference between the results of ray tracing and calculations based on the cross-sectional deficit. The experimental value will also be much the same, although it would be not quite as good as the value for $\alpha = 1.5^\circ$ [6].

The combined coupling efficiency (radiation source + coupling + CM)

We have now reached a point where we can calculate a combined coupling efficiency for the entire transmitter section of the system (the first main section shown in fig. 3a in part I of this article). The components here are a diode laser, a ball lens (the coupling element), a pigtail and a butt-joint multiplexer. The pigtail should be regarded as part of the CM, as in this part of the article.

The calculation consists in multiplying the curves for the best-case 'total' coupling efficiency (η , see fig. 10 in part I) by the coupling efficiency (${}^2\eta_i$) for an input channel of a CM. The path from the end face of the diode laser to the front face of the pigtail gives the coupling efficiency η (see fig. 6, part I). In this part of the article we have mentioned several times a coupling efficiency of ${}^2\eta_i$. The values of ${}^2\eta_i$ used for the multiplication follow from an equation based on the cross-sectional deficit in a two-channel CM, given by eq. (18).

The result of the multiplication, the combined coupling efficiency ${}^2\eta_{1,c+CM}$, is plotted in fig. 11 as a function of the product $a_{in} \times NA_{in,co}$, which is treated as an independent variable here. The input pigtailed (P_i , see part I, fig. 3a) to which this variable relates were previously characterized as a 'boundary condition' for the system, to be chosen more or less at will by the designer.

The multimode character of the radiation transfer can again be seen from the curves found for the com-

[13] A. J. A. Nicia, Appl. Opt. 22, 1801, 1983.

bined efficiency. The two curves, which are stepped, apply to a strongly astigmatic type of laser.

The figure shows that reducing the astigmatism of the laser can considerably improve the coupling efficiency. As indicated, when currently available fibres are used (*Pr*) the coefficient for the overall transmission loss (${}^2\alpha_i$) will decrease from 2.28 dB to 1.26 dB if the astigmatism can be reduced from 30 μm

to 15 μm . The dashed curve for ${}^2\eta_i$ — which represents the maximum attainable value for the combined coupling efficiency (at $\eta = 100\%$) — then gives a value locally about 0.35 dB better, implying only a slight difference in transmission loss.

The curve for ${}^2\eta_i$ in the formulation just given (eq. 18) does not have the radius-aperture product of fig. 11 as independent variable but $\cos\psi_{\text{max}}$. The angle

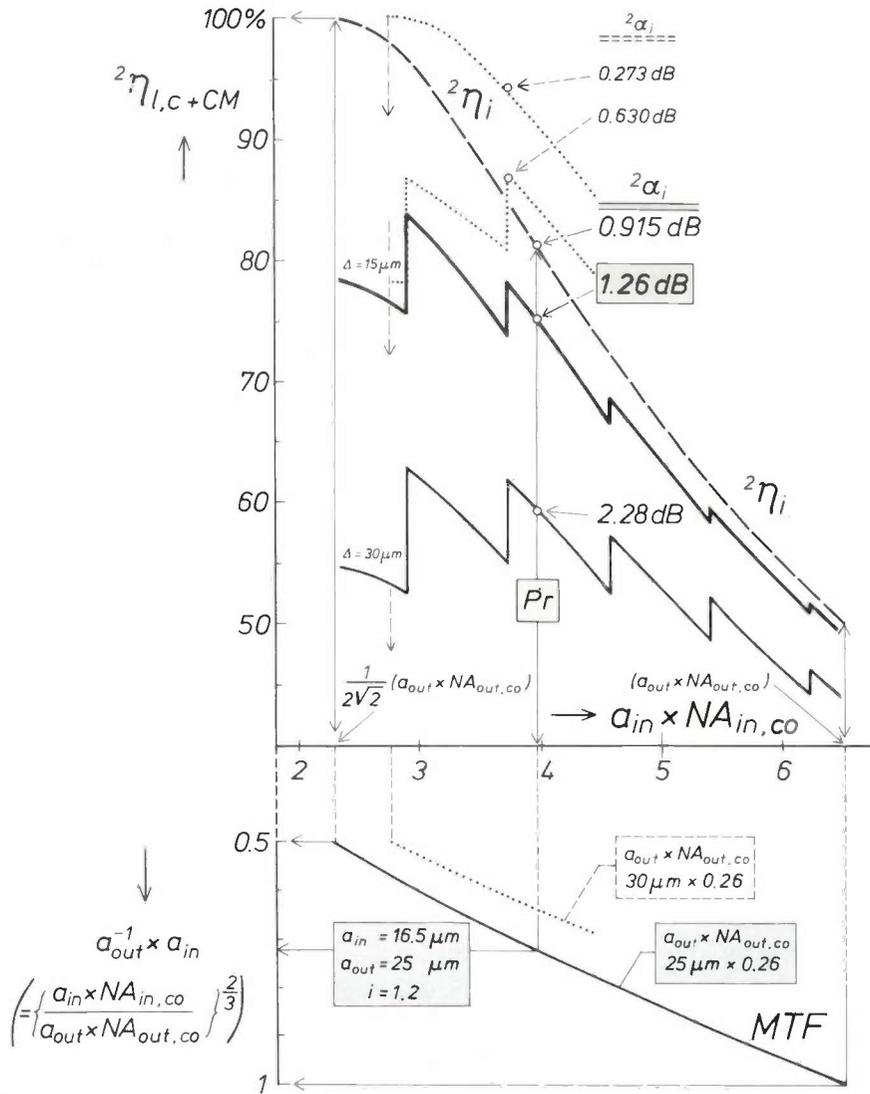


Fig. 11. Combined coupling efficiency ${}^2\eta_{i,c+CM}$, calculated per channel, for the entire transmitter section (*l* diode laser with strong astigmatism, *c* ball lens, *CM* two-channel butt-joint multiplexer including the appropriate input pigtails) for an optical communication system. The first 'free' condition is plotted along the horizontal axis. This is the product $a_{in} \times NA_{in,co}$, where a_{in} is the core radius of the pigtail P_i selected ($i = 1, 2$; see part I, fig. 3a) and $NA_{in,co}$ is its maximum numerical aperture. The efficiency is plotted twice, with the astigmatism of the diode laser as parameter ($\Delta = 15 \mu\text{m}$ and $\Delta = 30 \mu\text{m}$). The dashed curve (${}^2\eta_i$) represents the calculated coupling efficiency ${}^2\eta_{i,c+CM}$ for $\eta = 100\%$ (η relates to the path between the end face of the diode laser and the front face of the pigtail; see part I, figures 6 and 10). ${}^2\alpha_i$, transmission-loss coefficient for input channel *i*, in decibels. *MTF* ancillary curve for radiation transfer satisfying eq. 16a (for maximum 'throughput', the case of equality; see text) and depicting the relation which then exists between a_{in}/a_{out} and $a_{in} \times NA_{in,co}$. The product of radius and aperture at the output side of the multiplexer, $a_{out} \times NA_{out,co}$, is the second 'free' boundary condition; it acts as a parameter of *MTF* and sets the limits of the maximum longest possible path along the $a_{in} \times NA_{in,co}$ -axis. The limits indicated apply to the practical case (*Pr*), where $a_{out} = 25 \mu\text{m}$, $NA_{out,co} = 0.26$ and $a_{in} = 16.5 \mu\text{m}$, $NA_{in,co} = 0.24$. The dashed curves relate to a slightly different choice of the product $a_{out} \times NA_{out,co}$ ($30 \mu\text{m} \times 0.26$); their behaviour suggests a further increase of efficiency.

ψ_{\max} defined in the contact plane TP determines the size of the polished-away part of the core cross-section of an input pigtail, calculated in that plane. The problem is now to link $\cos\psi_{\max}$ and $a_{\text{in}} \times NA_{\text{in,co}}$.

In principle another, more or less open, 'boundary condition' is available to the designer: the radius-aperture product at the output end of the CM, $a_{\text{out}} \times NA_{\text{out,co}}$. The choice of this product can make the desired link. In fig. 11 this choice of $a_{\text{out}} \times NA_{\text{out,co}}$ is evaluated for the case Pr ; the result is the ancillary curve for maximum radiation transfer (MTF), which describes the relation between $a_{\text{out}}^{-1} \times a_{\text{in}}$ and $a_{\text{in}} \times NA_{\text{in,co}}$. Mathematically this relation is given by:

$$a_{\text{out}}^{-1} \times a_{\text{in}} = \left(\frac{a_{\text{in}} \times NA_{\text{in,co}}}{a_{\text{out}} \times NA_{\text{out,co}}} \right)^{2/3}, \quad (20)$$

an expression which is equivalent to the sufficient condition (16a) for unperturbed radiation transfer, provided the marginal case — equality — is used for the inequality. Eq. (20) clearly demonstrates how the choice of the two products radius \times aperture establishes the ratio $a_{\text{out}}^{-1} \times a_{\text{in}}$. The desired connection with the angle ψ_{\max} is finally found by considering that the variations which the designer can in principle adopt will always involve the relation:

$$a_{\text{out}} = a_{\text{in}} (1 + \cos\psi_{\max}), \quad (21)$$

where $0.5 \ll a_{\text{out}}^{-1} \times a_{\text{in}} \ll 1$. Fig. 11 also shows how the second 'boundary condition' ($a_{\text{out}} \times NA_{\text{out,co}}$) also determines the limits on the $a_{\text{in}} \times NA_{\text{in,co}}$ -axis within which there can be said to be a combined efficiency. Variation of this second boundary condition (from 6.5 to 7.8 μm , the dashed case) shows that some gain in efficiency can indeed be found here. In this case the gain is only obtained because the curve for ${}^2\eta_i$ has a better shape.

Optics and efficiency in the prism multiplexer

The path of the rays in a PM with a roof prism suitable for two input channels was discussed in the section dealing with the two types of multiplexer at the beginning of this article. Between the ball lenses (IL and OL , fig. 5) the two beams may be regarded as radiation transmitted inside the guiding cores of two hypothetical pigtails. Immediately behind the prism the emergent radiation forms a spatial pattern which in many respects is not really very different from the distribution of the radiant power over the cross-section TP , the contact plane in the CMs. The efficiency considerations and equations already given for CMs will therefore describe the effects in PMs with two or

four input channels reasonably well, provided a few modifications are made. An example of these modifications is the replacement of the quantities a_{in} and a_{out} in the equations by $f_{\text{IL}} \times NA_{\text{IL}}$ and $f_{\text{OL}} \times NA_{\text{OL}}$.

It will be clear that the refraction angle (γ , fig. 2b) of the prism in the PM is relatively small, say 6.5° in practice. To satisfy eq. (2d) the glass used will have to be such that n_{RP} is approximately equal to 2. The path of rays already constructed in fig. 2b and the equations mentioned there make it possible first of all to calculate the angle of incidence ϕ_{in} , given the distance (2e) between the axes of the input channels and the focal length (f_{IL}) of the ball lens IL . The distance r_{out} mentioned in fig. 2b corresponds in the CMs to the distance $\frac{1}{2}R_2$, as shown for example in fig. 10.

From the choice of r_{out} , with which a_{out} in the CMs is connected, r_{in} can be calculated, again from an equation mentioned in fig. 2b. The equation for s in fig. 2b finally gives the value, determined by r_{in} and ϕ_{in} , that enables the position of the prism to be defined as well.

DEMULTIPLEXING

Classification

As mentioned incidentally in the theoretical treatment, the demultiplexing component, which is part of the receiver section of the system (see fig. 3a, part I), has to be able to act as a spatial filter with the highest possible wavelength selectivity. This is essential because the process of detection — by an avalanche diode or a pin photodiode in each output channel^[14] — cannot distinguish between the different carrier wavelengths in a compound signal. The photodiodes absorb the radiant energy as highly broadband devices. If a signal corresponding to only one wavelength (or colour) is to be detected in each channel, any other signals at other wavelengths in the output channel must first be sufficiently suppressed to allow the desired signal to be detected. The demultiplexers are therefore more complicated in design than our multiplexers.

Since energy losses must be kept to the minimum in the demultiplexers, as in the multiplexers, it seems most appropriate to base filtering by wavelength either on interference or on dispersion (refraction)^[15]. Interference is a field effect, whereas dispersion is a material effect, which occurs when the relevant material constant — the refractive index — depends on the wavelength.

[14] L. J. M. Bollen, J. J. Goedbloed and E. T. J. M. Smeets, The avalanche photodiode, Philips Tech. Rev. 36, 205-210, 1976.

[15] See the article by W. J. Tomlinson quoted in note [3] of part I.

The monochromatic radiation signals that can be obtained in this way are identified by their direction in space. The identification can again be performed effectively in one of two ways, either directly, e.g. by diffraction from a Rowland grating — a simple method — or indirectly by means of radiant power differences due to transmission (and reflection); this is rather more complicated. In addition to the direction of the radiation, the degree of polarization provides a further criterion for wavelength discrimination. This was only considered in passing in our investigation, however.

These two methods of discriminating between monochromatic radiation signals, by direction and by radiant power, and the two different principles of de-

multiplexing, by interference and by dispersion, produce a classification scheme (*Table III*) for the four possible types of demultiplexer. The core component of a type of demultiplexer is stated for each group. *Fig. 12* gives a general picture of four corresponding designs. These four basic types of demultiplexer can also be used in principle for systems with single-mode fibres. In the treatment given here, however, the choice has been limited to multimode fibres.

The prism demultiplexer (RP)

In the first version of a demultiplexer, the prism type (*fig. 12a*), it can be seen that only one lens is required (*IL, OL*). This collimates the input beam and at the same time decollimates the output beams (the pair λ_1 and λ_2). This 'Littrow configuration' is attractive because of its simplicity and its exceptionally small dimensions.

The prism in this configuration does however have to have a different shape from that in a prism multiplexer and it must also have a reflective end face. By virtue of the 'symmetrical folding' of the ray paths inside the prism, a refraction angle ($\frac{1}{2}\gamma$) of only half the original value is sufficient (as indicated by dotted lines in the figure), without any decrease in the total deviation.

The directional discrimination resulting from the deviation must of course offer sufficient space in the image plane of *OL* (the dimension Δx in *fig. 12a*) for the output pigtails (Q_1, Q_2) to be positioned in such a way that each individual signal (λ_1, λ_2) arrives with the minimum insertion loss in the appropriate output pigtail, without significant crosstalk.

The wavelength dependence characterizing the refractive index of the glass was not relevant in the treatment of prism multiplexers, but is essential to the operation of the demultiplexers considered here. It is the main factor determining the distance (Δx) between the end faces of the output pigtails of the demultiplexer, through the equation:

$$\Delta x = f_{OL} \times \int_{\lambda=\lambda_1}^{\lambda_2} \frac{d\phi_{out}(\lambda)}{d\lambda} d\lambda, \tag{22}$$

Table III. Theoretical classification of optical demultiplexers by wavelength-filtering principle and means of discrimination. Four types are possible (*fig. 12*); the respective core components are a roof prism (*RP*), a Rowland grating (*RG*), a glass junction (*GJ*), or an interference filter (*SF*). The type *RG* (grey area) is the best in practice.

Principle Discrimination	Dispersion (material)	Interference (field)
Direction	<i>RP</i>	<i>RG</i>
Radiant power (direction)	<i>GJ</i>	<i>SF</i>

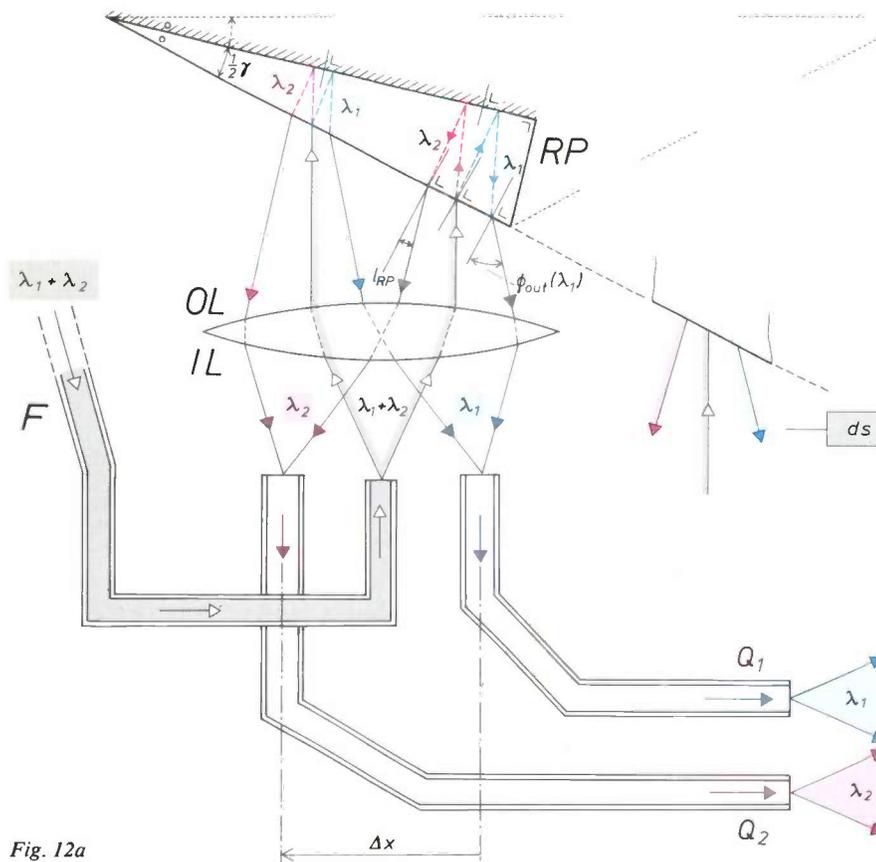


Fig. 12a

where f_{OL} is the focal length of the lens; the derivative $d\phi_{out}(\lambda)/d\lambda$ — the ‘angular dispersion’ — gives the magnitude of the directional change in the output signal, calculated per unit change in wavelength. The direction is defined with respect to the normal (I_{RP}) to the plane of incidence. A total wavelength change of from say λ_1 to λ_2 corresponds to a change in refractive index that in turn establishes the behaviour of the angle ϕ_{out} in eq. (22), so that the distance Δx can then be calculated.

The angular dispersion is the product of two derivatives:

$$\frac{d\phi_{out}(\lambda)}{d\lambda} = \frac{d\phi_{out}}{dn_{RP}} \frac{dn_{RP}(\lambda)}{d\lambda} \quad (23a)$$

The second factor is the purely material-related property of the prism as mentioned earlier; the first is a trigonometric function, which takes a relatively simple form (because of the minimum deviation necessary here as well^[16]):

$$\frac{d\phi_{out}}{dn_{RP}} = 2 \frac{\sin(\frac{1}{2}\gamma)}{\{1 - n_{RP}^2 \sin^2(\frac{1}{2}\gamma)\}^{1/2}} \quad (23b)$$

This result for the Littrow configuration in fig. 12a follows essentially from the connection between the deviation angle, the refraction angle and the refractive index (eq. 2a), mentioned earlier for prism multiplexers.

To find out more easily whether a practical demultiplexer of the type *RP* would have an angular dispersion of interest, it is assumed that the situation in fig. 12a also satisfies (or nearly satisfies) the special symmetry condition:

$$\gamma = \phi_{out} \quad (23c)$$

In the multiplexers, as we have seen, this extra symmetry establishes a relation between the refractive index and the refraction angle (eq. 2d). From this relation and eq. (23b) we can write eq. (23a) in the form:

$$\frac{d\phi_{out}(\lambda)}{d\lambda} = \frac{2\sin(\frac{1}{2}\gamma)}{1 - 2\sin^2(\frac{1}{2}\gamma)} \frac{dn_{RP}(\lambda)}{d\lambda} \quad (23d)$$

which is more suitable for numerical calculation.

Unfortunately, even with highly dispersive types of glass, the angular dispersion in a prism as in fig. 12a remains below 0.1 to 0.2 milliradians per nanometre

[16] Minimizing the deviation reduces the astigmatism in the image projected on to the ends of the output pigtailed.

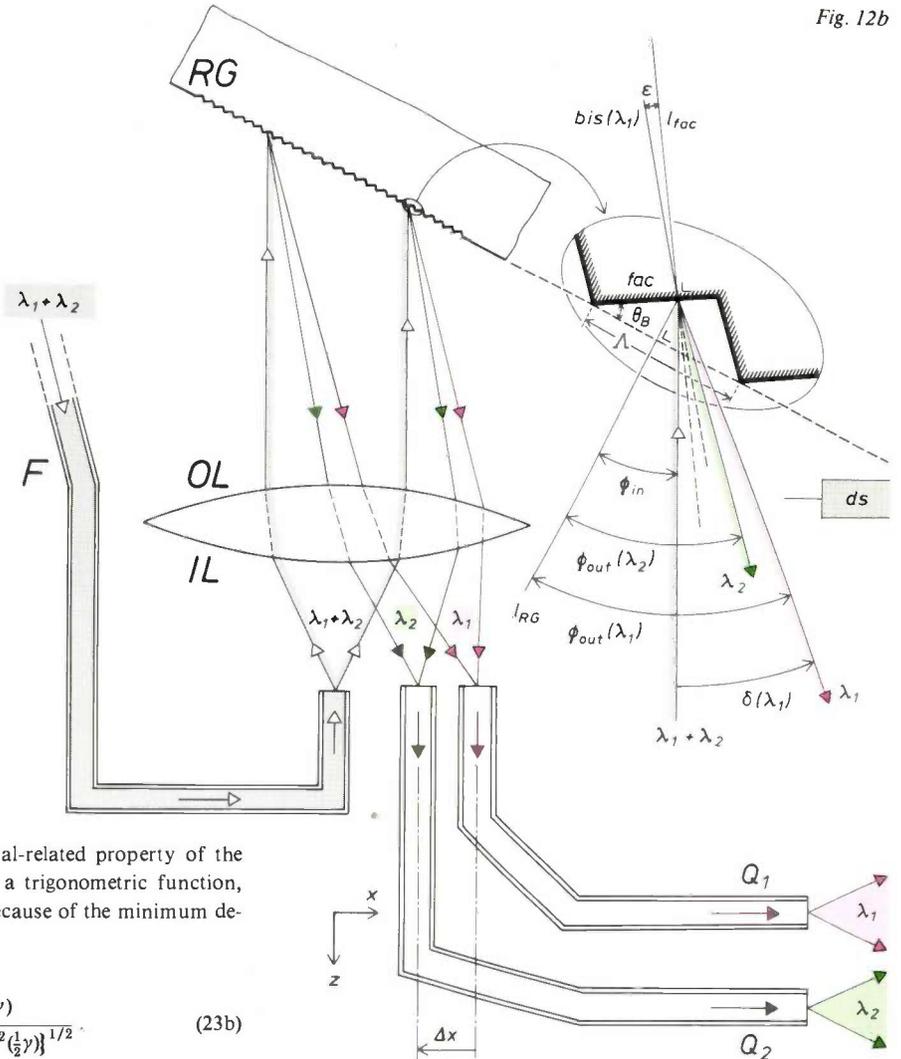


Fig. 12b

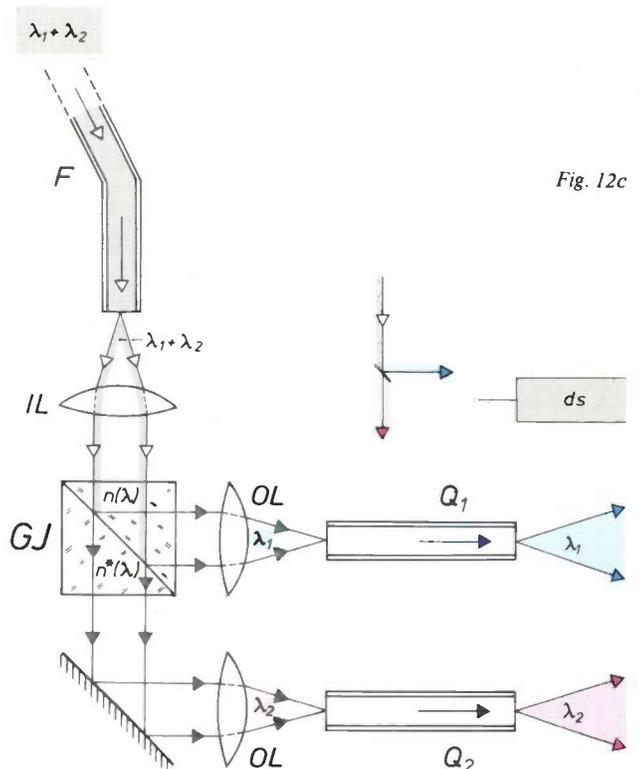


Fig. 12c

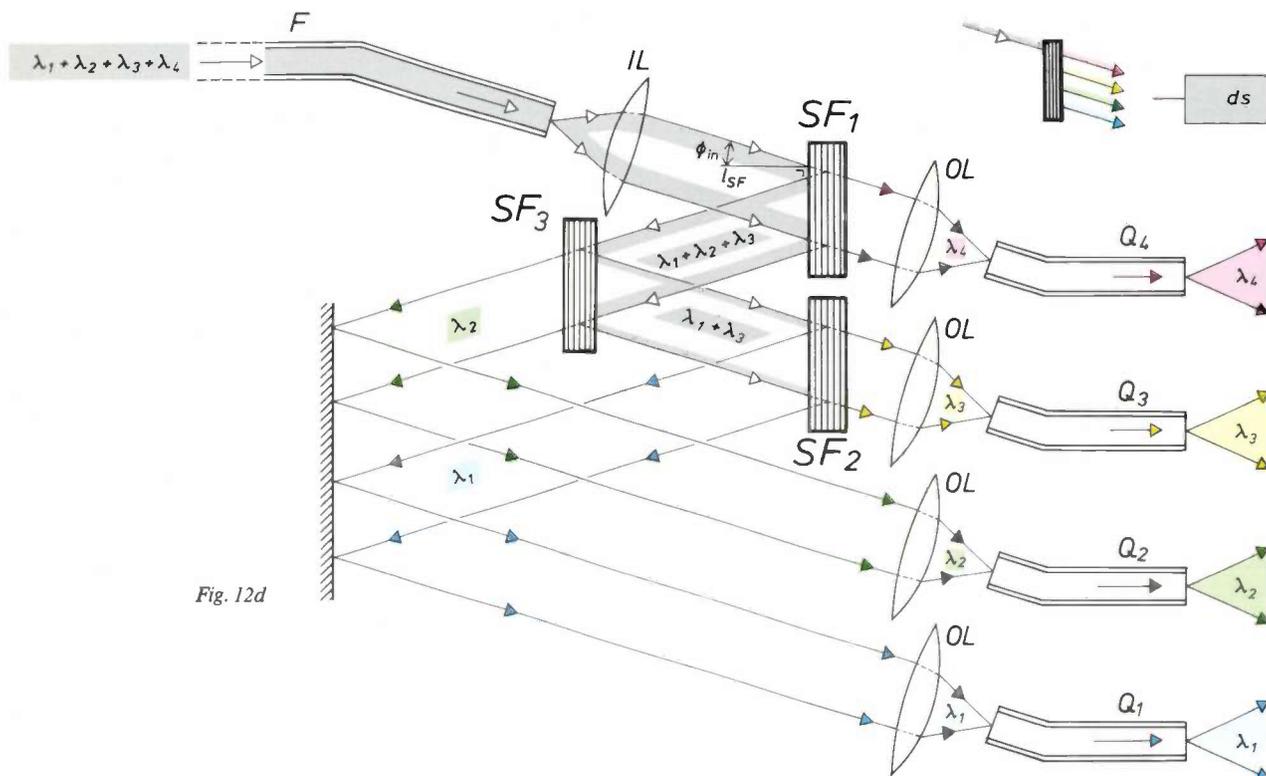


Fig. 12d

Fig. 12. The four types of demultiplexer (symbol *ds*) in the classification scheme given in Table III. The symbol *ds* refers to the directional discrimination selected for the various signal channels, defining the spatial separation function on which the operation of the demultiplexer is based. *a*) Prism demultiplexer. *RP* roof prism. *b*) Ball-lens demultiplexer with grating. *RG* Rowland grating. *c*) Glass-junction demultiplexer. *GJ* glass junction. *d*) Demultiplexer with interference filters. *SF*₁, *SF*₂, *SF*₃ interference filters. *F* input pigtail. *Q*₁, *Q*₂, *Q*₃, *Q*₄ output pigtails. *IL*, *OL* ancillary lenses. For the other symbols see text. The experimental demultiplexers were all of the design type *b*), with four or six output pigtails (fig. 15). The optical axis coincides with the *z*-axis. The centre-points of the entrance planes of the output pigtails (*Q*_{*i*}) all lie on the *x*-axis.

of difference in wavelength. For common diode lasers (see the wavelength data in fig. 4, part I) the difference in wavelength for two signals cannot be more than a few hundred nanometres. If the demultiplexer is not to be unacceptably large, the focal length of the lens is also subject to an upper limit (about 5 mm). Even if a larger demultiplexer could be considered, there would be the further complication that the optics could no longer be just a single element.

Substituting such values in eq. (22), we find that Δx in a demultiplexer of the type *RP* could be no more than a few tens of microns. This seems barely adequate or even too small to leave sufficient room for the output pigtails *Q*_{*i*}. Fortunately the type of demultiplexer with a grating as core component offers much more promise.

The ball-lens demultiplexer with grating (*RG*)

In our study the type of demultiplexer based on diffraction from a Rowland grating (*RG* in Table III) was given a decided preference. The main reason is its better angular dispersion; values of a few milliradians

per nanometre wavelength difference are not difficult to achieve, which promises an improvement with respect to prism demultiplexers of more than one order of magnitude. The Littrow configuration is again used for this version (fig. 12*b*).

As may be seen from the detail enlargement of the front face of the grating in fig. 12*b*, the grooves cut into it, seen in transverse section, all have the same special shape — a tilted and inverted letter V. The grating constant (*A*), the total width per groove, must be of the order of 1 μm (corresponding to about 1000 grooves per mm). The surface in the grooves is carefully finished so that each facet (*fac*) where the radiation ($\lambda_1 + \lambda_2$) to be demultiplexed is incident also acts as a perfect mirror.

All these facet mirrors are at a fixed angle relative to the front face of *RG* — the 'blazing angle' (θ_B). In this way the returning radiation, e.g. the signal at a carrier wavelength λ_1 , can be concentrated in the first main maximum of the entire diffraction spectrum of *RG* for that particular wavelength. For the other main maxima (and the secondary peaks) much less radia-

tion remains, or even none. This arrangement obviously gives greatly improved coupling efficiency at the output side of the demultiplexer; for in a version like that of fig. 12*b*, radiation in the second main maximum, or the third, etc. certainly cannot reach the right pigtails (Q_1 , Q_2) at the output, and would thus be lost.

To describe this special operation of *RG*, three equations are required. The first is the 'grating equation', which establishes the direction of the returning radiation (in this case corresponding to the first main maximum). This direction is expressed as a function of wavelength by:

$$\phi_{\text{out}}(\lambda) = \arcsin\left(\frac{\lambda}{\Lambda} - \sin \phi_{\text{in}}\right). \quad (24a)$$

This equation shows that the relation between the wavelength and the grating constant in fact determines the magnitude of ϕ_{out} — and hence the demultiplexing of the different carrier waves in the incident signal. It can be seen in fig. 12*b* that the two angles, $\phi_{\text{out}}(\lambda)$ and ϕ_{in} , are defined with respect to the normal (l_{RO}) to the front face of *RG*. The incident radiation ($\lambda_1 + \lambda_2$) has the direction of the optical axis of the lens. The angle between this optical axis and the front face of *RG* is equal to $\frac{1}{2}\pi + \phi_{\text{in}}$.

The second equation determines the difference in direction between the incident and the returning radiation:

$$\delta(\lambda) = \phi_{\text{out}}(\lambda) - \phi_{\text{in}}. \quad (24b)$$

The variation of this difference in direction as a function of wavelength is also of crucial importance in calculating the distance between the end faces of the pig-tails (Q_1 and Q_2). In practice, if ϕ_{in} is about 30° , the deviation angle $\delta(\lambda)$ can be nearly 15° . This means that a 'true' Littrow configuration — $\phi_{\text{out}}(\lambda)$ equal to ϕ_{in} — will not do after all.

The third equation shows the extent to which the diffraction of the radiation can become sufficiently 'concentrated', as the ordinary laws of reflection are obeyed as well. The maximum concentration occurs when the bisector, the line *bis*(λ_1) in fig. 12*b*, of the angle between the incident radiation and the returning radiation coincides with the normal l_{fac} to the reflecting facet. The (small) angle ε is a measure of this approximate coincidence, and hence of the maximum concentration of the radiant power of the carrier. The relevant equation is:

$$\varepsilon = \frac{1}{2}\phi_{\text{out}}(\lambda) + \frac{1}{2}\phi_{\text{in}} - \theta_{\text{B}}. \quad (24c)$$

In practice it is found that it is possible at a wavelength of 825 nm, say, to design the demultiplexer in

such a way that the angle ε hardly differs from zero.

The angular dispersion $d\phi_{\text{out}}(\lambda)/d\lambda$ can be calculated from eq. (24a). In the case of the Littrow configuration the calculation is particularly easy, and the result is:

$$\frac{d\phi_{\text{out}}(\lambda)}{d\lambda} = \frac{2}{\lambda} \tan \phi_{\text{in}}. \quad (25)$$

To determine the distance between the end faces of the pigtails Q_1 and Q_2 it is again necessary to perform an integration over the wavelength difference and a multiplication by the focal length of the lens. Although the reality will no doubt differ somewhat from the Littrow configuration, it gives the great advantage of a linear dispersion relation, so that the output pig-tails can be evenly spaced.

Demultiplexers with glass junction or interference filter (GJ, SF)

In the versions *GJ* and *SF* of Table III, radiant-power differences make it possible to discriminate between the various monochromatic radiation signals (fig. 12*c*, *d*). Both types were carefully considered at the start of our study, but as the work proceeded a clear preference developed for the grating type (*RG*) discussed above.

The type with a glass junction^[17] offered little promise of practical implementation at the time; suitable types of glass were not available, and this is still the case. The two refractive indices in fig. 12*c* would have to be so wavelength-dependent that one carrier (λ_1) would be totally reflected, while the other (λ_2) would pass through the glass junction with no reflection at all.

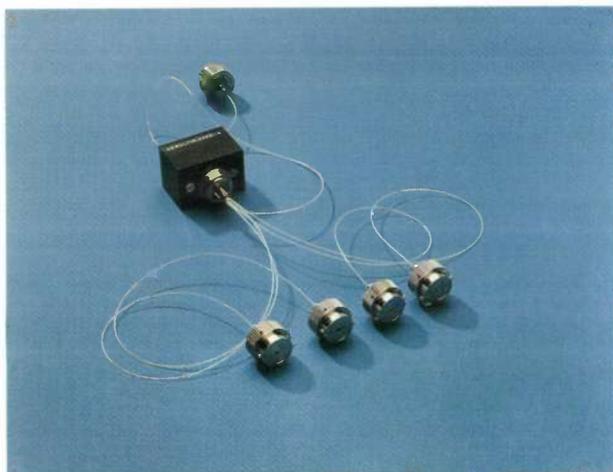
The last type in Table III is based on interference, like the type *RG*. The number of core components (*SF* in fig. 12*d*) can be one less than the number of carriers to be separated. With regard to lenses (*IL*, *OL*) and geometry a number of variations in the configuration are possible, of course, and have indeed been proposed in the literature. The design shown here is therefore only one of many possible versions.

The core component, an interference filter (or stack element)^[10], consists of a stack of thin transparent dielectric layers, with alternate high and low refractive-index values. The total number of layers in a stack can be as many as forty. The layers should not absorb much radiation, so that virtually all of the radiant power should be present at every point where there is both reflection and transmission. Theoretically such a demultiplexer ought therefore to have an efficiency of nearly 100%.

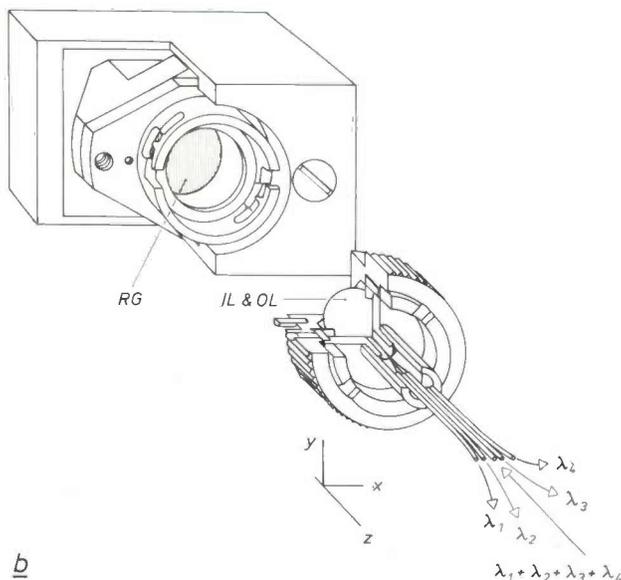
^[17] P. Di Vita, European patent No. 0003 575.

Experimental demultiplexers

For the purposes of this study a number of experimental demultiplexers were made and tested, all of the type *RG* (fig. 12*b*). Most of them had four output pigtails (instead of two), and some even had six^[18]. If larger numbers were used, say ten pigtails, the design principle would lose its attraction. Fig. 13 shows a demultiplexer of the type *RG* suitable for wavelength separation in a signal containing four different carriers. The perspective drawing (fig. 13*b*) gives an idea of



a



b

Fig. 13. *a*) One of the experimental demultiplexers, of the type containing a Rowland grating and a ball lens. The version shown here was designed for demultiplexing a combined signal with four carrier wavelengths. The one input channel and the four output channels are provided with 'half-bayonet' couplings for connection to the fibre in which long-distance transmission takes place and the four signal detectors in the receiver section of the communication system (see part I of this article, fig. 3*a*). *b*) Partly cut-away view of the demultiplexer shown in the photograph. *RG* Rowland grating. *IL & OL* ball lens for collimating the input beam ($\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4$) and at the same time focusing each of the output beams on to an output pigtail (Q_i , see fig. 12*b*). The optical axis lies along the *z*-axis. The four output pigtails lie in the (*x,z*)-plane. The *y*-axis is parallel to the grooves in the surface of *RG*.

some of the details of the construction. The mounted Rowland grating can just be seen in the central module; one half of a bayonet coupling (shown disconnected from the central module) contains the ball lens and the carrier element for guiding the five pigtails. The largest dimension in such demultiplexers is less than 6 cm.

Design and 'boundary conditions'

The fourth type (*SF* in Table III) is not used in practice because of the technological problems. In this type the number of core components increases more or less equally with the intended number of output channels, as we have already seen in fig. 12*d*. The type *RG*, on the other hand, can always function with only one core component (which has the added advantage of being commercially available). Although the type *SF* does have the advantage of a high efficiency, as noted earlier, this could well be cancelled out in practice by the input losses, since losses do increase with the number of core components. Conversely, it therefore seems that the type *RG*, with only one core component, will offer a higher effective efficiency.

In designing the demultiplexers it has to be borne in mind that the choice of diode lasers is limited and that the various carriers may have some instability as a function of time. A third point concerns the slight but unavoidable irregularity in the positioning of the pigtails. The small variation in carrier wavelengths is essentially connected with the characteristics of the diode laser. If the frequency characteristic of the output pigtails in the demultiplexer were perfectly 'flat', this wavelength variation would not cause any increase in radiation losses. It seems that the simplest way to make this characteristic sufficiently flat is to use pigtails for the output channels with an oversized core radius and with a refractive-index profile that is not parabolic but stepped — the centring of the radiation in the pigtails is then less critical, and this results in a sufficiently flat 'passband'.

As mentioned in connection with eq. (25) for the angular dispersion, the required distance Δx between the axes of two neighbouring output pigtails is closely related to the difference $\Delta\lambda$ between the carrier wavelengths assumed in the two pigtails (see fig. 12*b*, $\Delta\lambda = \lambda_2 - \lambda_1$). This distance is proportional to the focal length of the ball lens (*OL*), of course. In fig. 13 the wavelengths assumed were 780 nm, 810 nm, 840 nm and 870 nm; $\Delta\lambda$ is then 30 nm^[18]. The focal length, which is slightly wavelength-dependent, was slightly less than 5 mm. Taking these two values into consideration, we find that for the demultiplexer in fig. 13 a distance Δx of about 200 μm would be necessary if the

operation is to follow the grating equation (24a) reasonably closely. In such a configuration the pigtailed lie close to the optical axis, with the advantage that the aberrations of the ball lens are less important.

In the focal plane of the ball lens an image of the core cross-section of the input fibre (F) will appear. This image will be very little larger than the core cross-section ($2a_{in}$), at least if the radiation is sufficiently monochromatic. If the core diameter ($2a_{out}$) of the output pigtailed were exactly the same as that of F , the frequency characteristic could not have a flat top. Taking a *larger* value for $2a_{out}$ does therefore provide a top with a flat section (the 'passband'). The width ΔB of the passband can be expressed as a fraction of the wavelength difference $\Delta\lambda$, which is given by the simple relation:

$$\frac{\Delta B}{\Delta\lambda} = \frac{2a_{out} - 2a_{in}}{\Delta x} \quad (26)$$

It can be seen that the width of the passband increases as the difference in diameter $2a_{out} - 2a_{in}$ is made larger. A smaller Δx means a larger ΔB . Technically, the smallest practical value of Δx is about $2a_{out}$; the thickness of the cladding of the output pigtailed, which will be at least about $20\mu\text{m}$, would then be completely negligible after all.

The external limit to the increase in a_{out} will eventually be the size of the radiation-sensitive surface of the photodiodes, which act as detectors after the demultiplexer (see fig. 3, part I). This is because the core of the pigtail and the window of the photodiode have to follow the 'trumpet geometry'. Taking the data on available multimode fibres into account (Table I, part I), we may expect $\Delta B/\Delta\lambda$ to have at the most a value of about 0.35. With a channel spacing $\Delta\lambda$ of 30 nm, the passband could then reach a width (ΔB) of 10.5 nm. This seems amply sufficient to eliminate the adverse effect of the variation in wavelength, at least with standard lasers.

In practice the passband of the demultiplexer in fig. 13 is slightly narrower and the width is rather more than 7 nm. This is because one of the design priorities was to follow the grating equation closely. This means that the smallest acceptable value for Δx is not $150\mu\text{m}$ but $200\mu\text{m}$. Keeping close to the grating equation (eq. 24a) implies that the choice of grating constant should correspond to the 'central' wavelength of the demultiplexer (825 nm in the case of fig. 13), while satisfying eqs (24b) and (24c) reasonably well at the same time, of course. Energy losses in the demultiplexer, as well as its sensitivity to the polarization of the passing radiation, can best be minimized in this manner^[19].

Dimensioning and construction

The block diagram in fig. 14a shows that the dimensioning of the demultiplexer was based on the two main design quantities: the width of the passband (ΔB) and the grating constant (Λ). The choice of the diode laser, which was discussed in part I of this article, really decides the first quantity. The choice of grating, which essentially determines the second quantity, is limited, as we saw, to the commercially available types with specially shaped grooves. The ball lens (IL & OL) was selected from Philips types, all made from the same type of optical glass (Schott, LASF9, as mentioned in the section on the fabrication of the prism multiplexer). The block diagram shows that the design is completed when the value of Δx_{min} derived from ΔB and $2a_{out}$ is smaller than the distance Δx given by the angular dispersion ($d\phi_{out}(\lambda)/d\lambda$) and the focal length (f).

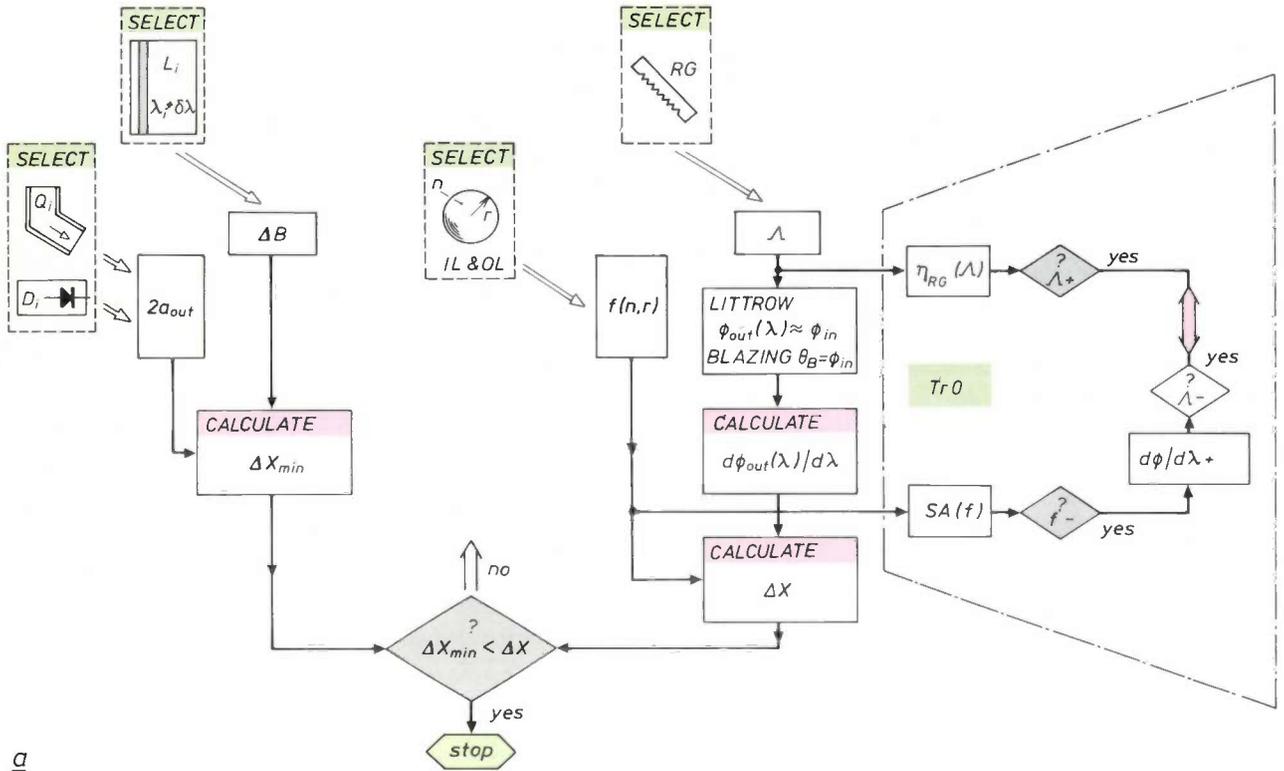
In fact a compromise has to be made, as shown in the section TrO inside the dashed lines, on the right in fig. 14a. This is necessary because it is important to control the spherical aberration (SA) of the ball lens — otherwise the image in the focal plane would be blurred. Making the focal length smaller helps to reduce the spherical aberration. To maintain Δx it is then necessary to increase the angular dispersion correspondingly, which in turn implies making the grating constant *smaller*. This change affects the coupling efficiency (η_{RG}) of the grating (which increases with increasing Λ)^[19]. This conflict of interests has led to a compromise in which greater weight is attached to a good coupling efficiency.

The list in fig. 14b contains a number of design parameters chosen for the experimental demultiplexer in fig. 13. Details of the output pigtailed and the four carrier wavelengths with their maximum acceptable variation are also shown.

Test measurements on the ball lens demonstrated that giving it an antireflection coating reduces the reflection losses from no less than 4×0.39 dB to about 0.1 dB. (The factor 4 points to the fact that the ball lens, owing to the Littrow configuration, introduces four interfaces with reflection losses into the ray diagram.) In particular, reflection of radiation from the front face (IL , see fig. 12b) could have awkward consequences. The arrangement of the five pigtailed is so compact that some of the reflected radiation would be captured directly by the two central output pigtailed, on either side of the input pigtail. Crosstalk would therefore occur in these signal channels. Measure-

[18] H.-G. Finke, A. J. A. Nicia, D. Rittich, *Nachrichtentech. Z.* 37, 346-351, 1984.

[19] E. G. Loewen, M. Nevière and D. Maystre, Grating efficiency theory as it applies to blazed and holographic gratings, *Appl. Opt.* 16, 2711-2721, 1977.



a

grating [a]	Λ	0.83 μm (= 1200 grooves/mm)
	λ_{centr}	0.825 μm
	ϕ_{in}	29°40'
	θ_{B}	36°50'
	η_{RG}	84%
	$d\phi_{\text{out}}(\lambda)/d\lambda$	1.38 mrad/nm

ball lens	r	4.5 mm
	n	1.85
	$f \left(= \frac{n}{2(n-1)} r \right)$	4.9 mm

pigtail Q_i (1st generation multimode)	$2a_{\text{out}}$	100 μm
	cladding	2 x 20 μm
diode laser L_i	λ_i [b]	780, 810, 840, 870 nm

b

[a] made by Bausch and Lomb, Rochester (N.Y.), U.S.A.

[b] acceptable variation: maximum 7 nm.

Fig. 14. a) Calculations for dimensioning type-RG demultiplexers (ball lens with grating, see Table III and figures 12b and 13). The diode lasers (L_i), the output pigtailed detectors (D_i), the Rowland grating (RG) and the ball lens ($IL \& OL$) are all more or less free variables that form the basis of the design. The distance between neighbouring output pigtailed detectors, which determines the dimensioning, has a minimum value (ΔX_{min}), fixed by constructional constraints, and a value (ΔX) stemming from the angular dispersion ($d\phi_{\text{out}}(\lambda)/d\lambda$). Calculation and comparison of the two values leads either to completion ('yes') of the dimensioning, or to the necessity of a repetition ('no'), with modified variables, of course. TrO a compromise, containing the trade-off between coupling efficiency ($\eta_{\text{RG}}(\lambda)$) and spherical aberration ($SA(f)$). The other symbols are defined in the text. b) Selected data for the dimensioning and design of a type-RG demultiplexer with four output channels (see fig. 13).

ments of this unwanted coupling have shown that a total of 9.4% of the power of the reflected radiation appears in the crosstalk (the two central channels taken together). Because of the marked curvature of the reflecting front face, coupling to the two other output pigtailed is impossible, and these are therefore free from crosstalk. If the ball lens did not have an antireflection coating, about 8.6% of the incident radiant power would be lost by reflection from the front face, which corresponds to the loss of 0.39 dB mentioned earlier, and also follows from the calculation of the reflection coefficient for normal incidence^[18]. The crosstalk level in each of the two central output pigtailed would then amount to about 0.4% of the incident radiant power, which means that the antireflection coating gives an improvement of 30 or 40 times.

The demultiplexer was fabricated in much the same way as the prism multiplexer described earlier in this article. The ball lens, for instance, is mounted in a housing much like the half-bayonet coupling described earlier. This part of the demultiplexer can be connected firmly and reproducibly by bayonet catches to the central module, which contains the grating.

As in fig. 5, fig. 13*b* shows only one structural component (the fibre carrier rod) for which an ultrafine finish is essential. The pigtailed are bonded into the various grooves of the carrier rod (fig. 15). The dimensions of the grooves and their spacings determine the spatial configuration of the pigtailed, which leaves the designer a wide margin of freedom to match the demultiplexer to the available lasers (their carrier wavelengths, for example, do not have to be evenly spaced). The end faces of the pigtailed are polished.

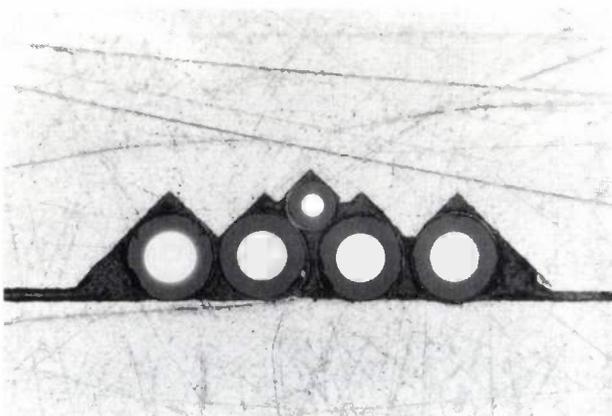


Fig. 15. End face of a fibre carrier rod, in one of the experimental demultiplexers of the type *RG*. The micrograph, made with a scanning electron microscope, shows the one input pigtail, in the centre, and an array of four output pigtailed on opposite sides. The groove structure has an ultrafine finish, to ensure accurate positioning of the fibre end faces. The bonding of the various fibres is also carried out with extreme care. The output pigtailed, including the cladding, have a diameter of 140 μm .

The fibre carrier rod can be adjusted by screws to bring the end faces accurately in the focal plane of the ball lens.

The grating is separately adjustable; it can be pivoted about two axes (x , y) and rotated about the optical axis (z) of the structure. The object of this adjustment is to ensure that the passband is the same (and wide enough) in all the output channels.

Measurements

Measurements were performed on the prototype demultiplexers to obtain some idea of the losses and the separation between the passbands. The band separation determines how selective the operation of the demultiplexer is. The dependence of the results on the ambient temperature was also investigated.

To start with the last point: temperature fluctuations of -20 to $+70^\circ\text{C}$ have very little effect on the losses. Devices for controlling the ambient temperature, which would considerably increase the cost of the system, are therefore unnecessary. Only if a wavelength variation of at least 8 nm caused a laser to start to operate in the steep edge of the passband, well away from the centre, would a temperature rise of about 50°C produce a change of any significance in the losses (up to 0.4 dB). Such a large variation in the carrier wavelength is very unlikely, however.

The radiation source used for the measurements was again a halogen lamp, this time combined with a monochromator of very high quality. The spectral width of the test radiation was set at 0.1 nm, so that the spectral variation was exceptionally small. The method used for the measurements was otherwise not very different from the method described in the section dealing with the loss measurements on the butt-joint multiplexer. Fig. 16 gives an example of a transmission characteristic determined in this way for a demultiplexer of the type *RG* with four output pigtailed. Calculated for each channel, the coupling efficiency is only 1.2 to 1.3 dB below the theoretical maximum of 0 dB. This loss is mainly caused by the grating (0.8 to 1 dB); as we have seen, 0.1 dB can be attributed to the ball lens, and the remainder (0.2 to 0.3 dB) is a reflection loss of the radiant power at the end faces of the output pigtailed. Channel separation is very satisfactory, with the level of crosstalk transmission in the region of -30 dB. In this case also the experimental results agree well with the ray-tracing results (the continuous curves in fig. 16). The ray-tracing results confirmed, incidentally, that the ball lens, provided it has an antireflection coating, does not introduce insertion losses in the passband. This justifies the position given to the calculated passband at the measured level, at -1.2 to -1.3 dB.

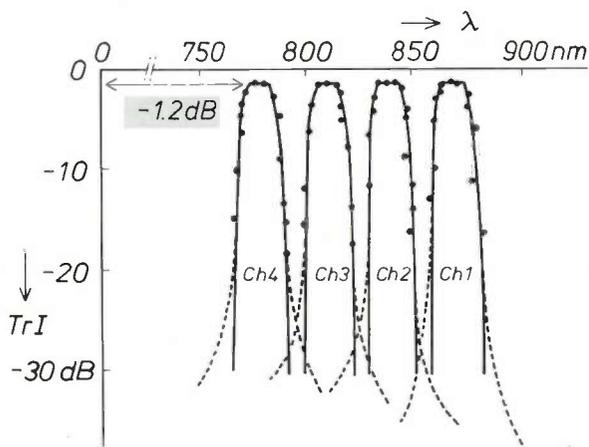


Fig. 16. Example of a transmission characteristic for a demultiplexer of the type *RG* (Table III) with four output pigtails. Vertical: *TrI*, transmission of radiation, expressed in decibels. Horizontal: λ , wavelength of the test radiation, in nm. The experimental results (dashed) agree well with the calculated values (continuous curves). The separation of the four channels (*Ch_i*; $i = 1, 2, 3, 4$) becomes difficult in the region of -30 dB owing to crosstalk. Each channel has a maximum transmission only about 1.2 dB below the theoretical maximum. The ball lens of the demultiplexer is provided with an effective antireflection coating.

The measurements were repeated for the same demultiplexer with an uncoated ball lens. The coupling efficiency was found to be about 2.7 dB below the maximum, and the crosstalk, at least in the two central output pigtails, rose to a level between -20 and -15 dB.

To improve the results as shown in fig. 16 it will be necessary to increase the coupling efficiency (η_{RG}) of the grating. This cannot be done merely by making better gratings. It might be possible to obtain improvements if a fundamental physical change can be accepted, which would undoubtedly complicate the system. One such option would be to introduce linear polarization of the incident radiation [9] [19].

Summary. Part II of this study of a WDM system for optical communication describes colour (or wavelength) multiplexing and demultiplexing components. These comprise butt-joint multiplexers (two and four input channels) and prism multiplexers (two input channels), ball-lens demultiplexers with a grating (four output channels), prism demultiplexers (two output channels) and demultiplexers with glass junction or interference filters. Design specifications, fabrication methods and test measurements are discussed. The maximum dimension of all types is 60 mm; bayonet connectors ensure easy interchangeability. The number of components requiring ultrafine finishing is minimized (1). The coupling efficiency of the butt-joint multiplexers is calculated by two methods (cross-sectional deficit and ray tracing) with aperture matching in the contact plane. The combined efficiency for the transmitter section of the system is also calculated. Comparison of calculated and measured efficiency values shows that the differences are small to negligible. Transmission-loss coefficients (four channels) are 1.5 to 2 dB per input channel; with two input channels about 1 dB is possible for the entire transmitter section. The ball-lens demultiplexers with a grating in a Littrow configuration and with spectrum concentration ('blazing') are the best. Their transmission loss is about 1.2 dB per channel; crosstalk is 30 dB or less. Antireflection coating of the ball lens is essential; temperature stabilization is unnecessary.

Scientific publications

These publications are contributed by staff from the laboratories and other establishments that form part of or are associated with the Philips group of companies. Many of the articles originate from the research laboratories named below. The publications are listed alphabetically by journal title.

Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 5100 Aachen, Germany	A
Philips Research Laboratory, Brussels, 2 avenue Van Becelaere, 1170 Brussels, Belgium	B
Philips Natuurkundig Laboratorium, Postbus 80 000, 5600 JA Eindhoven, The Netherlands	E
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N
Philips Research Laboratories, Cross Oak Lane, Redhill, Surrey RH1 5HA, England	R
Philips Research Laboratories, Sunnyvale P.O. Box 9052, Sunnyvale, CA 94086, U.S.A.	S

J. Benschop & J. Braat	E	Gradient-index objectives for CD applications	Appl. Opt. 26	1195-1200	1987
J. C. Barbour, A. E. T. Kuiper, M. F. C. Willemsen & A. H. Reader	E	Thin-film reaction between Ti and Si ₃ N ₄	Appl. Phys. Lett. 50	953-955	1987
G. W. 't Hooft, W. A. J. A. van der Poel, L. W. Molenkamp & C. T. Foxon	E,R	Near-band-gap luminescence from a GaAs-AlGaAs interface	Appl. Phys. Lett. 50	1388-1390	1987
P. J. Severin	E	Single-mode-fibre passive components made in the fused-head-end technique	Electron. Lett. 23	449-450	1987
K. S. Chung	E	Zero-crossing estimation in signal generation of narrowband digital angle-modulated signals	Electron. Lett. 23	697-699	1987
J. P. Landesman, P. Friedel & M. Taillepiéd	L	XPS study of the chemical cleaning of epitaxial Ga _{0.47} In _{0.53} As(100) surfaces	Europhys. Lett. 3	1143-1149	1987
R. H. Coursant, L. Eyraud*, P. Eyraud*(*INSA, Villeurbanne), M. Fink (Univ. Paris) & J. M. Tellier	L	Lead titanate and lead metaniobate piezoceramics characterization. Application to ultrasonic trans- ducer design	High Tech Ceramics, P. Vincenzini (ed.), Elsevier Science, Amsterdam	1627-1636	1987
G. de With	E	Translucent Y ₃ Al ₅ O ₁₂ ceramics: Something old, something new	High Tech Ceramics, P. Vincenzini (ed.), Elsevier Science, Amsterdam	2063-2075	1987
H. Sari, S. Moridi, L. Desperben (Alcatel Thomson Radioteleph., Genevilliers) & P. Vandamme (CNET, Lannion)	L	Baseband equalization and carrier recovery in digital radio systems	IEEE Trans. COM-35	319-327	1987
A. Fihel & H. Sari	L	Performance of reduced-bandwidth 16 QAM with decision-feedback equalization	IEEE Trans. COM-35	715-723	1987
H. R. Claessen & P. van der Zee	E	An accurate DC model for high-voltage lateral DMOS transistors suited for CACD	IEEE Trans. ED-33	1964-1970	1986
A. J. E. M. Janssen	E	Comments on 'Characterizing the radar ambiguity functions'	IEEE Trans. IT-33	298	1987
I. J. M. M. Raaijmakers, A. H. Reader & H. J. W. van Houtum	E	Nucleation and growth of titanium silicide studied by <i>in situ</i> annealing in a transmission electron microscope	J. Appl. Phys. 61	2527-2532	1987
S. Hasegawa*, T. Tsukao* (*Univ. Kanazawa) & P. C. Zalm	E	Bonding and electronic structures of amorphous Si _{1-x} H _x	J. Appl. Phys. 61	2916-2920	1987

- H. J. W. van Houtum, I. J. M. M. Raaijmakers & T. J. M. Menting (*High Tech. School, Eindhoven*) *E* Influence of grain size in the transformation temperature of C49 TiSi₂ to C54 TiSi₂ J. Appl. Phys. 61 3116-3118 1987
- F. J. A. den Broeder, H. C. Donkersloot, H. J. G. Draaisma* & W. J. M. de Jonge* (**Univ. of Technol., Eindhoven*) *E* Magnetic properties and structure of PD/Co and Pd/Fe multilayers J. Appl. Phys. 61 4317-4319 1987
- K. Erdmann*, P. Deppe*, M. Rosenberg* (**Univ. Bochum*) & K. H. J. Buschow *E* NMR and Mössbauer study of R₂TM₁₄B with R = Y, La, Th, and TM = Fe, Co J. Appl. Phys. 61 4340-4342 1987
- P. J. Kelly & M. S. S. Brooks (*Eur. Inst. for Transuranium Elements, Karlsruhe*) *E* Electronic structure and ground-state properties of the actinide dioxides J. Chem. Soc. Faraday Trans. II 83 1189-1203 1987
- P. J. Schoenmakers, P. E. Rothfus (Philips Centre for Manufacturing Technol., Eindhoven) & F. C. C. J. G. Verhoeven *E* Calculation of pressure, density and temperature profiles in packed-column supercritical fluid chromatography J. Chromatography 395 91-110 1987
- J. P. Farges, C. Schiller & W. J. Bartels *E,L* Growth of large diameter dislocation-free indium phosphide ingots J. Cryst. Growth. 83 159-166 1987
- J. Haisma, B. H. Koek, J. W. F. Maes, D. Mateika, J. A. Pistorius & P. J. Roksnoer *E,H* Hetero-epitaxial growth of GaAs on garnets J. Cryst. Growth 83 466-469 1987
- P. H. L. Notten *E* Hypochlorite reduction at GaAs: A multifaceted reaction J. Electroanal. Chem 224 211-223 1987
- C. Crevecoeur & H. J. de Wit *E* The anodization of heated aluminum J. Electrochem. Soc. 134 808-816 1987
- R. Grössinger*, R. Krewenka*, X. K. Sun*, R. Eibler*, H. R. Kirchmayr* (**Tech. Univ. Wien*) & K. H. J. Buschow *E* Magnetic phase transitions and magnetic anisotropy in Nd₂Ne_{14-x}Co_xB compounds J. Less-Common Met. 124 165-172 1986
- R. Grössinger*, R. Krewenka*, H. R. Kirchmayr* (**Tech. Univ. Wien*), S. Sinnema*, Y. Fu-Ming*, H. Ying-Kai*, F. R. de Boer* (**Univ. Amsterdam*) & K. H. J. Buschow *E* Magnetic anisotropy in Pr₂(Fe_{1-x}Co_x)₁₄B compounds J. Less-Common Met. 132 265-272 1987
- M. Erman, N. Vodjdani, P. Jarry, D. Graziani & H. Pinhas *L* Optical and electrooptical analysis of GaAs inverted rib phase modulators grown by vapor phase epitaxy J. Lightwave Technol. LT-4 1524-1533 1986
- D. B. M. Klaassen, D. M. de Leeuw & T. Welker *E,A* Degradation of phosphors under cathode-ray excitation J. Lumin. 37 21-28 1987
- P. C. M. Gubbens*, A. M. van der Kraan*, J. J. van Loef* (**Interuniv. Reactor Inst., Delft*) & K. H. J. Buschow *E* Crystal field effects in Tm₂M₁₇ compounds studied by ¹⁶⁹Tm Mössbauer spectroscopy J. Magn. & Magn. Mater. 67 255-259 1987
- G. A. C. M. Spierings & J. van Dijk *E* The dissolution of Na₂O-MgO-CaO-SiO₂ glass in aqueous HF solutions J. Mater. Sci. 22 1869-1874 1987
- S. M. Wolfrum *E* Drying and sintering of Al₂O₃ compacts made by sol-gel processing J. Mater. Sci. Lett. 6 706-708 1987
- J. S. Cohen & M. Winnink (*Univ. Groningen*) *E* Stable thermodynamic states J. Math. Phys. 28 1394-1397 1987
- J. Braat *E* Polynomial expansion of severely aberrated wave fronts J. Opt. Soc. Am. A 4 643-650 1987
- W. G. M. van den Hoek, C. A. M. de Vries & M. G. J. Heijman *E* Power loss mechanisms in radio frequency dry etching systems J. Vac. Sci. & Technol. B 5 647-651 1987
- C. Colinet*, A. Pasturel* (**Lab. Thermodyn., Saint Martin d'Hères*) & K. H. J. Buschow *E* Study of the enthalpies of formation in the Gd-(Fe, Co, Pd, Pt) systems Metall. Trans. 18A 903-907 1987
- A. G. van Nie *E* An investigation into the termination resistance of thin film resistors Microelectron. & Reliab. 27 423-428 1987
- P. A. Breddels & J. H. C. Mulkens *E* The determination of the Frank elastic constant for twist deformation of 4'-n-pentyl-4-cyanobiphenyl (5CB) using a conoscope Mol. Cryst. & Liq. Cryst. 147 107-112 1987

- | | | | | | |
|---|------|---|--|-----------|------|
| J. F. M. Westendorp*, F. W. Saris*
(*FOM, Amsterdam), B. Koek,
M. P. A. Vieggers & I. Fenn-Tye
(AERE Harwell, Didcot) | E | Ion beam mixing of thin W, Ta and Au films in Cu | Nucl. Instrum. & Methods Phys. Res. B26 | 539-546 | 1987 |
| A. J. E. M. Janssen | E | On certain integrals occurring in the analysis of a frequency-domain, power-compensated adaptive filter | Philips J. Res. 42 | 131-171 | 1987 |
| J. W. M. Bergmans | E | On the design of smearing filters for ISDN transmission systems | Philips J. Res. 42 | 172-208 | 1987 |
| J. W. M. Bergmans | E | Partial response equalization | Philips J. Res. 42 | 209-245 | 1987 |
| B. de Mooij & K. H. J. Buschow | E | A new class of ferromagnetic materials $RFe_{10}V_2$ | Philips J. Res. 42 | 246-251 | 1987 |
| E. Arnold, H. Baumgart, B. Khan & S. Ramesh | N | Laser-beam-induced recrystallization of silicon and its application to silicon-on-insulator technology | Philips J. Res. 42 | 253-280 | 1987 |
| J. W. M. Bergmans | E | Performance consequences of timing errors in digital magnetic recording | Philips J. Res. 42 | 281-307 | 1987 |
| J. W. M. Bergmans | E | A method of designing robust linear partial response equalizers | Philips J. Res. 42 | 308-338 | 1987 |
| D. B. de Mooij, J. L. C. Daams & K. H. J. Buschow | E | A metastable compound in the Y-Fe-B system | Philips J. Res. 42 | 339-349 | 1987 |
| R. Coehoorn, C. Haas*, J. Dijkstra*, C. J. F. Flipse* (*Univ. Groningen), R. A. de Groot (Univ. Nijmegen) & A. Wold (Brown Univ., Providence, RI) | E | Electronic structure of $MoSe_2$, MoS_2 , and WSe_2 . I. Bandstructure calculations and photoelectron spectroscopy | Phys. Rev. B 35 | 6195-6202 | 1987 |
| R. Coehoorn, C. Haas (Univ. Groningen) & R. A. de Groot (Univ. Nijmegen) | E | Electronic structure of $MoSe_2$, MoS_2 , and WSe_2 . II. The nature of the optical band gaps | Phys. Rev. B 35 | 6203-6206 | 1987 |
| J. G. Kloosterboer & G. F. C. M. Lijten | E | Thermal and mechanical analysis of a photopolymerization process | Polymer 28 | 1149-1155 | 1987 |
| J. B. H. Peek | E | Some features of the Compact Disc digital audio system | Proc. 1st Int. Conf. on Industrial and applied mathematics, Paris-La Vilette, 1987 | 215-234 | 1987 |
| S. Gourrier, P. Rabinzohn & Rocher | C. L | Etat de l'art de la technologie des dispositifs et circuits intégrés GaAs | Proc. 5èmes Journées Nationales Micro-ondes, Nice 1987 | 14-15 | 1987 |
| B. Byzery & M. Binet | L | Caractérisation automatique (puissance, bruit) des transistors hyperfréquences | Proc. 5èmes Journées Nationales Micro-ondes, Nice 1987 | 51-52 | 1987 |
| T. Ducourant | L | Les convertisseurs A/N ultrarapides sur GaAs | Proc. 5èmes Journées Nationales Micro-ondes, Nice 1987 | 77-78 | 1987 |
| M. Roussel | L | Les mesures sous pointes en hyperfréquences: Un outil de caractérisation pour les MMIC | Proc. 5èmes Journées Nationales Micro-ondes, Nice 1987 | 165-167 | 1987 |
| J. W. Broer | E | New horizons for older communicators | Proc. 34th Int. Tech. Commun. Conf., Denver, CO, 1987 | WE85-WE88 | 1987 |
| K. A. Schouhamer-Immink | E | Graceful degradation of digital audio transmission systems | Proc. AES Conv., London 1987 (pre-print) | 15pp. | 1987 |
| P. Boher, F. Pasqualini, J. Schneider & Y. Hily | L | Multipolar plasma passivation scheme for GaAs: <i>In situ</i> ellipsometry control and devices electrical results | Proc. CIPG '87, Antibes 1987 | 120-122 | 1987 |
| P. Autier, J. B. Theeten & J. M. Auger | L | Reactive ion etching of GaAs, AlGaAs and InP monitored by <i>in-situ</i> laser interferometry | Proc. CIPG '87, Antibes 1987 | 158-160 | 1987 |
| P. Boissenot, P. Frijlink, P. Rabinzohn & D. Selle | L | Reactive ion etching of Germanium in a SF_6 discharge: Statistical experimental design approach. Gravure ionique reactive du germanium en plasma SF_6 | Proc. CIPG '87, Antibes 1987 | 167-169 | 1987 |

H. Sari & G. Karam	L	Compensation of modem imperfections by adaptive filtering	Proc. GLOBECOM '86, Houston, TX, 1986	511-516	1986
P. Pesqué, M. Fink (<i>Univ. Paris</i>) & O. Bonnefous	L	Diffraction effects on the measurement of phased array element acousto-electric response	Proc. IEEE Ultrasonics Symp., Williamsburg, VA, 1986	669-674	1986
O. Bonnefous, P. Pesqué & X. Bernard	X. L	A new velocity estimator for color flow mapping	Proc. IEEE Ultrasonics Symp., Williamsburg, VA, 1986	855-860	1986
J. P. André, J. N. Patillon, J. L. Gentner, E. P. Menu, D. Moroni & G. M. Martin	L	High quality lattice matched InGaAs/InP heterostructures prepared by atmospheric pressure MOVPE for high speed photodetectors	Proc. Int. Symp. GaAs and Related Compounds, Las Vegas, NV, 1986	417-422	1987
C. Rocher, J. Maluenda, B. Gabillard, T. Ducourant, M. Prost & M. Rocchi	B. L	Evaluation of the theoretical maximum fabrication yield of GaAs 1K bit SRAM's	Proc. Int. Symp. GaAs and Related Compounds, Las Vegas, NV, 1986	509-514	1987
V. Bodart, P. Houdy, J. P. Weber, C. Hily & P. Ruterana	L	C-W multilayers for soft X-rays mirrors. Fabrication using <i>in situ</i> kinetic ellipsometry	Proc. Int. Symp. on Trends and new applications in thin films, Strasbourg 1987	15-20	1987
P. van de Weijer & R. M. M. Cremers	M. E	Pulsed optical pumping in low-pressure mercury discharges	Radiative processes in discharge plasmas, J. M. Proud & L. H. Luessen (eds), Plenum, New York	65-93	1986
D. R. Wolters & H. L. Peek	E	Fowler-Nordheim tunneling in implanted MOS devices	Solid-State Electron. 30	835-839	1987
P. W. J. M. Boumans & J. J. A. M. Vrakking	M. E	Detection limit including selectivity as a criterion for line selection in trace analysis using inductively coupled plasma-atomic emission spectrometry (ICP-AES) — a tutorial treatment of a fundamental problem of AES	Spectrochim. Acta 42B	819-840	1987
R. J. M. Zwiers, J. J. M. Braat & G. C. M. Dortant	E	A replicated bi-aspherical readout lens for optical disc systems	SPIE 645	53-57	1986
P. J. Schoenmakers & F. C. C. J. G. Verhoeven	E	Supercritical-fluid chromatography — prospects and problems	Trends Anal. Chem. 6	10-17	1987
J. G. Siekman	E	Polishing of brittle amorphous alloy	Wear 117	359-374	1987

Contents of Philips Telecommunication and Data Systems Review 45, No. 3, 1987

- H. van Kampen: Meeting point PABX (pp. 2-3)
- G. M. J. Havermans & R. J. Mulder: SOPHO-S, an IS-PABX based on ISDN standards (pp. 4-16)
- G. M. J. Havermans & E. Radant: SOPHO-S2500 in the German ISDN-pilot project (pp. 17-22)
- J. P. Maat: The application of SOPHO-S2500 in fully integrated networks (pp. 23-27)
- J. M. Broekhuizen: DPNSS 1 in SOPHO-S (pp. 28-34)
- P. B. Hesdahl: The digital voice/data-terminal SOPHO-SET S375 D (pp. 35-41)
- P. C. Hubers: New features for analogue SOPHO-S telephone sets (pp. 42-47)
- D. Anderson: SOPHO-K12: a high performance low cost key telephone system (pp. 48-53)
- C. M. Klik: Support for X.21 Data Switching with SOPHO-S (pp. 54-58)
- F. J. L. Dams: Data switching experience in a laboratory network (pp. 59-63)
- D. van der Lugt: SOPHO-TEXT switch in communication networks (pp. 64-68)
- E. Grohmann, W. Lösel & W. Stahl: Radio Base Station RBS 901 for the Nordic Mobile Telephone system (pp. 69-81)
- R. Brie & O. Le Riche: A new generation of Philips modems (pp. 82-88)
- G. Schipper: HCS115 Videotex terminal (pp. 89-96)