

PHILIPS TECHNICAL REVIEW

VOLUME 40
1982



PHILIPS

Philips Technical Review is published by Philips Research Laboratories, Eindhoven, The Netherlands, and is devoted to the investigations, processes and products of the laboratories and plants that form part of or cooperate with enterprises of the Philips group of companies. In the articles the associated technical problems are treated along with their physical or chemical background. The Review covers a wide range of subjects, each article being intended not only for the specialist in the subject but also for the non-specialist reader with a general technical or scientific training.

The Review appears in an English and a Dutch edition, both identical in contents. There are twelve numbers per volume, each of about 32 pages. A yearly index is given for each volume and an index covering ten volumes appears every five years (the latest one is included in this volume).

Editors: Dr J. W. Broer
Dipl.-Phys. R. Dockhorn, Editor-in-chief
Dr E. Fischmann
Ir J. A. Klaassen
T. M. B. Schoenmakers, assistant editor
Dr J. L. Sommerdijk
Dr J. Ubbink
Ir F. Zuurveen

English edition: D. A. E. Roberts, B.Sc., M. Inst. P., A.I.L. (Redhill, Surrey)

© N.V. Philips' Gloeilampenfabrieken, Eindhoven, The Netherlands, 1983.
Articles may be reproduced in whole or in part provided that the source 'Philips Technical Review' is mentioned in full; photographs and drawings for this purpose are available on request. The editors would appreciate a complimentary copy.

Subject index, Volume 40, 1982

	Page		Page
Abstraction	225	O-BUS system for flexible public transport by on-call buses	231
Alkali-antimonide films for photocathodes, growth of	19	Oscilloscope, 60 MHz, with P ² CCD for digital image storage	55
Assembly robot, experimental	33	Parallel programs	254
Audio long-play disc, optical, see Compact Disc Digital Audio		P ² CCD in 60 MHz oscilloscope with digital image storage	278
Bits for thermocompression bonding, miniature	200	P ² CCD in 60 MHz oscilloscope with digital image storage	55
Buses, on-call, for flexible public transport	231	PHIDAS	245
CAD/CAM, data-base management with PHIDAS system	245	PHIDIAS	262
Camera window for ultrasoft X-rays from space	12	Phosphors for colour TV, pigmentation of	48
Chemical vapour deposition of wear-resistant coatings on tool steel	204	Photocathodes, growth of alkali-antimonide films for	19
Compact Disc Digital Audio:		Pigmentation of phosphors for colour television	48
editorial	149	Printed boards, attachment of leadless electronic components to	342
general	151	Printing, ink-jet	192
system aspects and modulation	157	Quartz crystals, TTC-cut	1
error correction and concealment	166	Raman spectrometric analysis of glass fining	310
digital-to-analog conversion	174	Refrigerator-freezer with heat pipe	350
Curvature of reflecting surfaces, instrument for measuring	338	Resonators with TTC-cut quartz crystals	1
D.C. motor, linear, with permanent magnets	329	Robot, experimental assembly	33
Diamond, spark machining of	202	Software (special issue):	
Diamond die	133	editorial	217
Digital-to-analog converter for Compact Disc	174	introductory article	219
Distributed computations on arrays of processors	270	abstraction	225
Echo reduction in Teletext by automatic equalizer on single chip	319	O-BUS: a system for flexible public transport by means of on-call buses	231
Electrochemiluminescence in electrolyte-free solutions	69	INDA, a software tool for the production engineer	237
Electrostriction, electromechanical transducers based on	358	a data-base management system for CAD and CAM	245
Fining of glass	310	parallel programs	254
Glass fibre, electron-microprobe analysis of	349	software aspects of the PHIDIAS system	262
Glass-fibre circuits, components for	46	distributed computations on arrays of processors	270
Glass fining	310	transformation methods for improving parallel programs	278
Heat pipe in refrigerator-freezer	350	see also Telephone cables	
Heat pipe in solar collector	181	Solar collector, evacuated tubular, with heat pipe	181
High-speed spark-machining equipment, three special applications	199	Spark-machining equipment, high-speed, three special applications	199
Holst, Gilles, pioneer of industrial research in the Netherlands	121	Speech: studies with SPARX system, and manipulation of speech sounds	134
INDA, a software tool for the production engineer	237	Spheres, aluminium, made by spark machining	202
Ink-jet printing	192	Surfaces, reflecting, instrument for measuring curvature of	338
Intonation control in synthetic speech	134	Telephone cables, multiwire, computer-aided research on	85
Lacquers, photopolymerizable, for LaserVision discs	298	Teletext, echo reduction by automatic equalizer on single chip	319
LaserVision discs:		Thin-layer cells for electrochemiluminescence	69
manufacture by photopolymerization process	287	Tool steel, application of wear-resistant coatings by CVD	204
photopolymerizable lacquers for	298	Transducers, electromechanical, with no hysteresis	358
Linear d.c. motor with permanent magnets	329	Transformation methods for improving parallel programs	278
Magnetic fluid, making tracks on video tape visible with	129	TRAPATT oscillator	99
Microprobe, electron, analysis of glass fibre	349	TTC-cut quartz-crystal resonators	1
Microwave measurement of moisture content of process materials	112	Video disc, long-play, see LaserVision discs	
Microwave TRAPATT oscillator	99	Video tape, making tracks visible with magnetic fluid	129
Moisture content of process materials measured by microwave method	112	Wear-resistant coatings on tool steel, deposition by CVD	204
		Window for a camera for ultrasoft X-rays from space	12
		Wire-drawing die	133

Author index, Volume 40, 1982

	Page		Page
Asselman, G. A. A. and A. J. van Mensvoort		Kruseman Aretz, F. E. J.	
A refrigerator-freezer with heat pipe	350	Abstraction	225
Barth, P. J. , see Voorman, J. O.		Legierse, P. E. J. , see Haverkorn van Rijsewijk, H. C.	
Bleeker, J. A. M. , W. H. Diemer, A. P. Huben and H. Huizenga		Lippits, G. J. M. , see Kloosterboer, J. G.	
Camera window for ultrasoft X-rays from celestial sources	12	Martin, A. J.	
Bloem, H. , J. C. de Grijs and R. L. C. de Vaan		Distributed computations on arrays of processors	270
An evacuated tubular solar collector incorporating a heat pipe	181	—, see Kessels, J. L. W.	
Brice, J. C. and W. S. Metcalf		Meinders, H. C. , see Kloosterboer, J. G.	
Quartz-crystal resonators using an unconventional cut	1	Melis, J. H. A.	
Burnett, D. J.		O-BUS: a system for flexible public transport by means of on-call buses	231
INDA, a software tool for the production engineer	237	Mensvoort, A. J. van , see Asselman, G. A. A.	
Carasso, M. G. , J. B. H. Peek and J. P. Sinjou		Metcalf, W. S. , see Brice, J. C.	
The Compact Disc Digital Audio System	151	Meyer, K. H. , see Honds, L.	
Carl, K. , J. A. M. Dikhoff and W. Eckenbach		Meyer, W. and W. Schilz	
The pigmentation of phosphors for colour television	48	Microwave measurement of moisture content in process materials	112
Casimir, H. B. G.		Newton, B. H. , see Davies, R.	
Gilles Holst, pioneer of industrial research in the Netherlands	121	Nicia, A. J. A. and C. J. T. Potters	
Davies, R. , B. H. Newton and J. G. Summers		Components for glass-fibre circuits	46
The TRAPATT oscillator	99	Nie, A. G. van , see Foederer, A. F.	
Diemer, W. H. , see Bleeker, J. A. M.		Nooteboom, S. G. , see Hart, J. 't	
Dikhoff, J. A. M. , see Carl, K.		Peek, J. B. H. , see Carasso, M. G.	
Dolizy, P.		Plassche, R. J. van de , see Goedhart, D.	
Growth of alkali-antimonide films for photocathodes	19	Potters, C. J. T. , see Nicia, A. J. A.	
Dollekamp, H. , L. J. M. Esser and H. de Jong		Rijckaert, A. M. A.	
P ² CCD in 60 MHz oscilloscope with digital image storage	55	Making the tracks on video tape visible with a magnetic fluid	129
Döring, M.		Scha, R. J. H.	
Ink-jet printing	192	Software	219
Eckenbach, W. , see Carl, K.		Schaper, H. , H. Köstlin and E. Schnedler	
Esser, L. J. M. , see Dollekamp, H.		Electrochemiluminescence in electrolyte-free solutions	69
Fischer, W. E.		Schilz, W. , see Meyer, W.	
A data-base management system for CAD and CAM	245	Schnedler, E. , see Schaper, H.	
Foederer, A. F. , J. L. M. Hagen and A. G. van Nie		Schnell, A.	
An instrument for measuring the curvature of reflecting surfaces	338	Electromechanical transducers with no hysteresis	358
Goedhart, D. , R. J. van de Plassche and E. F. Stikvoort		Schouhamer Immink, K. A. , see Heemskerk, J. P. J.	
Digital-to-analog conversion in playing a Compact Disc	174	Sinjou, J. P. , see Carasso, M. G.	
Grijs, J. C. de , see Bloem, H.		Sintzoff, M.	
Hagen, J. L. M. , see Foederer, A. F.		Transformation methods for improving parallel programs	278
Hanenberg, J. G. van den and J. Vredenburg		Snijder, P. J. , see Voorman, J. O.	
An experimental assembly robot	33	Stikvoort, E. F. , see Goedhart, D.	
Hart, J. 't , S. G. Nooteboom, L. L. M. Vogten and L. F. Willems		Straten, P. J. M. van der and G. Verspui	
Manipulation of speech sounds	134	Chemical vapour deposition of wear-resistant coatings on tool steel	204
Haverkorn van Rijsewijk, H. C. , P. E. J. Legierse and G. E. Thomas		Summers, J. G. , see Davies, R.	
Manufacture of LaserVision video discs by a photopolymerization process	287	Thomas, G. E. , see Haverkorn van Rijsewijk, H. C.	
Heemskerk, J. P. J. and K. A. Schouhamer Immink		Timmermans, J. , see Hoeve, H.	
Compact Disc: system aspects and modulation	157	Vaan, R. L. C. de , see Bloem, H.	
Hoeve, H. , J. Timmermans and L. B. Vries		Veldhuis, J.	
Error correction and concealment in the Compact Disc system	166	Computer-aided research on multiwire telephone cables	85
Honds, L. and K. H. Meyer		Verspui, G. , see Straten, P. J. M. van der	
A linear d.c. motor with permanent magnets	329	Verweij, H.	
Huben, A. P. , see Bleeker, J. A. M.		The fining of glass	310
Huizenga, H. , see Bleeker, J. A. M.		Vledder, H. J. , see Klein Wassink, R. J.	
Immink, K. A. Schouhamer , see Heemskerk, J. P. J.		Vogten, L. L. M. , see Hart, J. 't	
Jong, H. de , see Dollekamp, H.		Voorman, J. O. , P. J. Snijder, J. S. Vromans and P. J. Barth	
Kessels, J. L. W. and A. J. Martin		An automatic equalizer for echo reduction in Teletext on a single chip	319
Parallel programs	254	Vredenburg, J. , see Hanenberg, J. G. van den	
Klein Wassink, R. J. and H. J. Vledder		Vries, L. B. , see Hoeve, H.	
The attachment of leadless components to printed boards	342	Vromans, J. S. , see Voorman, J. O.	
Kloosterboer, J. G. , G. J. M. Lippits and H. C. Meinders		Waumans, B. L. A.	
Photopolymerizable lacquers for LaserVision video discs	298	Software aspects of the PHIDIAS system	262
Köstlin, H. , see Schaper, H.		Wijers, J. L. C.	
		Three special applications of the Philips high-speed spark-machining equipment	199
		Willems, L. F. , see Hart, J. 't	

Quartz-crystal resonators using an unconventional cut

J. C. Brice and W. S. Metcalf

The first quartz crystals to be used for frequency stabilization were predominantly 'Y cuts': they were cut perpendicularly to the y-axis of the crystal. The 'AT', 'BT' and other cuts were discovered later as a result of the search for orientations for which the frequency is immune to temperature changes. These cuts are 'singly' rotated with respect to the Y-cut. By cutting in a 'doubly rotated' orientation, one extra degree of freedom is obtained and more requirements can be met. The 'TTC cut' discussed in this article is a doubly rotated cut. The authors show that the use of such crystals provides a simple and inexpensive means of frequency stabilization that gives an excellent performance even under highly adverse conditions.

Introduction

In many varieties of electronic equipment, ranging from clocks and counters through radio, telephone, television and computers to satellite navigation systems, quartz resonators are used for highly accurate selection of the frequency of oscillators and filters. In the 1920s, when quartz was first used for this purpose, an accuracy of 1 part in 10^4 or 10^5 could be obtained — one or two orders of magnitude better than the conventional *LC* circuits could provide. Since then quartz resonators have steadily been improved and accuracies of 1 part in 10^8 are now quite feasible under good conditions. A clock of this accuracy would be right to one second in three years. By taking elaborate precautions significantly greater accuracy is attainable. This article describes devices that allow great accuracy to be obtained in simple systems under adverse conditions.

A quartz resonator in an oscillator circuit locks the frequency of oscillation to the resonant frequency of the quartz plate. The accuracy of locking is determined by the *Q* (quality factor) of the resonator: short-term frequency fluctuations are inversely proportional to the *Q*. For long-term stability, however, the stability

of the resonant frequency itself, under varying operating conditions, is of more importance. It is on this point that significant advances have been made in the last ten years.

The main factor affecting the resonant frequency is usually a change of temperature. In the 1930s it was found that the temperature dependence of the frequency depends on the orientation in which the quartz plate is cut from the crystal. Among the many cuts that have since been investigated, the 'AT cut' has maintained its dominant position, with both the first and the second temperature coefficients equal to zero at room temperature. Thus, in a range of some 30 °C around room temperature an AT-cut crystal changes its resonant frequency by only about 4 parts in 10^6 . However, this is not good enough for many modern applications. A mobile transmitter-receiver system, for instance, may have to function in a Canadian winter at -30 °C or at noon in the Sahara at +50 °C. An AT crystal changes its frequency by more than 1 part in 10^5 in this temperature range, and other cuts do not do any better. Thus it is necessary either to measure the temperature and apply a frequency correction, or to prevent the temperature of the crystal from changing. The second method is easier. A simple and widely applied scheme is to use a small oven that

J. C. Brice, M.A., Ph.D., is with Philips Research Laboratories (PRL), Redhill, Surrey, England; W. S. Metcalf, M.A., B.Sc., is with Cathodeon Crystals Ltd, Linton, Cambridge, England.

always keeps the crystal at a constant temperature well above the highest temperature required, for instance at 90 °C in a range of 4 °C; the complication of having to cool the crystal is thus avoided. 'Adapted' AT crystals are often used in this way. By slightly changing the angle of cut from the ideal AT one, it is possible to make the first temperature coefficient zero at say 90 °C. The second coefficient, however, does not then vanish at that temperature.

In this article we describe the results of our investigations on resonators with an unconventional cut that appears to be much better suited for oven-controlled oscillators: the TTC cut (TTC = thermal-transient compensated). This cut is one of many that have been proposed recently [1]. It is closely related to the 'SC cut', which offers 'strain compensation'. Quartz is such a well-investigated material by now that the properties of new cuts can be readily predicted by computer calculations. A series of calculations made at Philips Research Laboratories, Redhill (PRL), have confirmed earlier claims made for the TTC cut and have suggested other benefits.

The benefits predicted are, in short, that the frequency is virtually independent of temperature in the range 80-100 °C and of mechanical strains in the crystal. The strain compensation eliminates the frequency transients due to thermal strains, which to some degree always occur in the temperature cycling that goes with oven-control. TTC crystals, however, differ from AT crystals in many other aspects. For instance, they vibrate mechanically in a different way and the electrical characteristics are different; the effects of cutting, lapping and etching also depend on the angles of cut. An experimental investigation was therefore required to determine whether any of these effects would prevent the predicted advantages from being realized in practice. The answer is no, as was shown by a joint programme of research, started in April 1978 by PRL and Cathodeon Crystals Ltd, Linton; some of the funds and facilities were provided by the UK Government. One and a half years later, in November 1979, the production of a family of TTC devices was started.

In this investigation research and technology have gone hand in hand. In order to carry out satisfactory experiments, it was necessary first to solve the problem of crystal manufacture. During the early stages of the joint programme Cathodeon Crystals produced over 500 crystal units to establish the required technology. The new cut is a 'doubly rotated' cut, rather than the 'singly rotated' AT cut (*fig. 1*). This means that, once the crystal axes x , y , z have been established, there are two angles (ϕ and θ) to be adjusted instead of one (θ) before plates can be cut from the crystal. This required

a complete redesign of the jigs in the cutting machine. The new designs [2] were however immediately useful for small-scale production, and once the advantages of the TTC cut were established they could readily be applied to large-scale production, which was therefore started immediately.

Before discussing our results we shall discuss the general physical and technical background of quartz resonators [3] and pay specific attention to the AT and TTC cuts [1].

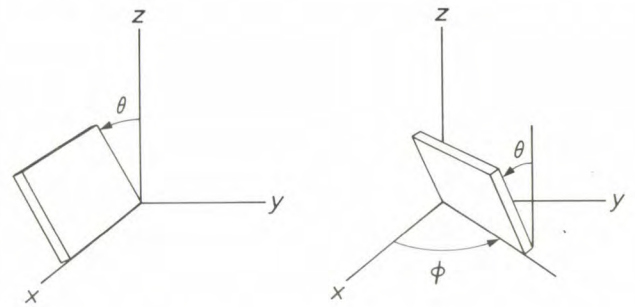


Fig. 1. Singly and doubly rotated cuts in quartz. The axes x , y , z are the crystal axes as usually defined in quartz (*fig. 2*). A singly rotated cut (*left*) is obtained by rotating the 'Y cut' — the cut perpendicular to the y -axis — about the x -axis; a doubly rotated cut by rotating a Y cut first about the z -axis and then about its new intersection with the x, y -plane (the x' -axis). The conventional AT cut is singly rotated ($\theta = 35.35^\circ$), the new TTC cut is doubly rotated ($\phi = 21.90^\circ$, $\theta = 33.93^\circ$).

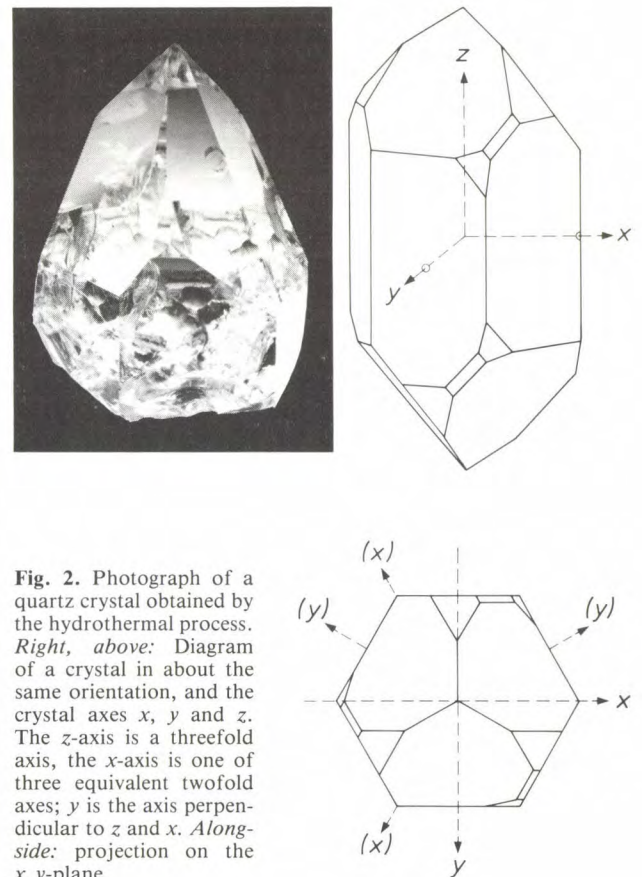


Fig. 2. Photograph of a quartz crystal obtained by the hydrothermal process. *Right, above:* Diagram of a crystal in about the same orientation, and the crystal axes x , y and z . The z -axis is a threefold axis, the x -axis is one of three equivalent twofold axes; y is the axis perpendicular to z and x . *Alongside:* projection on the x, y -plane.

Quartz resonators

The two main properties that have made quartz such a successful material are its piezoelectricity, allowing mechanical vibrations to be coupled to electrical vibrations, and its ability to sustain mechanical vibrations with very little attenuation.

The piezoelectricity of quartz is directly related to its crystal structure. Fig. 2 shows a crystal as it is grown by the hydrothermal process^[4], and a schematic drawing of a crystal with the x -, y - and z -axes that are normally used. The z -axis is a threefold axis, the x -axis is one of three equivalent two-fold axes. Because the z -axis is threefold only, the structure has no centre of symmetry, so that it is possible to distinguish say the $+x$ -axis from the $-x$ -axis. This implies piezoelectricity: when the crystal is deformed in any direction not parallel to the z -axis, the atomic nuclei and their surrounding electrons are marginally separated, so that the crystal polarizes, resulting in opposite charges on some of its faces. Conversely, when it is subjected to an electric field, it deforms.

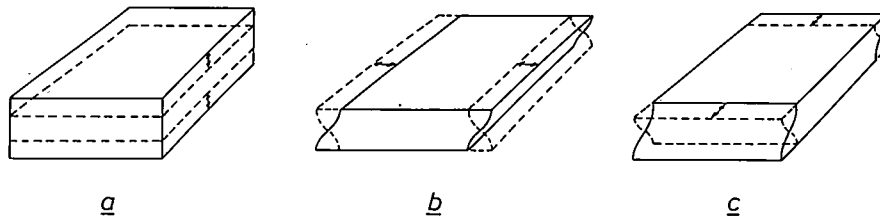


Fig. 3. The three predominant modes of vibration in a thin plate of quartz for a simple case. The vibrations are acoustic standing waves across the thickness of the plate; the displacements in mode a are longitudinal, those in modes b and c are shear waves. For each mode the fundamental resonance vibration is shown, corresponding to one half-wavelength across the thickness of the plate. In a doubly rotated cut the vibrations are not purely shear waves, e.g. in the c mode of a TTC cut the particles move at an angle of 12° to the plane of the plate.

When the plate is large and thin, there are three predominant modes of mechanical vibration, corresponding to standing waves across the thickness of the plate. They are labelled a, b and c, in order of decreasing wave velocity. In a simple cut, e.g. the AT cut, the displacements in mode a are longitudinal (i.e. perpendicular to the plate), those in modes b and c are shear displacements (see fig. 3). In the TTC cut, this is only approximately the case; the displacements in the b and c modes are about 10% longitudinal.

The 'antiresonant' frequencies f_m of mode m ($= a, b$ or c) are given by

$$f_m = Mv_m/4h, \quad (1)$$

where $2h$ is the thickness of the plate, M an odd integer (1, 3, 5, ...) and v_m the acoustic wave velocity of mode m :

$$v_m = (c_m/\rho)^{1/2}, \quad (2)$$

where ρ is the density of the crystal and c_m the elastic

stiffness for the strains occurring in mode m . The quantities c_m and v_m for each mode depend on the cut. Table I gives values of v_a , v_b and v_c for the AT cut and the TTC cut. Eq. (1) can be derived from the condition that the mechanical stress must be zero at the faces, and results from the conclusion that the thickness of the plate must contain an odd number of half-wavelengths. Oscillations corresponding to even values of M are of no interest; they correspond to equal charges on the faces instead of opposite charges, and cannot be excited electrically. In an AT-cut crystal, the c mode is the only one that can be excited electrically, since the 'piezoelectric coupling factors' k_a and k_b are zero for AT; in a TTC crystal all three modes can be excited. The term 'antiresonant' refers to the electrical behaviour.

A quartz resonator is a plate of quartz with electrodes on its faces. A mechanical resonance has a profound effect on the electrical impedance of the resonator. For a given alternating current in the resonator, the voltage across it is the sum of the usual

Table I. Wave velocities for the modes a, b and c in TTC- and AT-cut quartz plates.

	TTC	AT
v_a	6700 m/s	7000 m/s
v_b	4000	3800
v_c	3600	3200

[1] A complete review of the new cuts is given by A. Ballatto in: W. P. Mason and R. N. Thurston (eds), Physical acoustics XIII, p. 115; Academic Press, New York 1977.

[2] R. B. C. Maddox (PRL, Redhill) designed the jigs and supervised their construction. J. Dowsett, R. Butters and P. Morley (Cathodeon, Linton) were responsible for the manufacture of the devices.

[3] An account of early work in the field is given in R. A. Heising, Quartz crystals for electrical circuits, Van Nostrand, New York 1946. This excellent book was reprinted by Electronic Industries Association, Washington 1978.

The subject has also been reviewed by J. C. B. Missel and by W. Parrish, Philips tech. Rev. 11, 145, 323 and 351, 1949/50, and 12, 166, 1950/51.

[4] A review of hydrothermal growth is given in J. C. Brice, Repts Prog. Phys. 40, 567, 1977.

'dielectric' voltage V_d and a 'piezoelectric' voltage V_p , due to the deformations of the crystal. If the frequency is well away from any resonance, V_p can be neglected. When a point of resonance is approached, however, mechanical vibrations induce appreciable values of V_p . Away from resonance, V_p and V_d are not related in phase. At some frequency just below mechanical resonance, they compensate one another, resulting in *zero impedance*. At the point of mechanical vibrations itself, however, mechanical vibrations, associated with a voltage V_p , can occur without any current (if losses are neglected, that is), and so *the impedance is infinite*.

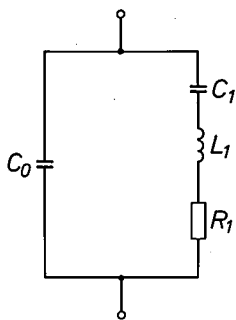


Fig. 4. Equivalent circuit of a quartz plate resonator near one of its resonant frequencies. C_0 is the static capacity between the electrodes on the faces. R_1 represents the losses. C_1 and L_1 are given by the relations:

$$C_1 = 8 C_0 k_m^2 / \pi^2 M^2,$$

$$L_1 = h^2 / 2 C_0 k_m^2 v_m^2.$$

$$(L_1 C_1 = 1 / \omega_a^2, \omega_a = 2\pi f_a = \pi M v_m / 2h.),$$

k_m is the piezoelectric coupling factor, giving the degree of coupling between the mechanical quantities (stress and strain) and the electrical quantities (electric field and dielectric displacement) that describe the quartz crystal when it is oscillating in mode m ; M , h and v_m as in eq. (1).

This electrical behaviour — corrected for losses — is simulated in detail by the equivalent circuit of *fig. 4*. C_0 is the static capacity, C_1 and L_1 relate to the mechanical properties (c_m, ρ) and the piezoelectric coupling factor k_m for the mode used; R_1 represents the losses. *Fig. 5* gives the resistive and reactive components R and X of the impedance Z as a function of frequency f near resonance. The impedance has a minimum ($X = 0$) at 'resonance' ($f = f_r$), a maximum at 'antiresonance' (mechanical resonance, $f = f_a$). The curves and the equivalent circuit only apply of course in the neighbourhood of a single resonant frequency, well away from other resonances; at every other resonance the behaviour is qualitatively repeated.

The crystal can be used in a series circuit or in a parallel circuit (*fig. 6*). The load capacitor C_L may be an integral part of the source or it may be used for

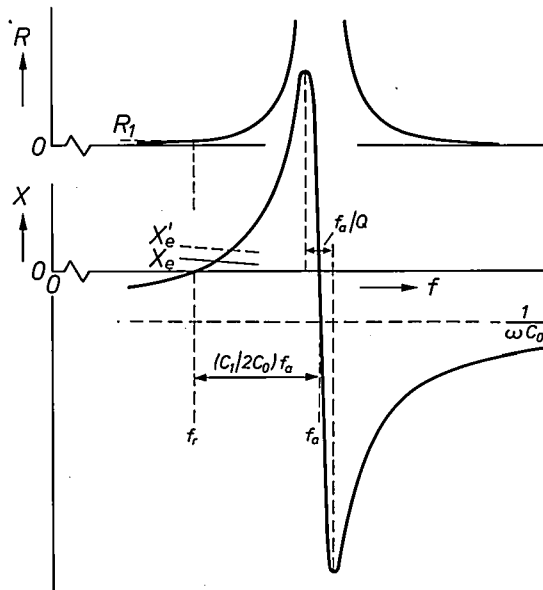


Fig. 5. Resistive and reactive components R and X of the impedance Z of the equivalent circuit (*fig. 4*) as a function of frequency near a single mechanical resonance. There are two points of resonance: at the 'resonant frequency' f_r ($X = 0$, $|Z|$ is at a minimum) and at the 'antiresonant frequency' f_a (the 'switch-over point' of X ; $|Z|$ and R have a maximum). The distance between the two resonance points is $(C_1/2C_0)f_a$, and is very small with respect to the frequency if C_0 is much larger than C_1 .

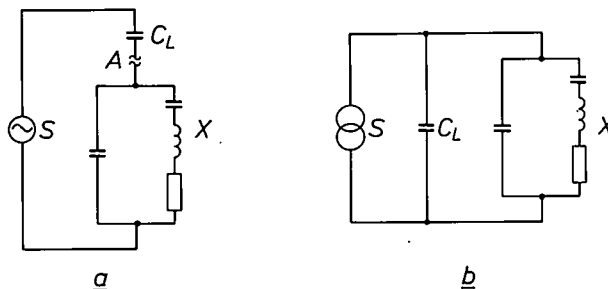


Fig. 6. Oscillator circuits. *a*) Series circuit. *b*) Parallel circuit. S source (often some kind of low- Q oscillator itself). X crystal. C_L load capacitor. In a series circuit S has to be of low internal impedance, in a parallel circuit it must be of low admittance.

final trimming of the frequency or for adjusting the degree of coupling. If we consider a series circuit, the crystal is used with a low-impedance source, and operated near f_r . The circuit can be considered as a 'feedback loop'. The condition for oscillation is that, if the loop is broken, say at A in *fig. 6a*, the input and output signals at the break are equal in phase and amplitude. If the Q is high, a slight frequency change strongly affects the phase. The frequency therefore adjusts itself to a value where X in *fig. 5* is equal and opposite to some reactance X_e determined by the 'external' circuit.

There are now two ‘measures of quality’ for the resonator. In the first place we have the usual quality factor Q :

$$Q = \omega_r L_1 / R_1 = 1 / \omega_r R_1 C_1 \quad (\omega_r = 2\pi f_r = 1 / \sqrt{L_1 C_1}),$$

giving the accuracy with which the resonator locks the oscillating frequency to its frequency of resonance; Q is inversely proportional to the losses in the system (Q is the ratio of the energy stored in the mechanical vibration to the energy lost per cycle). If the internal friction accompanying the vibrations were the only cause of losses, Q could be of the order of 10^7 . In actual devices energy may also be lost to the gas in the encapsulation, to the mounting system and to other modes of vibration. Thus, for the best devices, the encapsulation should be evacuated, the crystal should be ‘decoupled’ from the mounting, and the mode used should not couple to other modes of vibration.

Secondly, the steeper the $X_r f$ -curve (fig. 5) near f_r , or the closer f_r to f_a , the greater the immunity of the oscillation frequency to changes in the external circuit (i.e. in X_e), e.g. due to temperature changes in external circuit elements. From the equivalent circuit we can show that

$$(f_a - f_r) / f_r = C_1 / 2C_0, \quad (dX/d\omega)_{\omega=\omega_r} = 2L_1.$$

These quantities are of course directly related ($L_1 C_1 = 1/\omega_r^2$). Thus C_0/C_1 (or $C_0 L_1 \omega_r^2$) is a good

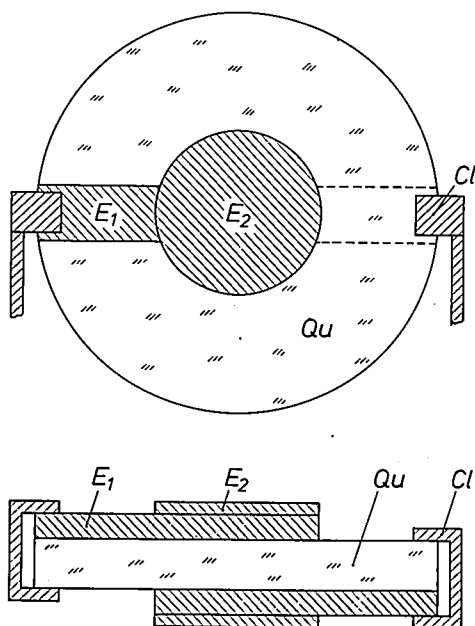


Fig. 7. A quartz plate with electrodes, shown in plan and elevation; Qu quartz plate, Cl mounting clip or clamp; E_1 is the electrode proper, E_2 is extra electrode material added to tune the device to the desired frequency. Not to scale; compare with fig. 11.

measure of the immunity to changes in the external circuit. When this quantity is translated back into physical quantities, we find (see the caption to fig. 4):

$$C_0/C_1 = \pi^2 M^2 / 8k_m^2.$$

Thus, cuts with small piezoelectric coupling factors are desirable but very small values of k_m produce a device whose circuit frequency cannot be adjusted.

A quartz plate with electrodes is shown schematically in fig. 7. Its frequency is mainly determined by the thickness of the plate (eq. (1)), but slightly decreased by the mass of the electrodes. The portion of the plate under the electrodes thus has a different resonant frequency from the non-electroded part. This helps to decouple the central, oscillating part of the plate from the mounting system as the energy becomes trapped under the electroded (low-frequency) area, and this leads to an improved quality factor. This effect can be enhanced by making the outer portions of the plate thinner.

The AT cut and the TTC cut

Special cuts have often been selected on the basis of the flatness of their frequency-temperature curve at a particular temperature. Thus ‘BT cuts’ (operating in the b mode) have been widely used; their f, T -curves have a maximum — i.e. the first temperature coefficient $T_f^{(1)}$ is zero — at a convenient temperature. The AT cut (operating in the c mode) has been so successful because the first two temperature coefficients, $T_f^{(1)}$ and $T_f^{(2)}$, are zero at room temperature. The f, T -curve is therefore extremely flat, since it has also a point of inflection at room temperature.

At PRL computer programs have been written that calculate ‘third-order’ f, T -curves about a selected temperature T_0 as defined by the first three terms of the power series:

$$(f - f_0) / f_0 = T_f^{(1)}(T - T_0) + T_f^{(2)}(T - T_0)^2 + T_f^{(3)}(T - T_0)^3 + \dots,$$

where $f_0 = f(T_0)$, and $T_f^{(n)} = \frac{1}{n!} \frac{\partial^n f}{\partial T^n}$. The

temperature coefficients for f are derived from those for c_m , q and h by differentiating eq. (1) with respect to T and using eq. (2), remembering that c_m , q and h are all functions of temperature. For $T_f^{(1)}$ we find:

$$T_f^{(1)} = \frac{1}{2} T_c^{(1)} - \frac{1}{2} T_q^{(1)} - T_h^{(1)},$$

where $T_c^{(1)}$ is the first temperature coefficient of c , etc. The expressions for the second and third coefficient become very tedious, but can easily be handled by the computer program.

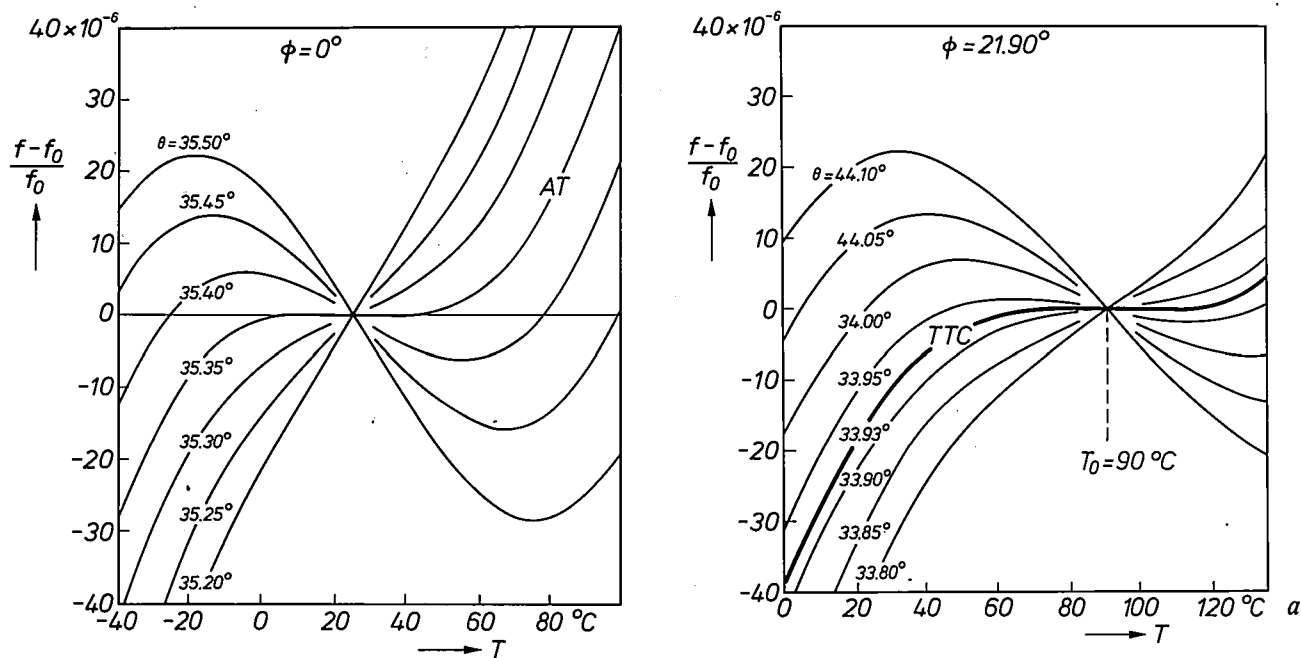


Fig. 8. Calculated frequency-temperature curves for cuts close to the AT cut proper. The relative frequency deviation, $(f - f_0)/f_0$, from the resonant frequency f_0 at the temperature $T_0 = 27^\circ\text{C}$ is plotted against temperature for cuts with $\phi = 0$ and various values of θ . All curves have an inflection point at 27°C . The AT curve ($\theta = 35.35^\circ$) is also flat at this temperature.

The AT cut is a singly rotated cut: $\phi = 0$ (see fig. 1). Fig. 8 gives a set of calculated third-order curves for $T_0 = 27^\circ\text{C}$, $\phi = 0$, with various values of θ . All of these curves have an inflection point at 27°C ; the AT cut ($\theta = 35.35^\circ$) is the one that is also flat at this temperature.

When ϕ is increased, the common temperature of inflection increases. Fig. 9a gives a set of curves for $\phi = 21.90^\circ$; the common inflection point occurs at $T = 90^\circ\text{C}$, a convenient oven temperature. The curve that is also flat at this temperature is our TTC cut ($\theta = 33.93^\circ$), which thus has the same advantage at 90°C as AT has at room temperature. Fig. 9b shows the sensitivity of the curve to variations in ϕ .

This cut is also *strain-compensated*. In the ϕ, θ -diagram of fig. 10a, the line A represents the values of ϕ and θ where $T_f^{(1)}$ and $T_f^{(2)}$ are equal to zero at the same temperature; this temperature, as a function of ϕ , is given in fig. 10b. The line B represents the values where the frequency is expected to be independent of an in-plane strain. This line is somewhat less certain than A, since the relevant material parameters are not so well known, but it is a good representation of what is known from experiments^[1] and calculations. The TTC cut is at the intersection of A and B. Figs 8 and 9 are for 'ideal' devices with very thin electrodes, whereas fig. 10 is for practical devices with electrodes about $0.5\ \mu\text{m}$ thick.

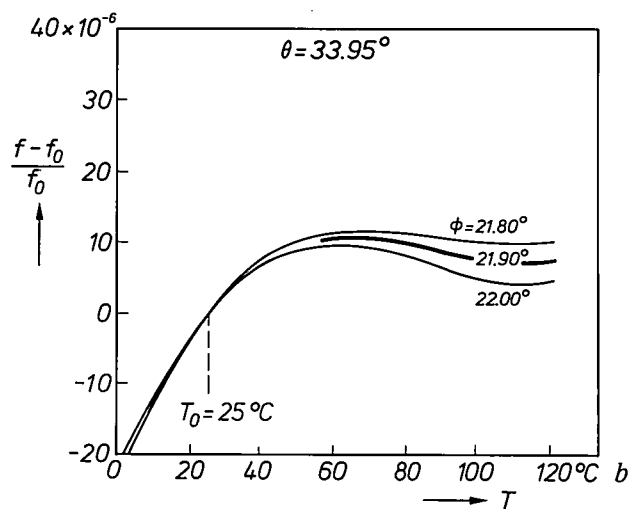


Fig. 9. a) Calculated frequency-temperature curves for $\phi = 21.90^\circ$ and various values of θ . The common inflection point is at $T_0 = 90^\circ\text{C}$. The curve for $\theta = 33.93^\circ$, the TTC cut, has a very wide temperature range of constant frequency about 90°C . b) Calculated curves for $\theta = 33.95^\circ$ and three values of ϕ . The curves in (b) have been centred on 25°C to separate them in the area of interest. The thick curves have been checked experimentally.

The advantages of the TTC cut

These features of the TTC cut are very valuable for oven-controlled resonators. If it is decided to use oven control, say at 90°C , it is of course possible to use an AT crystal with a slightly different value of θ , so that a turning point occurs at 90°C ($\theta \approx 35.60^\circ$, see fig. 8). The f, T -curve is then flat at the temperature of operation. The special feature of AT (inflection point at 27°C) is then of no advantage, but this seems to be the best approach if only singly rotated cuts are permissible, and it has long been the conventional ap-

proach. If double rotation is allowed, however, the TTC cut offers a much flatter curve at 90 °C.

This flatness may be exploited to save oven power — this might be important for portable communication equipment powered by batteries. If an oven is to consume little power, it must have good thermal insulation, but this implies slow cycling with large temperature variations; thus, if large variations are allowed, power may be saved. By increasing the permissible temperature changes from 4 °C (AT) to 20 °C (TTC), the power required for a simple oven design decreases from 5 W to 1 W.

For oven-controlled oscillators, strain compensation is also extremely important, as the oven temperature tends to cycle rapidly over a small range (e.g. 0.2 °C every 30 seconds). This implies that heat flows in and out of the quartz plate, giving temperature gradients, which are always accompanied by strains. Thus, if the frequency is *not* independent of strain, as in the AT case, it fluctuates much more strongly than

would be expected from its 'static' f, T -curve. This effect is avoided by strain-compensation.

Immunity to strain has other advantages. If the frequency *does* depend on strain, it will not be perfectly constant if the resonator is subject to vibrations, as is often the case in practice. The frequency may also change in the long run ('ageing') because of the relaxation of strains introduced into the mounting when the resonator was assembled. Finally, we should mention the 'nonlinear effects' introduced by strain-dependence: at high powers, the strains occurring during the oscillation itself will affect the elastic 'constants', which are then no longer constant during the oscillation; this may result in frequency changes and distortion of the sine wave, giving noise. All these effects should be eliminated or reduced by strain compensation.

In the following we shall show that these expectations are borne out in practice.

The experimental investigation

We have made TTC-cut resonators that operate at frequencies between 5 and 16 MHz in the fundamental mode ($M = 1$ in eq. (1)), and resonators operating at frequencies up to 42 MHz in the third overtone ($M = 3$). In all cases the c mode of vibration is used. In these resonators, unlike AT-cut ones, the unwanted b mode *can* be excited electrically, and is separated from the c mode by only 10% in frequency. A simple filter is usually added to prevent the circuit from oscillating in the b mode. The b mode has a very high temperature coefficient.

In this article we shall concentrate on the characteristics of 10 MHz resonators. We shall compare them with AT devices in terms of the 'quality criteria' Q and C_0/C_1 , the insensitivity of the frequency to strain, thermal shock and ageing, and to direct voltages across the electrodes. Except for the orientation of the cut, both sets of devices were made by identical and fairly simple techniques. Where the comparison is at 85 °C, adapted AT crystals with a turnover point at 85 °C were used. The results on the 10 MHz resonators are also representative for the other resonators.

Manufacture

The crystal from which the plates are to be cut is mounted on a jig with two degrees of freedom in orientation with respect to the saw blade of the cutting machine. The proper orientation is verified by X-ray diffraction to an accuracy of about 0.01° in both θ and ϕ . The electrodes (fig. 7) affect the temperature coefficient slightly, so that, for zero temperature coefficient, the values of θ and ϕ have to be

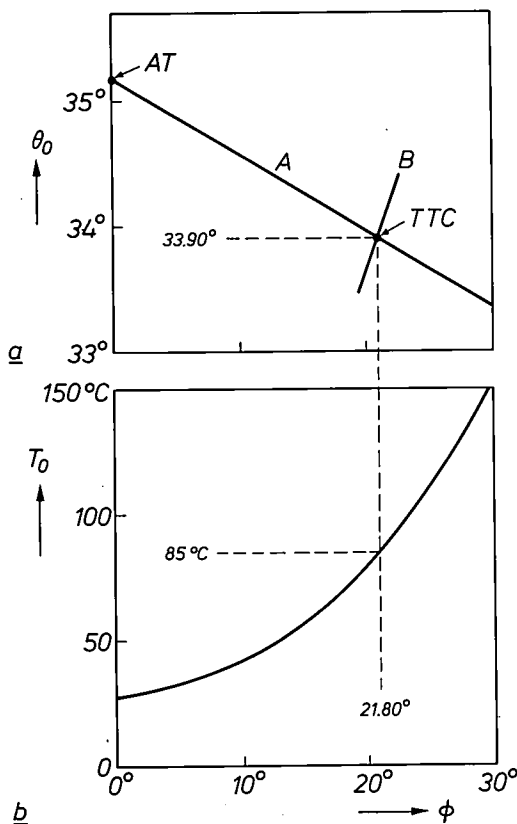


Fig. 10. The orientation angle θ_0 (above, line A) and the temperature T_0 (below) for which both $T_f^{(1)}$ and $T_f^{(2)}$ are equal to zero, as functions of the orientation angle ϕ . Line B in the upper diagram gives the θ -values for which the resonant frequency would be expected to be strain independent [1]. The cut at the intersection of A and B (TTC) is thus expected to have very good temperature compensation as well as strain compensation. Note that θ and ϕ are on different scales; the variations of θ along the line A are less than 2°. It should be noted that figs 8 and 9 apply to 'ideal' devices with very thin electrodes, but this figure applies to practical devices with electrodes about 0.5 μm thick.

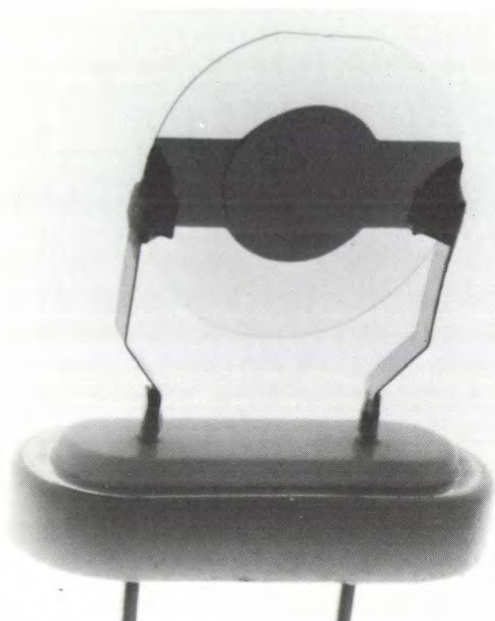


Fig. 11. A quartz plate mounted on its holder. The diameter of the plate is 12 mm.

slightly different from those for a blank as derived from fig. 9.

The cutting is a rather coarse procedure and the plates are lapped with abrasive powders between cast-iron laps to make a coarse adjustment of the thickness. The plates are contained in perforations in a plastic sheet. A radio receiver is connected to the upper and lower laps and when the correct frequency is reached, a signal is heard and the lapping is stopped. The slices are then stacked and ground circular, and separated for etching. The etching removes the damage done by the lapping process and prepares the surface for application of the electrodes by evaporation. Etching is continued until the resonant frequency is slightly higher than the desired value. After the deposition of the electrodes, which are of silver or gold, the discs are mounted on their holders (fig. 11) and a further small amount of electrode material is evaporated on to the central area of the electrodes until the resonator is tuned to the desired frequency. The frequency is monitored continuously during this process. Finally the devices are sealed either in evacuated glass envelopes or in gas-filled metal cases.

The main problems in this process are achieving and maintaining the desired angles and achieving the desired thickness after etching. Only a small amount of electrode material may be used for the final tuning, since it affects the temperature coefficient slightly.

The equivalent-circuit parameters

Table II presents the mean equivalent-circuit parameters obtained for a large number of 10 MHz resonators. The table shows that the TTC resonators

compare well with the AT resonators and have higher values of Q and C_0/C_1 . This would be expected for C_0/C_1 , since the piezoelectric coupling constant k_c in mode c is smaller for a TTC cut than it is for an AT cut.

Although the Q -values are not nearly as high as can be obtained, we think that the difference of a factor of 2 is significant, since both sets of resonators were produced in identical ways. The energy trapping is probably better for the mode of vibration used in the TTC cut than for the one in the AT cut.

The Q -values can be improved by making the outer part of the plates thinner (see p. 5) and by polishing the blanks. Q -values limited by the losses in the quartz can then be attained easily, but for most purposes a lower Q is satisfactory and the extra costs cannot be justified.

The f, T -curves; the first temperature coefficient

The temperature behaviour of the plates is as expected from the calculations; the thick curves in fig. 9 were experimentally obtained with plates with negligible electrode weights. Figs 9a and b show that $T_f^{(1)}$ is much more sensitive to θ than it is to ϕ ; TTC is also less sensitive to θ -variations than AT. From the estimated accuracy of $\pm 0.01^\circ$ in θ and from the curves, we would expect variations in $T_f^{(1)}$ of $\pm 3 \times 10^{-8}/^\circ\text{C}$ for the TTC plates, and $\pm 5 \times 10^{-8}/^\circ\text{C}$ for the AT plates. The variations that we actually found in batches of crystals are about $\pm 1 \times 10^{-7}/^\circ\text{C}$ for TTC crystals and $\pm 2 \times 10^{-7}/^\circ\text{C}$ for AT crystals.

Thermal shock

When a resonator is subjected to a sudden change of temperature a transient change of frequency, due to transient strains, is usually observed. When an AT crystal is temperature cycled, this effect upsets the frequency stability, even for cycles around room temperature. In our resonators, the transients take about 30 seconds to build up. This is roughly the time for a

Table II. The average equivalent-circuit parameters for a large number of TTC and AT 10 MHz resonators.

	TTC	AT
C_0	5.0 ± 0.3 pF	5.7 ± 0.02 pF
C_1	7.3 ± 0.4 fF	30.1 ± 1.4 fF
L_1	34.7 ± 0.3 mH	8.4 ± 0.2 mH
R_1	5.3 ± 0.5 Ω	2.6 ± 0.6 Ω
$Q = \omega L_1/R_1$	$(4.1 \pm 0.3) \times 10^5$	$(2.03 \pm 0.5) \times 10^5$
C_0/C_1	680 ± 60	190 ± 10
$k_c^{[5]}$	4.99%	8.80%
$\pi^2/8k_c^2$	495	159

change of case temperature to affect the quartz plate. The decay time is about 90 seconds.

The results of our thermal-shock experiments are summarized in *Table III*. It can be seen that the effect on the frequency is 10 to 20 times smaller for TTC crystals than for AT crystals. For the applications we have in mind, oven-controlled resonators, this is the most important advantage of the combination of strain compensation and temperature compensation that TTC cuts offer. *Fig. 12* shows a direct demonstration of this advantage.

Table III. The effect of thermal shock on the frequency. $(\Delta f/f)_m$ is the maximum relative frequency change occurring after a sudden change of temperature ΔT . At the test temperature (about 90 °C) $T_T^{(1)}$ for all resonators was zero to within $\pm 10 \times 10^{-9}/^\circ\text{C}$. The results are almost identical for AT devices if they are tested at 27 °C.

ΔT	$(\Delta f/f)_m$	
	TTC	AT
0.2 °C	10×10^{-9}	200×10^{-9}
2 °C	40×10^{-9}	500×10^{-9}

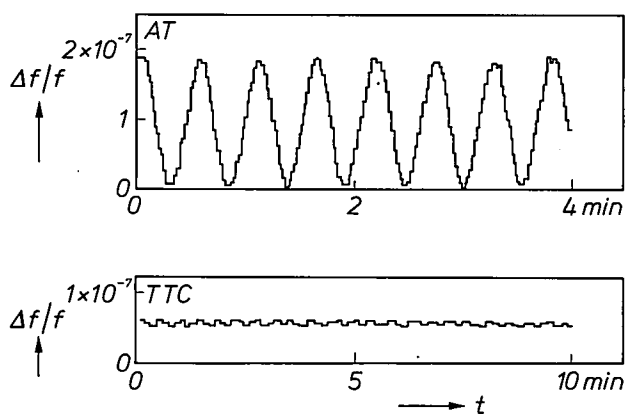


Fig. 12. Frequency stability of an AT crystal and a TTC crystal, temperature-cycled in a simple oven. Average oven temperature 71.5 °C. Cycling range 0.2 °C. Cycling period 35 s.

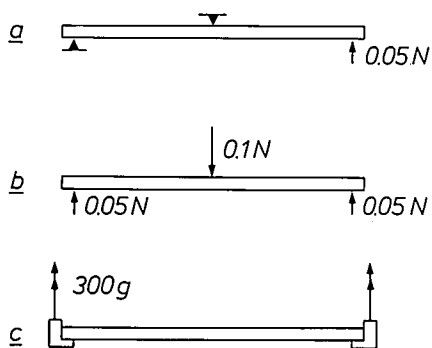


Fig. 13. a) Arrangement of symmetrical bending test, schematic. The upper, centre knife is in two sections, leaving room for the electrodes. The frequency change is measured for different orientations of the plate in the jig. b) The forces occurring in the test. c) If the mass of the plate is 30 mg, the forces in *b* are equivalent to an acceleration of 300g of the plate in its clamps.

Strain

Strain may occur not only because of temperature changes, but also more directly. For TTC crystals the frequency would be expected to be immune to strains in the plane of the plate. For a plate held at its edges (fig. 7), however, the most likely deformation to occur in practice is symmetrical bending, either by acceleration (vibration) of the resonator or by forces imposed directly due to the clamps. Special jigs were therefore devised in which plates were bent symmetrically as indicated in *fig. 13a*, for several orientations of the plate in the jig, and the change in resonant frequency was measured [6]. Typical results are shown in *fig. 14*. The maximum frequency change is at least five times smaller for TTC than for AT. In the AT crystal there is no orientation at which the effect is zero, whereas for TTC the effect fluctuates around zero. The AT crystals all behave rather similarly, with a period of about 90° in orientation. The TTC crystals do not produce a consistent pattern, however; the effects are

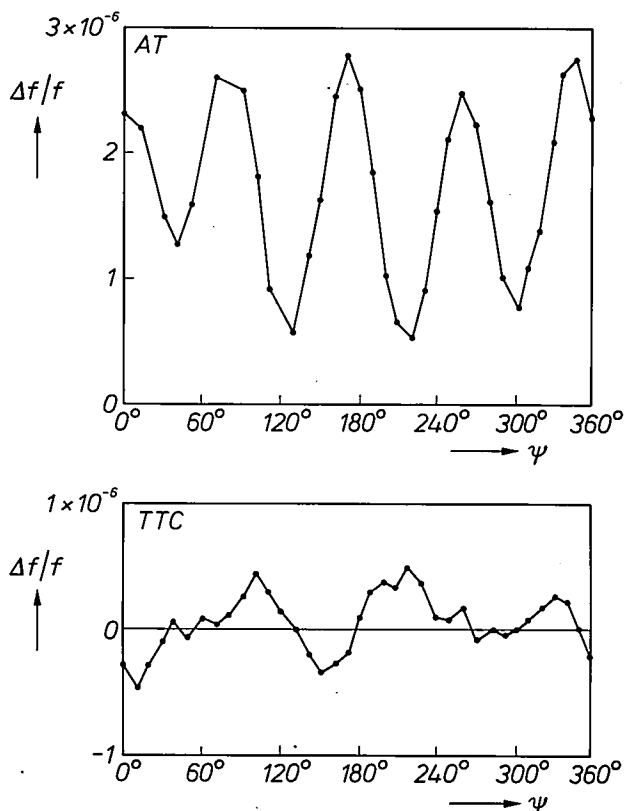


Fig. 14. Frequency changes with symmetrical bending for an AT-cut resonator and a TTC-cut 10 MHz resonator. The relative frequency change $\Delta f/f$ is plotted against the angle ψ between the line connecting the mounts or the points of application of the forces (fig. 13) and the projection of the *x*-axis on the plate (see fig. 1). Several AT crystals give much the same pattern. The TTC crystals do not give a constant pattern.

[5] From Table IV of Ballato's article [1].

[6] E. D. Fletcher and A. J. Douglas, Proc. 33rd Ann. Frequency Control Symp. 1979, p. 346.

believed to be secondary, associated with the exact positioning of the electrodes, etc.

The bending test of fig. 13*a* is equivalent to an acceleration of the plate in its mounts of about 300*g* where *g* is the gravitational acceleration (see fig. 13*b, c*; the mass of the plate is about 30 mg). An acceleration of 1*g* would give a relative frequency change of at most 1×10^{-9} . Since the acceleration does not exceed a few *g* in most vibrations, the frequency change will usually be negligible even for a desired accuracy of one part in 10^8 .

For asymmetrical bending more systematic changes are found; these are much the same for TTC and AT cuts (fig. 15). They are therefore probably associated with more complicated strains than simple in-plane strains. As the TTC curve of fig. 15, unlike the one in fig. 14, is reasonably reproducible, it can be used to select the optimum position for the mounting clips to make the device insensitive to these more complicated strains.

Ageing

It has been suggested (p. 7) that strain compensation might also reduce ageing effects for the frequency. For our 10 MHz resonators the TTC cuts are indeed consistently better than the AT cuts in this respect. We have applied a standard ageing test to a large number of resonators: they are kept at a temperature of 85 °C for three months and the frequency is monitored during that period. The frequencies of the TTC devices rose slowly during the entire period, but the rate of increase slowly decreased. The frequencies of the AT devices, on the other hand, initially fell rapidly, then more slowly until they also started to rise after six weeks. At the end of the three months the average ageing rates were $(1 \pm 0.3) \times 10^{-9}$ and $(6 \pm 4) \times 10^{-9}$ per day — equivalent to annual rates of about 0.4×10^{-6} and 2×10^{-6} — for the TTC and AT batches respectively.

As suggested earlier, the greater effects for the AT resonators can be explained by a greater sensitivity to the relaxation of strains in the mounting. There are of course other causes of ageing, and these should be much the same for both TTC and AT cuts, e.g. condensation of material on the electrodes or evaporation from them. As the plates are only about 10^6 atomic layers thick, a significant fraction of one layer of atoms would cause a frequency change of the order of 1 part in 10^6 .

The effect of a direct voltage

When a direct voltage is applied to a quartz resonator, the frequency changes abruptly because the d.c. field directly affects the elastic constants. It then slowly continues to drift, on account of the migration of small ions through the lattice (solid-state electrolysis); these ions create space charges that again modify the frequency. In our samples lithium appears

to be the migrating ion [7]. Lithium is a common impurity in quartz at concentrations of a few parts per million.

This is the one effect where TTC is at a disadvantage. In our AT devices, the frequency first drops 5 to

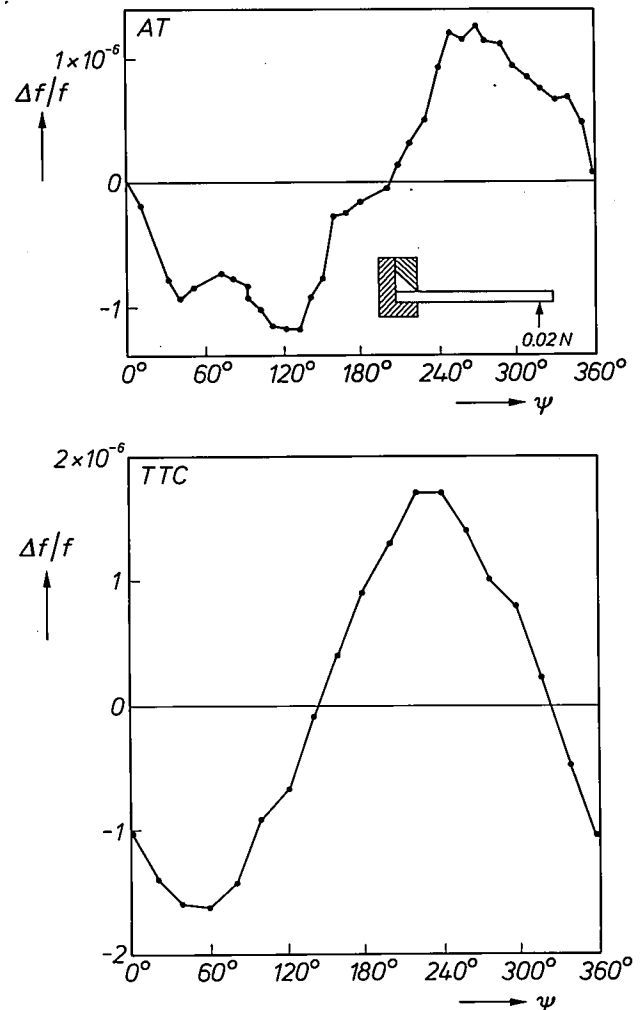


Fig. 15. Relative frequency changes for asymmetrical (cantilever) bending. Inset: experimental arrangement.

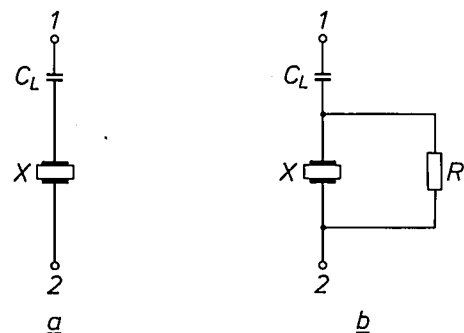


Fig. 16. Resonant part of an oscillator circuit. *X* quartz-crystal resonator, C_L series load capacitor. *a*) The conventional circuit. A direct voltage between 1 and 2 will appear across *X*, as the d.c. resistance of C_L is about $10^9 \Omega$, and that of *X* is about $10^{12} \Omega$. Such voltages may affect the frequency, especially in the TTC case. *b*) A resistance *R* of 1 M Ω will eliminate the direct voltage across *X*, without otherwise affecting the circuit.

[7] J. P. Stagg (PRL, Redhill) identified lithium as the mobile species.

8 parts in 10^9 per volt and then drifts down a further 1 part in 10^9 per volt in a few hours, independently of the polarity of the voltage. In the TTC devices the frequency changes are larger, persist much longer and are more complicated. The sign of the initial change depends on the polarity of the voltage, and the later changes do not have the nature of a persistent drift but are fluctuations. The initial change is 10 to 15 parts in 10^9 per volt, the subsequent fluctuations are 2 to 10 parts in 10^9 per volt, and changes can still be measured after 16 hours. The complications can be explained partly by the fact that the plane of the plate no longer contains a twofold axis. In a singly rotated cut the twofold x -axis in the plane of the plate ensures that any effect must be independent of the polarity of the voltage.

The effect is best circumvented by eliminating any significant direct voltages. Fig. 16a shows a typical conventional circuit. The d.c. resistance of the load capacitor C_L is 10^9 to $10^{10} \Omega$, and that of the crystal is 10^{12} to $10^{13} \Omega$. A direct voltage between points 1 and 2 therefore appears across the crystal. It can be easily eliminated, however, by adding a shunt resistance much smaller than the resistance of C_L but large enough not to affect the circuit in any other way (fig. 16b).

Conclusions

Our investigations show that the expected advantages of the TTC cut over the conventional AT cut for oven-controlled resonators are attained in practice. The frequency is much less sensitive to temperature variations in the range of convenient oven temperatures and to the thermal strains that occur in temperature cycling. The frequency is also less sensitive to vibrations, to changes in the external circuit and to ageing; the Q (the electrical quality factor) is better. Other measurements, not discussed in this article, have shown that in all other respects the TTC devices are at least as good as the AT devices made in the same way. In particular, the noise produced at frequencies away from the carrier is usually less and certainly not more than in comparable AT devices. Making use of TTC crystals will not result in a significant increase in price,

Table IV. Improvements in frequency stability, power consumption or range of stability that can be made by substituting the AT resonator by a TTC resonator in a simple oven.

Average oven temperature 85 °C	AT		TTC	
Power, maximum	8 W	8 W	1 W	1.5 W
Frequency variation	1 in 10^7	1 in 10^8	1 in 10^7	1 in 10^7
Ambient temperature range	0 to 60 °C	0 to 60 °C	0 to 60 °C	-40 to +60 °C

as they are only marginally more difficult to manufacture and it is simple and inexpensive to modify the oscillator circuit so that the b mode is suppressed and direct voltages do not appear across the crystals.

Two simple experiments demonstrate the possible improvements. In the first experiment, we made similar large changes in the external circuits of an AT oscillator and a corresponding TTC oscillator: a series load capacitor was decreased from 60 pF to 5 pF. In the AT oscillator this resulted in a frequency change of 9 parts in 10^4 . As would be expected from the C_0/C_1 -ratios (Table II), the change was much less for the TTC oscillator: 3 parts in 10^4 .

In the second experiment, we used an oven-controlled oscillator with an AT crystal; this had a maximum power consumption of 8 W and a frequency stability of 1 part in 10^7 over an ambient temperature range of 0-60 °C (Table IV, first column). By substituting a TTC resonator and adapting the oven where necessary, we could either improve the stability by a factor of 10, or reduce the required oven power by a factor of 8, or reduce the power a little less and extend the temperature range to -40 °C (see Table IV).

These simple experiments neatly summarize the ways of using the new cut. It can be used to improve the frequency stability of a system, to make a system suitable for tougher physical conditions and to save oven power.

In practice TTC-cut crystals are now coming into use in the most demanding situations, e.g. in frequency counters, satellite navigation systems, television transmitters and in high-frequency single-sideband radio systems.

This work was carried out with the support of Procurement Executive, U.K. Ministry of Defence, sponsored by DCVD.

Summary. It was expected that the properties of the 'thermal-transient compensated' (TTC) cut of quartz would make it very suitable for use in oven-controlled oscillators. It is a 'doubly rotated' cut, in which the two angles (ϕ and θ) characterizing the orientation are different from zero, unlike the widely used 'singly rotated' AT cut ($\phi = 0$). Computer calculations predict that the frequency of the TTC cut is very insensitive to temperature changes around 90 °C (a convenient oven temperature) and to mechanical strain. The strain compensation eliminates the frequency transients due to thermal strains occurring with oven control, and also eliminates some of the ageing and nonlinear effects. An experimental comparison of AT and TTC resonators for 10 MHz, identical except for the angles of cutting, bears out the expectations. The frequency-temperature curve of the TTC crystals is very much flatter at 90 °C, the sensitivity to thermal shock and symmetrical bending is very much less and the ageing effects are smaller. Two disadvantages of the TTC cut, the ease of excitation of an unwanted mode of oscillation (b) near the mode of operation (c) and a sensitivity to direct voltages, are easily circumvented by simple provisions in the circuit. The new cut can be used to improve the frequency stability of a system, to make it suitable for tougher physical conditions and to save oven power.

Camera window for ultrasoft X-rays from celestial sources

J. A. M. Bleeker, W. H. Diemer, A. P. Huben and H. Huizenga

The Netherlands astronomical satellite ANS carried instruments for observations in the X-ray bands (wavelengths 0.2-6 nm) and in the ultraviolet (150-330 nm). Astronomers are keenly interested in extending the X-ray observations to wavelengths of 20 nm or more, and in obtaining a good survey of the sky at these wavelengths, which was not possible with the instruments on the ANS. The 'Cosmic Ray Working Group', based at the Huygens Laboratory of Leiden University, has built a camera for this purpose. The article below discusses the ultra-thin plastic window of this camera.

When the ANS satellite was built, Philips were closely concerned with its construction. In the present work the Philips share is confined to a contribution to the development of the camera window. The research on this project was commissioned by the Cosmic Ray Working Group. The editors are grateful to members of this Group for their cooperation in outlining the background to this research.

Introduction

For studying astronomical objects in the wavelength range of soft and very soft X-rays (1-25 nm) the 'Cosmic Ray Working Group' at Leiden has developed a focal-plane camera. The camera comprises a position-sensitive proportional counter for X-ray photons, on which an image of a section of about 2 by 2 degrees of the sky is produced by means of an X-ray optical system. It is intended to incorporate the instrument in a satellite and to use it for mapping the entire celestial sphere in the stated wavelength range with an angular resolution of about one minute of arc. The results of laboratory measurements with a prototype of the camera show that it fully comes up to expectations [1].

The use of proportional counters for X-ray observations in space has until now been restricted to wavelengths shorter than about 8 nm; this limit is determined by the thickness of the entrance window of the counter. To shift this 'cut-off' beyond 20 nm, an ultra-thin plastic window had to be developed that would meet very exacting requirements for effective area, uniformity, mechanical strength, gas-diffusion leakage and X-ray transmission. This window, developed by the working group in cooperation with the X-ray Tubes Laboratory of the Philips Scientific and Industrial Equipment Division in Eindhoven, is the subject of this article [2]. First of all we shall deal

briefly with the significance of this spectral region in astronomy and also with the camera itself. The X-ray waveband of interest will be referred to as the XUV range (see *fig. 1a*).

Astrophysical objectives

Astronomers have been very interested in the XUV range ever since it was discovered in the seventies that the regions of space closest to us contain much less interstellar matter than had originally been thought, and is thus more 'transparent' to XUV rays. The well-known radio measurements on the 21 cm hydrogen line indicate a mean neutral gas density of 1 atom per cm^3 in our galactic system. With a density of this order the 'visibility' for X-radiation at say 10 nm is only about 100 light years (see *fig. 1b*). Not many interesting X-ray objects would be expected in a sphere of this radius. The density quoted, however, is an average over galactic distances (1000 to 10 000 light years), and various recent observations (including XUV radiation) suggest that the local density is only 0.06 hydrogen atoms per cm^3 . Consequently the estimated local visibility has become more than ten times greater, and the probability of observing interesting objects has increased more than a thousand times.

Objects of particular interest here are the masses of gas ejected by stars in the final stages of their existence. Heavy stars then become supernovae, leaving a neutron star or perhaps a 'black hole' behind as a nucleus. The remnants of such supernova explosions — familiar examples are the Crab nebula in the constel-

Dr Ir J. A. M. Bleeker is the head of the Cosmic Ray Working Group, Huygens Laboratory, Leiden; Drs W. H. Diemer and A. P. Huben are with the X-ray Tubes Laboratory of the Philips Scientific and Industrial Equipment Division, Eindhoven; Dr H. Huizenga was formerly with the Cosmic Ray Working Group.

lation Taurus and the Veil nebula in the constellation Cygnus — can continue to exist for a very long time. The lighter stars evolve in a much more gradual process into a ‘White Dwarf’, possibly via the planetary-nebula stage. In gas masses of both kinds temperatures of between 10^5 - 10^7 K are expected, which means that they will radiate most strongly in the XUV region (fig. 1a), and will be completely ionized (i.e. they will form plasmas).

Finally, models have been developed in recent years that relate the high local visibility for XUV rays to

very old supernova remnants. In some of these models the interstellar space consists largely of hot plasma regions ($\approx 10^6$ K) in which cold gas clouds (10 to 100 K) circulate, surrounded by a ‘warm’ transitional layer ($\approx 10^4$ K). The presence and the properties of the plasma regions are again very suitable for study by means of XUV observations.

To summarize, it is assumed that the XUV maps of the celestial sphere will contain spots, streaks and larger areas that may provide a great deal of information about the final stages in the evolution of stars and the consequences for the interstellar medium.

The camera

The camera is designed for a ‘grazing-incidence telescope’. This makes use of the total internal reflection that occurs when X-radiation is incident on a smooth metal surface at an angle smaller than the critical grazing angle. The critical angle is proportional to the wavelength. In this way it is possible to make an optical system that will form images for very soft radiation (see fig. 2).

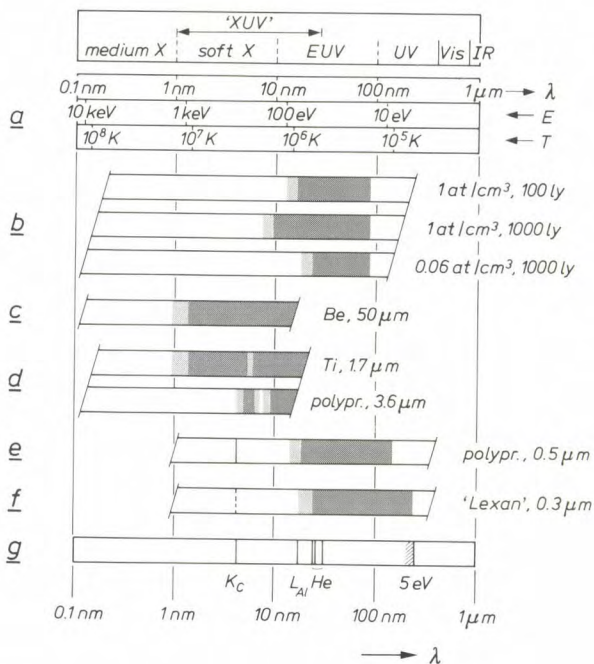


Fig. 1. a) Radiation from the X-ray range to the infrared range, characterized by the wavelength λ and the photon energy E ($E\lambda = hc = 1.2398 \text{ keV nm}$). XUV denotes the region from about 1 nm to about 30 nm. The temperature scale gives the temperature from the relation $E = kT$ ($1 \text{ keV} \hat{=} 1.16 \times 10^7 \text{ K}$). This relation gives the photon energy that is characteristic of the radiation emitted by a plasma of temperature T . b to f) Transmission of the interstellar medium and of various types of window. Dark grey: transmission less than 1%; light grey: between 1% and 10%; white: greater than 10%. b) Interstellar matter of different densities n_H (number of hydrogen atoms per cm^3) and column length (light years, ly). From 21 cm measurements an average n_H of 1 at/cm^3 is derived. At 10 to 20 nm a column of 1000 light years is then completely opaque (second strip), the ‘visibility’ is no more than about 100 light years (first strip). Recent observations, however, indicate that the local density is only about 0.06 at/cm^3 , which implies that at 10 to 20 nm a column of 100 light years is still transparent (third strip). c) Beryllium foil of 50 μm , widely used in rocket experiments for harder X-rays. d) The windows used in the ANS for the 3 nm and 5 nm ranges. e) Calculated transmission of 0.5 μm polypropylene. f) as (e), for 0.3 μm ‘Lexan’; this is the window discussed in the present article. In addition to a better XUV transmission it also has the desired UV absorption for photons of energy greater than 5 eV (the work function of some metals). g) Some special wavelengths and energies: from left to right: the K absorption edge of carbon (K_c), characteristic of most plastic windows (4.4 nm); the Al-L edge (L_{Al}), where the Al layer at the inside edge of the window causes a slight reduction in transmission (17 nm); some He lines at which transmission measurements have been performed (24.3 nm, 25.6 nm and 30.4 nm); and a typical value (5 eV) for the work function of some metals.

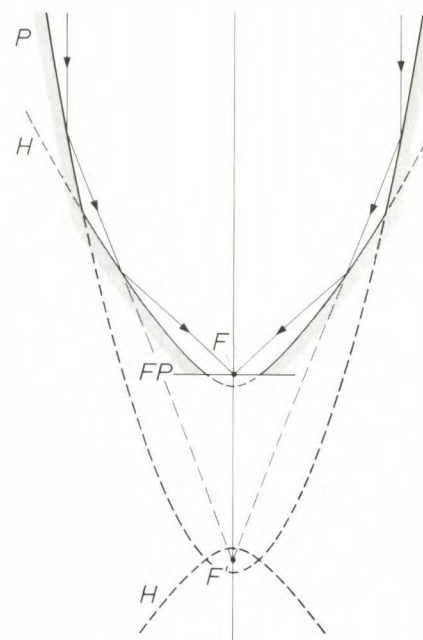


Fig. 2. Grazing-incidence telescope. The telescope reflector consists of a parabolic mirror and a hyperbolic mirror. The hyperbola corrects the image errors of the parabola. Rays parallel to the axis are reflected by the parabola P to the focal point F' of P . Since this is also the virtual focal point of the hyperbola H , the rays arrive at the real focal point F of H . An object at infinity produces an image in the focal plane FP , which is where the camera is located.

[1] A complete description of the camera is given in H. Huizenga, A focal plane camera for celestial XUV sources, Thesis, University of Amsterdam, 1980, and in J. A. M. Bleeker, H. Huizenga, A. J. F. den Boggende and A. C. Brinkman, IEEE Trans. NS-27, 176, 1980.
 [2] A more detailed description of the investigation is given in H. Huizenga, J. A. M. Bleeker, W. H. Diemer and A. P. Huben, Rev. sci. Instr. 52, 673, 1981.

The camera itself is shown schematically in *fig. 3*. An X-ray photon penetrating the entrance window *W* generates a primary electron cloud by photo-ionization and the subsequent collision ionization due to the ejected electrons. The electron cloud travels through a weak electric field in the drift space *D* towards a set of crossed grids (drift grid *G* and anode grid *A*). The strong electric field between the grids pulls the cloud through the drift grid, giving rise to an avalanche-like charge multiplication that generates 10^5 to 10^7 electrons and ions. A large number of 'charge sensors' *S* (charge-sensitive amplifiers) are connected to the grids at a uniform spacing. The charge distributes

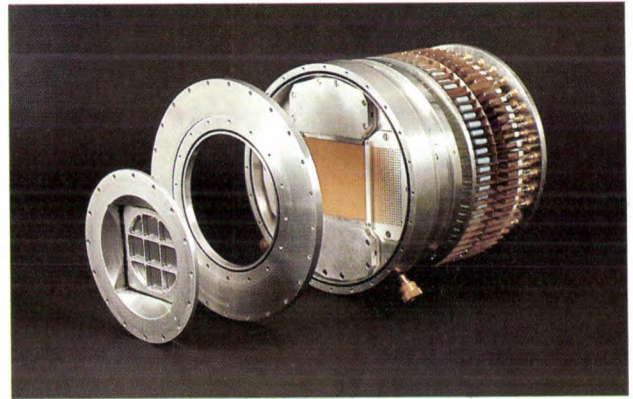


Fig. 4. The camera with top plate and window removed.

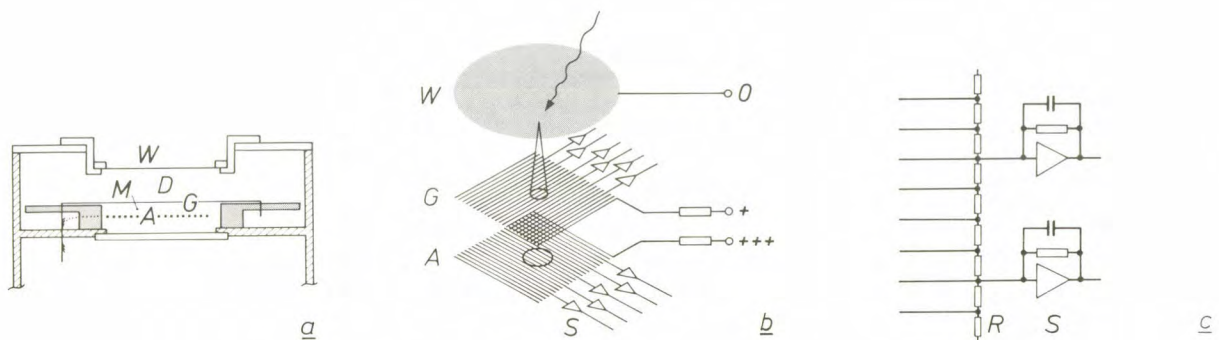


Fig. 3. The camera. *a*) Cross-section (schematic, not to scale). *b*) Principle of the spatial detection. *c*) Detail of a grid. *W* window. *D* drift space. *G* drift grid. *M* space where the charge multiplication takes place. *A* anode grid. *S* charge sensors (charge-sensitive amplifiers). *R* resistance circuits. The window aperture is 6.5 cm in diameter, the thickness of the drift space is 2 mm, the thickness of the avalanche space is 1 mm. There are 326 wires per grid, with a spacing of 0.2 mm.

itself over a large number of wires in each grid and is measured by two or three sensors. The output signals of the sensors enable the centre of the discharge to be determined to an accuracy better than the pitch of the wires. One grid thus determines the 'x-coordinate', the other the 'y-coordinate' of the photon. The camera not only determines the position of the photon, it also determines its energy; the energy is proportional to the sum of the three measured output signals. The sensors are connected to the transmitting equipment of the satellite, which transmits the magnitude of the energy, the x-coordinate, the y-coordinate and the time of reception of every single photon. Each magnitude is transmitted in the form of one or more words of 8 bits. The X-ray picture for any desired wavelength can be reconstructed on Earth from these data. The complete camera can be seen in *fig. 4*.

The proportional-counter technology has proved the most suitable for observing sources of some extent, such as planetary nebulae and supernova remnants^[3]; other technologies are available for point sources.

The window

The entrance window is a critical part of a proportional counter. A 50 μm beryllium foil — commercially available — is often used in counters for rocket-borne experiments in the medium X-ray range. This is completely opaque to XUV radiation, however (*fig. 1c*). In the ANS 3 nm radiation experiment the entrance window was a 1.7 μm titanium foil^[4], which has a 'dip' in the absorption spectrum at this wavelength (*fig. 1d*). Evaporation techniques have been used to produce very much thinner metallic foils, e.g. 100 nm, but these are not strong enough for windows.

For measuring soft X-rays the counters are now fitted with plastic windows. To remove the positive charge from the drift space the window has to be given a metallic coating on the inside, and this must be thin enough not to interfere with the XUV transmission. The application of this coating, however, introduces a fresh difficulty, since the work function of metals is in the region of 5 eV (about 250 nm) and the coating thus makes the counter sensitive to

UV photons of wavelength shorter than 250 nm. This means that the UV 'tail' in the spectrum of sources such as visible stars can give too large a background signal. The window must therefore absorb this short-wavelength UV radiation completely (see fig. 1g). Ultra-thin metallic foils are sometimes used as filters for this purpose. Finally, in satellite experiments, with a limited gas supply that must last for a long time, the diffusion leakage of plastic windows constitutes a serious problem. In rocket-borne experiments this problem is less serious.

Table I gives the materials that we have evaluated. Thin films of polypropylene are made by stretching thicker films in a special way. Films of this type are widely used in rocket experiments. They do not, however, meet our leak specification (see page 17) and in addition their UV transmission is too high for wavelengths longer than 160 nm (fig. 1e).

'Parylene N' is very suitable for thin films on a substrate, e.g. for insulation. These films are grown on the substrate directly by vacuum vapour-phase polymerization. To obtain a free film, it therefore has to be separated from the substrate. Separating the film without damaging it proved to be difficult with the commercially available film-on-substrate products. We also made ultra-thin (0.15-0.30 μm) uniform films of this material ourselves, using special intermediate layers^[5]; these films could be removed easily, but they also failed to meet our leak specification. In addition, the UV transmission of parylene N in the 140 to 180 nm band is too great.

Table I. The window materials we have investigated.

Polypropylene	$(\text{C}_3\text{H}_6)_n$
'Parylene N'	$(\text{C}_8\text{H}_8)_n$
'Formvar'	$(\text{C}_5\text{H}_8\text{O}_2)_n$
'Lexan'	$(\text{C}_{16}\text{H}_{14}\text{O}_3)_n$
'Bioden'	acetyl cellulose
Collodion	cellulose nitrate

Eventually we did manage to make good windows of 'Lexan' by 'casting' the films on a water surface. This method avoids the film-release problem. It produces ultra-thin films with a relatively low density of pinholes. Leakage through the pinholes is eliminated by stacking several of the films one above the other. The method was first used elsewhere for making 'Formvar' windows^[6] (see Table I). 'Lexan' was then added as a filter material; it has the desired absorption in the UV range (fig. 1f). Later films were made of 'Lexan' alone, and used as filters on polypropylene windows^[7].

We first used this method for making five-layer films, 200 nm thick, and of 14 mm diameter, of the

materials 'Formvar', 'Lexan', 'Bioden' and collodion (Table I). 'Lexan' was clearly the best, with low gas-diffusion leakage and the desired UV absorption, and so we concentrated on this material in our later experiments.

Film casting

Fig. 5 shows the equipment used for film casting. It consists of a glass tank containing deionized water, a carriage, which can be driven by a small vibration-free electric motor along a rail at the edge of the tank, and, attached to the carriage, a brass plate with a turned-up edge (the 'spreader'). The plate has been made hydrophobic by polishing. To start with, the spreader is filled with about 0.75 ml of a solution of 2% of 'Lexan' in dichloromethane, with 0.01% of dioctyl phthalate added as a softener. With the carriage at one end of the tank, the plate is lowered until its edge is just below the surface. The solution floats on the surface, and adheres to the spreader and to the front end of the tank. The carriage is set in motion, and the spreader is drawn along the tank at a velocity of 1 cm/s. A film is then produced on the water surface. The first part of the film is generally very inhomogeneous, while the last part contains many holes and strips, since dried pieces of material come away from the plate; the best parts of the film are in the middle.

The films produced at a velocity of 1 cm/s have a thickness of about 70 nm. This is about the optimum that can be obtained with stacked films: thinner films contain an increasingly greater number of holes; in thicker films the number of holes is not significantly reduced.

After the film has been formed, a carrier frame with an aperture of 65 mm diameter is lowered on to it. The frame is then dipped below the surface and withdrawn, with the film adhering to it. As soon as frame and film are dry they are pressed on to the next film that has been drawn, bowing the frame slightly to prevent air inclusions. Repeating this process four or five times yields films about 300 nm thick. Before each repetition a small piece is removed from the new film with a smaller frame and inspected under a microscope.

The result is highly sensitive to unwanted effects. The operating conditions for the process, including

- [5] A comparison with other technologies will be found in the publications of note [1].
 [6] See for example A. C. Brinkman, J. Heise and C. de Jager, Philips tech. Rev. 34, 43, 1974.
 [7] These experiments were carried out at Philips Research Laboratories, Eindhoven.
 [8] B. L. Henke, Adv. in X-ray Analysis 8, 269, 1965; F. Williamson and C. W. Maxson, Rev. sci. Instr. 46, 50, 1975.
 [9] J. A. M. Bleeker *et al.*, Astron. Astrophys. 69, 145, 1978.

the geometry of the water tank and spreader, the material and the treatment of the spreader, the composition and the quantity of the solution, had to be carefully optimized beforehand. The water and the solution must be very clean. The process has to be carried out in a dust-free environment and in the complete absence of vibration and turbulence: the air circulation in the dust-free bench is shut off during the casting process and, as we indicated earlier, vibration from the carriage motor was eliminated beforehand.

film to the mesh assembly must be soft to avoid damaging the film. Before cementing, the meshes are also examined under a microscope.

The nickel mesh is finer than the spatial resolution of the camera, but the coarser meshes and the grid produce a shadow at the focal plane. In an operating camera the XUV image moves continuously over the focal plane. The shadow therefore causes a known modulation of the signal in time and place, which must be allowed for in the final result. Because of the

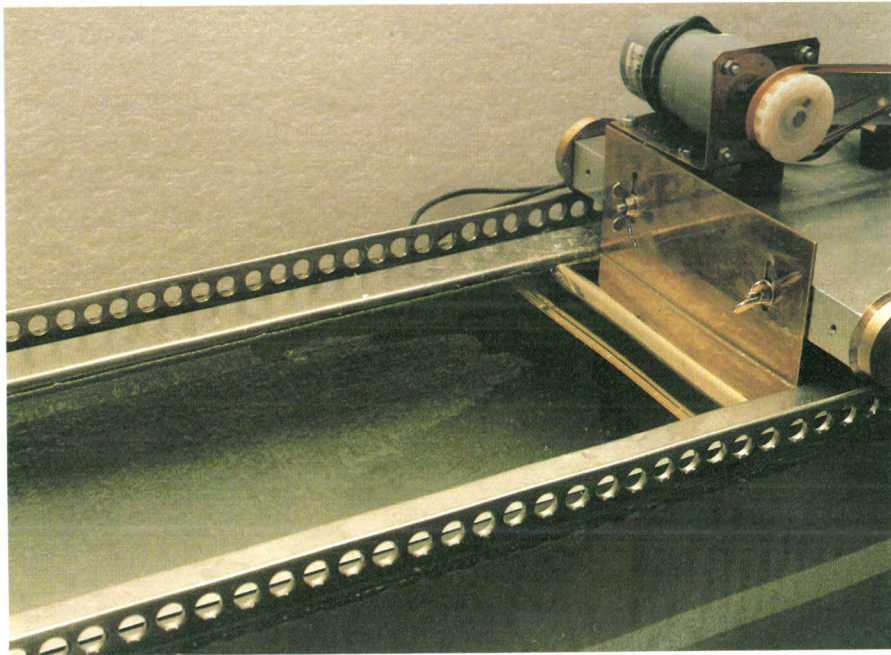


Fig. 5. Film-casting device. A film floating on the water surface is just discernible. The spreader is about 10 cm wide and fits between the rails with a small clearance. The water tank is 50 cm long, 15 cm wide and 20 cm deep.

Supporting meshes

Ultra-thin films of this size cannot stand up to any appreciable pressure difference; they have to be supported. The plastics we have used have a tensile strength of about 5×10^7 N/m, and it follows from a simple calculation that the holes of a supporting mesh must be smaller than 0.1 mm if the film is to withstand a pressure of 10^5 Pa (1 bar). We used a nickel mesh of the type used for storage grids in cathode-ray storage tubes; it is 5 μ m thick and has a pitch of 50 μ m. This fine mesh is in turn supported by two stainless-steel meshes 75 μ m thick and with a pitch of 1 mm, mounted on an aluminium grid with 1 mm bars at a spacing of 20 mm (*fig. 6*). The meshes and the aluminium grid are cemented together at the edges and at the bars of the aluminium grid. The adhesive attaching the

continuous movement the main effect of the shadow is a reduction of the transmission. For the meshes and the grid together the transmission amounts to about 0.40.

Conductive coating

The conductive coating on the inside of the film is applied in the form of an evaporated aluminium film, 7 nm thick. A film of this thickness does not impair the XUV transmission except for a narrow region around 17 nm (the *L* absorption edge of aluminium; see *fig. 1g*), where the absorption reaches a maximum of 25%. The sheet resistivity of the coating is about 30 Ω per square, which is low enough for effectively removing the positive charges from the drift space.

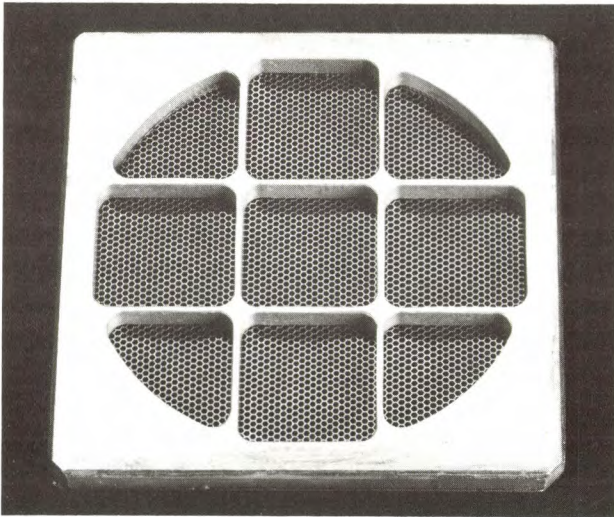


Fig. 6. The aluminium support grid with stainless steel meshes. The distance between the holes in the meshes is 1 mm. The meshes serve to support a nickel mesh (pitch 50 μm , thickness 5 μm), which in turn supports the 'Lexan' film.

Strength, gas-diffusion leakage and transmission

We have set up our own standard tests for the strength and gas-diffusion leakage of the windows. After a preliminary check on the leakage rate, the window is subjected to pressure cycling for about a hundred cycles, with the pressure varying from 0 Pa to a maximum that increases progressively from 20 kPa (0.2 bar) to 100 kPa (1 bar). The leakage rate for a pressure difference of 20 kPa after the cycling must not exceed 30 cm^3 of gas at 1 Pa ($3 \times 10^{-4} \text{cm}^3$ at 1 bar) at room temperature per cm^2 of the window and per second. About 25% of the windows made from selected films and meshes meet this specification.

This refers to films with a thickness of 300 nm. We also tried to make windows with films 150 nm thick,

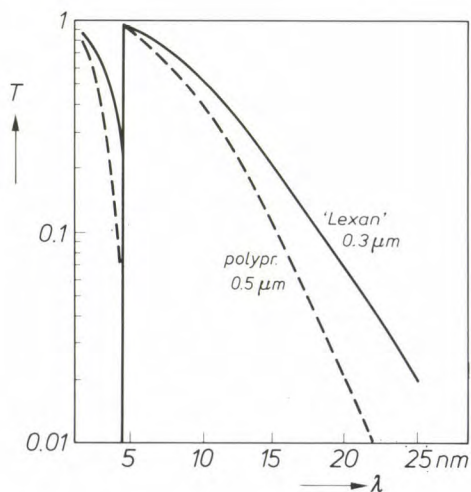


Fig. 7. Calculated transmission of 0.3 μm 'Lexan' and of 0.5 μm polypropylene.

but did not succeed. We believe that this is due to the mounting on the support structure rather than to the films themselves, since the excessive leakage was invariably found to be at the edges or near the bars of the aluminium grid.

Fig. 7 gives the transmission in the XUV range for a 0.3 μm 'Lexan' film and, for comparison, that of an 0.5 μm polypropylene film. The curves were calculated from the weighted sum of the absorption cross-sections of the constituent atoms. The curves in *fig. 1e* and *f* are very simplified summaries of *fig. 7*. We note that in the range 10 to 100 nm these results can only be very approximately correct, because above the carbon K absorption edge at 4.4 nm the L electrons start to play a dominant role in the absorption process; this is because the L electrons are not uniquely localized at

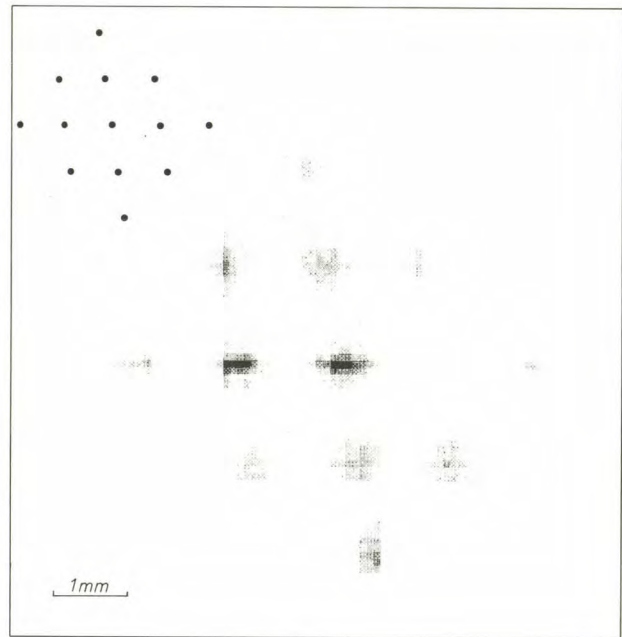


Fig. 8. XUV image at 25 nm of a multi-pinhole mask (see inset), taken with a camera fitted with a 'Lexan' window. The image consists of elements (pixels) of 0.1 mm. The spatial resolution is about 0.6 mm (full width half maximum). The image 'spots' do not have equal brightness because of non-uniformity in the incident beam and shadowing by the window-support structure. Because of this shadowing one hole is completely lost.

the atoms, as was assumed in the summation of absorption cross-sections. From provisional measurements at 24.3 nm, 25.6 nm and 30.4 nm (lines from the Lyman series for ionized helium) it seems to follow that at these wavelengths the calculated values ought to be reduced by a factor of 0.7. It may be supposed that the correction will become less significant at shorter wavelengths. To arrive at the transmission of the complete window, the reduction due to the meshes and the grid of the window-support structure must be taken into account.

The window in use

A series of measurements on the camera shown in fig. 4 have shown that the image formation is very uniform over the full aperture (33 cm^2), apart from the shadowing caused by the meshes and the grid, and that the signal in each pixel is proportional to the number and energy of the incident photons. The performance of the camera at 25 nm with the 'Lexan' window is demonstrated in fig. 8. The same camera, but with a polypropylene window, is also suitable for the 0.2-10 nm wavelength range ^[1].

Concluding remarks

The camera was originally developed for possible employment in an EXUV satellite (Extreme ultraviolet/X-ray sky-survey mission), for which a project study was undertaken by the European Space Agency (ESA). This project, with a planned launching date in 1986, was not selected as part of the mandatory ESA

science programme. The camera may possibly be used later in one of the space laboratories due to be launched by the Space Shuttle during the eighties. The films that we have developed will be used as filters in the EXOSAT (European X-ray Observatory Satellite, with a planned launching date in October 1982).

Summary. Interest in satellite observations in the very soft X-ray waveband (the XUV range) has greatly increased since the discovery that the visibility in this range is much greater than was first thought. For observing objects of same extent — supernova remnants, planetary nebulae and other hot plasma regions — the most suitable instrument has been found to be the proportional counter with spatial resolution, placed at the focal plane of an X-ray optical system. A focal-plane camera of this type has been developed by the Cosmic Ray Working Group at Leiden. The conventional entrance windows of proportional counters are not transparent in the XUV range. A special window has been developed in cooperation with the Philips X-ray Tubes Laboratory. It consists of a film of a polycarbonate, 'Lexan', $0.3 \mu\text{m}$ thick. It is produced by 'casting' a 70 nm film of this material on a water surface, and by laminating four or five such films. The film is supported in the window by a structure of meshes and a grid. One in four of the windows made from selected films and the meshes meet the specifications for strength and gas-diffusion leakage. The effective area of the window is 35 cm^2 .

Growth of alkali-antimonide films for photocathodes

P. Dolizy

Photocathodes are widely used nowadays in many practical devices such as image tubes, radiation detectors and photomultipliers. A material often used for the photocathode film is an alkali antimonide. The development of most conventional photocathodes is still more or less empirical. This article describes studies on alkali-antimonide films that may lead to a better understanding of the growth processes and hence to improvements in the manufacturing technology, with better reproducibility and photosensitivity.

General features of photocathodes

A photocathode is a metallic or semiconducting cathode that will give photoemission: it emits electrons on illumination if the photon energy is equal to or greater than a certain value, the photoelectric threshold energy [1]. If a photocathode forms part of an electrical circuit, its photoemission produces a measurable current. The ratio of this current to the incident luminous flux, the photosensitivity, is determined by the efficiency of each of the three stages in the photoemission process: the excitation of photoelectrons, their movement through the material to the vacuum interface, and their escape into the vacuum.

The excitation of electrons by photons of a given energy depends on the extinction coefficient, the refractive index and the thickness of the photocathode film. On their way to the surface the excited electrons may lose energy through interactions with other electrons or with lattice phonons and through recombinations with holes. This means that only some of the excited electrons will reach the surface. The escape probability $P(W,x)$ of a photoelectron excited at a distance x from the vacuum interface can be described to a first approximation by

$$P(W,x) = P(W,0) \exp(-x/L), \quad (1)$$

where W is the photon energy and L is called the escape depth. This equation is only valid if the thickness of the photocathode film is large compared to the escape depth: For photon energies lower than 2.5 eV the energy losses of the photoelectrons are mainly due

to the electron-phonon interactions. At these energies both $P(W,0)$ and L decrease with decreasing W , and $P(W,0)$ becomes zero at the photoelectric threshold energy.

For the conversion of visible light, including the near infrared and the near ultraviolet, the highest photosensitivity is obtained with the well-known silver-oxygen-caesium cathodes and with photocathodes based on semiconductors such as various alkali antimonides or some of the III-V-compounds [1]. The alkali antimonides absorb strongly in the visible range. The thickness of photocathode films of these materials is therefore made small, so that optical interference takes place in the film. An appropriate surface treatment can be applied to ensure that the threshold energy is not much higher than the band gap.

In conventional photocathodes the bulk material often has a cubic crystal structure and is a p-type semiconductor. To lower the threshold energy the material is coated with a thin n-type film. The presence of such a film produces band-bending in the valence and conduction bands, as shown schematically in *fig. 1* for the alkali antimonide Na_2KSb . The energy difference between the vacuum level and the bottom of the bulk conduction level, the 'effective electron affinity' decreases considerably here (0.7 eV), because of the presence of n-type surface states. As a result the threshold energy decreases from 2.0 eV to 1.3 eV.

P. Dolizy is with Laboratoires d'Electronique et de Physique Appliquée (LEP), Limeil-Brévannes, Val-de-Marne, France.

[1] See for example A. H. Sommer, *Photoemissive materials*, Wiley, New York 1968, and W. E. Spicer, *J. appl. Phys.* 31, 2077, 1960.

For some photoemitters the effective electron affinity may even become negative: the vacuum level is then lower than the bottom of the conduction band in the bulk of the material. Electrons in this band then require no extra kinetic energy to escape into the vacuum. Materials with a negative effective electron affinity can have a high photosensitivity because of their large escape depth; one example is p-GaAs coated with Cs [2]. In this article, however, such materials will not be discussed [3].

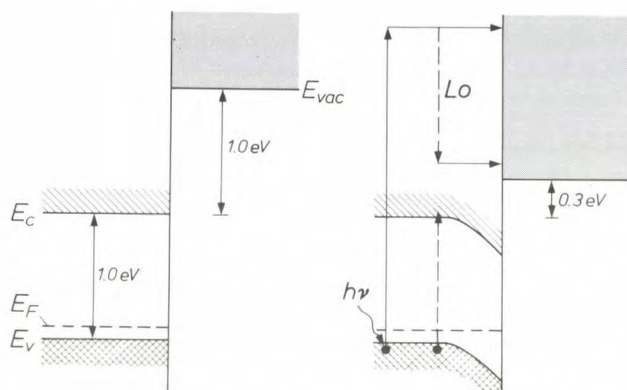


Fig. 1. Energy bands at the surface of an Na_2KSb film. E_v upper edge of the valence band in the bulk of the material. E_F Fermi level. E_c lower edge of conduction band in the bulk. E_{vac} potential energy of an electron in vacuum. For p-type material (left) the photoelectric threshold energy, $E_{vac} - E_v$, is equal to 2.0 eV. Adsorption of caesium at the surface gives rise to n-type surface states. This has the effect of bending the valence and conduction bands downward near the surface (right), so that the effective electron affinity, $E_{vac} - E_c$, has fallen to 0.3 eV and the threshold energy is only 1.3 eV. L_o loss of electron energy as a result of interaction with a valence electron. An electron excited by a photon with an energy $h\nu$ of about 2.5 eV can either escape directly to the vacuum, or may first create a new electron-hole pair, and then escape.

Photocathodes have many practical applications, e.g. in image tubes, radiation detectors and photomultipliers. In most applications the photocathodes consist of a polycrystalline film deposited on a substrate. Various substrates can be used: metals such as palladium, alloys such as nichrome, or insulators such as glass. In addition, the substrate surface can be large or small, flat or curved, without affecting the deposition of the photoemissive film. The light may be incident from the vacuum interface (the reflection mode), or from the substrate interface (the transmission mode). In the reflection mode the substrate does not have to be transparent and the film thickness is not critical. In the transmission mode, however, the substrate must be transparent and the film thickness has a critical optimum value: on the one hand the thickness should be not much larger than the escape depth, whereas on the other hand the film should be thick enough to ensure a sufficient absorption of the incident radiation via the optical interference inside the film. In most applications photocathodes are used in the transmission mode.

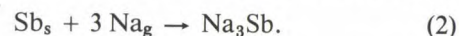
The scope of the investigations

Most conventional photocathodes are produced more or less empirically. The growth is usually monitored only by measuring the photosensitivity, in the mode to be used, at various stages of the growth. This does not give much information about the various processes that occur during the growth. An understanding of these processes may be very useful, however, for improving the reproducibility of the growth and the photosensitivity of the photocathodes, or for automating the growth procedure. In our investigations on conventional alkali-antimonide photocathodes we therefore carried out a number of additional measurements during the growth.

We applied the alkali-antimonide films in a vacuum chamber by depositing antimony atoms on a substrate in the presence of one or more alkali metals. To find out how the film thickness increased during the growth we measured the optical reflection and transmission. We used atomic absorption spectrometry (AAS) to determine the alkali partial vapour pressures, and also measured the photosensitivity in both the reflection and the transmission modes, to find out how the escape depth and escape probability varied during the growth. These measurements were found to be of great value for examining and improving the photocathode films during the growth in practical tubes. Before considering these studies more closely, let us first look at the growth of the alkali-antimonide films in more detail.

Growth of alkali-antimonide films

Alkali-antimonide films for photocathodes are usually evaporated from separate sources ('dispensers') for antimony and alkali. The dispensers are located inside the evacuated tube being manufactured, and they are heated electrically under external control. A thin film of antimony evaporated on to the substrate reacts exothermally with the alkali-metal vapour to produce various semiconductor compounds. The chemical reaction between solid antimony and sodium vapour, for example, can be represented by:



This equation provides no information about the existence of intermediate compounds that could be formed during the reaction.

The actual composition of an alkali-antimonide film depends on the partial alkali pressures in the vapour, the vapour temperature and the substrate temperature, as well as on the substrate state, the impurity content and the growth rate. Since alkali metals react very rapidly with oxygen, carbon dioxide and

water, it is essential that the growth takes place in a very clean vacuum envelope. The number of compounds formed increases with the atomic number of the alkali metal [4] and with the number of different metals in the vapour. The control of the partial alkali pressures is therefore very important for obtaining a photocathode film of the desired composition. At a partial pressure p the number of atoms n_i incident on the surface per cm^2 per second is given by [5]

$$n_i = p/(2\pi mkT_v)^{\frac{1}{2}}, \quad (3)$$

where m is the mass of an alkali-metal atom, k is Boltzmann's constant and T_v is the temperature of the vapour. Some of the adsorbed atoms will leave the surface again, however. Assuming that this desorption is a first-order process, the number of atoms n_1 leaving the surface per cm^2 per second is [5]

$$n_1 = n_{\text{ad}} A \exp(-E_b/kT_s), \quad (4)$$

where n_{ad} is the number of atoms adsorbed per cm^2 of the surface and A is the partition function, which is nearly constant to a first approximation and has a value between 10^{14} and 10^{16} s^{-1} . E_b is the binding energy between the metal and the substrate; T_s is the substrate temperature. When the evaporated film is in equilibrium with the vapour, the value of n_1 is equal to that of n_i , so that n_{ad} is given by

$$n_{\text{ad}} = \frac{p/(2\pi mkT_v)^{\frac{1}{2}}}{A \exp(-E_b/kT_s)} \quad (5)$$

If p , T_v , A , E_b and T_s are known, this relation can be used to calculate the number of atoms adsorbed on the substrate. In *fig. 2* the calculated monolayer fraction f for potassium on glass is plotted against the substrate temperature, for three different combinations of the partial vapour pressure of potassium and the binding energy. The value of f decreases strongly with increasing substrate temperature. The temperature above which f is virtually zero is considerably higher for a higher partial vapour pressure of potassium and higher binding energy.

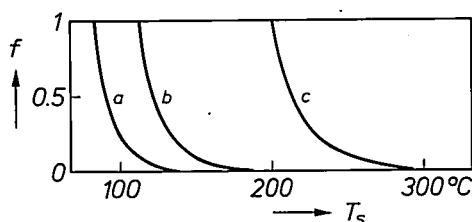
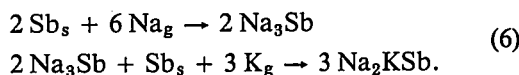


Fig. 2. Calculated monolayer fraction f for potassium evaporated on to a glass substrate, as a function of the substrate temperature T_s . The three curves correspond to three different combinations of the partial potassium vapour pressure p_K and the binding energy E_b between potassium and the glass. *a)* $p_K = 1.3 \times 10^{-4}$ Pa, $E_b = 92$ kJ/mol. *b)* $p_K = 1.3 \times 10^{-3}$ Pa, $E_b = 92$ kJ/mol. *c)* $p_K = 1.3 \times 10^{-4}$ Pa, $E_b = 125$ kJ/mol.

In small vacuum chambers the speed of the vacuum pump must be taken into account [6]. If the binding energy is low, the alkali partial pressure is also affected by the pumping time. The binding energy depends on the interaction between the alkali metal and the substrate. For example, sodium reacts more strongly than potassium or caesium with a glass substrate. Other important factors for the binding energy are the purity of the substrate and the presence of other materials on the surface, such as alkali metals and antimony, which have been adsorbed in a previous stage of the processing. If the vapour contains different alkali metals, the adsorption and desorption are associated with chemical exchanges of the alkali atoms on the substrate surface.

The composition of the compounds formed depends on the adsorption of the various alkali metals and on the type of reaction occurring at the surface. In some cases different reaction mechanisms give the same final compounds. A bialkaline antimonide of composition $(A,B)_3\text{Sb}$, for example, can be obtained from the $A_3\text{Sb}$ compound by a partial substitution of A by B, or by the addition of solid Sb and gaseous B via a diffusion process. The addition mechanism for the growth of an Na_2KSb film can be described as follows:



The measurements of the partial alkali-vapour pressures and the optical properties of the film during the growth have shown that such addition processes are always accompanied by substitutions inside the film. At a typical temperature of about 200°C and typical partial pressures of 10^{-4} Pa potassium is readily substituted by sodium [7].

An important property of alkali-antimonide films is their crystalline nature. To obtain a high photosensitivity, the microcrystals must be larger than the escape depth of the photoelectrons excited inside them. The material then behaves as a single crystal for the photoelectrons. Since the escape depth must be of the same order as the film thickness, the growth conditions should generally favour the formation of the largest possible microcrystals.

[2] J. van Laar and J. J. Scheer, Philips tech. Rev. 29, 54, 1968.

[3] More information on photoemitters with a negative 'effective electron affinity' can be found in the article by J. van Laar and J. J. Scheer [2], or in the six articles in Acta Electronica 16, No. 3, 1973 (pp. 215-271).

[4] E. Kinsky, Adv. in Electronics and Electron Phys. 33A, 357, 1972.

[5] H. Mayer, Vakuum-Technik 4, 1, 1955.

[6] J. P. Hobson, 1961 Trans. 8th Nat. Vacuum Symp., Washington D.C., Vol. 1, page 26.

[7] P. Dolizy, O. De Luca and M. Deloron, Acta Electronica 20, 265, 1977.

The texture of the films depends on their composition and crystal structure. A strong preferential orientation has been observed for hexagonal Na_3Sb , with the (00.1) plane parallel to the surface^[7]. On the other hand, films of hexagonal K_3Sb often have a poor crystal structure. An Na_2KSb film obtained by evaporating K and Sb on to a structured Na_3Sb film is also strongly oriented. In this case the (111) plane of the cubic crystal structure^[8] is parallel to the surface. Further evaporations of K and Sb lead to the formation of the hexagonal NaK_2Sb phase, coexisting with the initial Na_2KSb phase. This hexagonal phase again has a texture with the (00.1) plane parallel to the surface. It is assumed that during the growth the bulk of the films contains only one or two 'neighbour' compounds from the series Na_3Sb , Na_2KSb , NaK_2Sb and K_3Sb , when the antimony atoms are completely saturated by the alkali atoms.

Determination of properties during growth

In our investigations we determined the properties of interest during the growth of alkali-antimonide films on glass by measuring the atomic absorptions of the alkali vapours, the optical reflection and transmission of the film, and its photosensitivity in the reflection and transmission modes. We interpreted the results of these measurements in terms of the partial alkali-vapour pressures, the extinction coefficient, refractive index and thickness of the film, and also the escape depth and escape probability of the photoelectrons. These measurements will now be discussed and the derivation of the properties during the growth will be explained.

Optical measurements

Fig. 3 is a diagram of the equipment used for measuring the atomic absorptions of sodium and potassium in the alkali vapour. The growth takes place in a vacuum chamber containing an antimony bead, alkali dispensers, an anode and a glass substrate on which the photocathode film is deposited. Two modulated light beams, one with the correct resonance line for sodium, and the other with the line for potassium, are passed through the vacuum chamber. The transmitted light beam from both beams is detected by a simple arrangement of a focusing lens, an interference filter and a photocell. The electrical signal from this cell is supplied to a synchronous detector, which is only sensitive to signals at frequencies close to that of the modulation of the light beam. To measure the atomic caesium absorption, the transmitted resonance radiation was detected with a monochromator on account of unwanted signals from other elements.

We measured the reflection, transmission and photosensitivity during the growth with the arrangement shown in the diagram of fig. 4. Three modulated light beams were used in these measurements. The first beam is incident on the photocathode film from

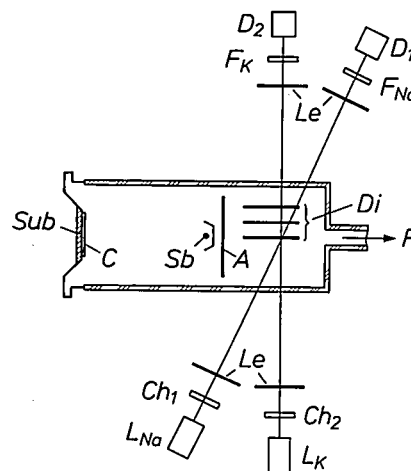


Fig. 3. Schematic arrangement for determining the partial sodium and potassium vapour pressures during the growth of an alkali-antimonide photocathode *C* in a vacuum chamber. Two light beams, originating from hollow-cathode lamps L_{Na} and L_{K} , with the appropriate sodium and potassium resonance lines are modulated by two different choppers Ch_1 and Ch_2 and passed through the vacuum chamber. The measured attenuation for the sodium or potassium resonance radiation is a measure of the partial sodium or potassium vapour pressure in the chamber. Le lenses. *Sub* glass substrate. *Sb* 'bead' of antimony. *A* anode. *Di* alkali dispensers. *P* vacuum pump. F_{Na} and F_{K} sodium and potassium interference filters. D_1 and D_2 detectors, each consisting of a photocell and a synchronous detector tuned to the modulation frequency of one of the two modulated light beams.

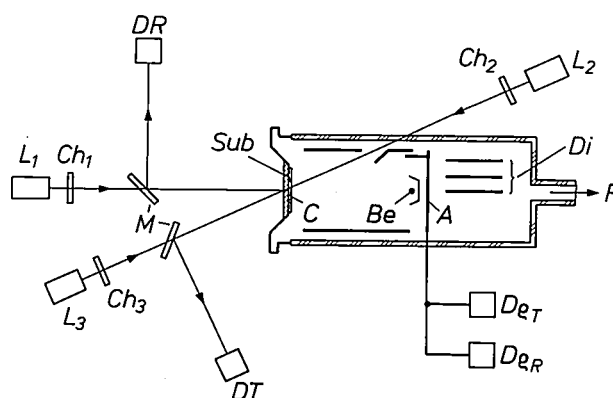


Fig. 4. Schematic arrangement for measuring the optical reflection and transmission modes and the photosensitivity in the reflection and transmission modes during the growth. Three light beams are used in the measurements; they originate from the sources L_1 , L_2 and L_3 , and each beam is modulated at a different frequency by the choppers Ch_1 , Ch_2 and Ch_3 . DR is a detector for the light originating from L_1 and reflected by the photocathode *C*. DT is a detector for the light originating from L_2 and transmitted by *C*. Dq_R and Dq_T are synchronous detectors for the current generated in the vacuum chamber by the light originating from L_2 and L_3 respectively. The photosensitivity of *C* in the reflection and transmission modes can be determined from these currents. *M* half-silvered mirrors. *Sub* glass substrate. *Sb* 'bead' of antimony. *A* anode. *Di* alkali dispensers. *P* vacuum pump.

the substrate side. Measurement of the intensity of the reflected light gives the reflectance of the film. The second beam is incident from the photocathode side, so that synchronous detection of the electric current generated gives the photosensitivity of the film in the reflection mode. The intensity of the light that passes through the film is also measured, giving the optical transmittance of the film. The third beam is again incident from the substrate side; synchronous detection of the current generated gives the photosensitivity in the transmission mode.

Derivation of the properties

The partial vapour pressures of sodium, potassium and caesium can be derived by comparing the observed attenuation of the resonant radiation with measurements on calibration cells with known alkali vapour pressures. Our experimental arrangement allows the partial vapour pressures to be determined simultaneously. An important experimental observation in this respect relates to the exchange of sodium and potassium. For example, when all the glass surfaces of a vacuum chamber have been covered with the two metals, the ratio of the partial pressures p_{Na} and p_K in the non-saturated vapour is only slightly dependent on which dispenser is emitting inside the chamber at that moment. This is illustrated in *fig. 5* for a glass chamber in which the substrate tempera-

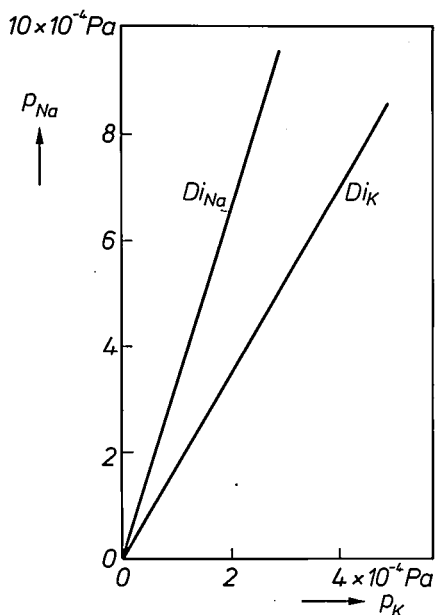


Fig. 5. Sodium partial vapour pressure p_{Na} plotted against the potassium partial vapour pressure p_K in a vacuum chamber in which all the surfaces are covered with the two metals. The substrate temperature is between 200 and 250 °C. If the sodium dispenser Di_{Na} alone or the potassium dispenser Di_K alone is emitting, the partial pressure of the other metal also increases strongly, because of substitution reactions at the walls. The ratio p_{Na}/p_K does not increase greatly when the potassium dispenser is switched off and the sodium dispenser is switched on.

ture is between 200 and 250 °C. The ratio p_{Na}/p_K does not increase much when the potassium dispenser is switched off and the sodium dispenser is switched on. In both situations the ratio does not deviate much from 2. As indicated before, this ratio is strongly dependent on the nature of the surfaces.

The extinction coefficient k and the refractive index n can be derived from the reflection measurements [9]. Since the film thickness d is smaller than the wavelength λ of the incident light, the reflectance R is a complicated function of k , n , d and λ . If R is plotted against d , or against the deposition time t , which amounts to the same thing, for given values of k , n and λ , a number of maxima and minima are obtained; see *fig. 6*. For a given wavelength the value of R at

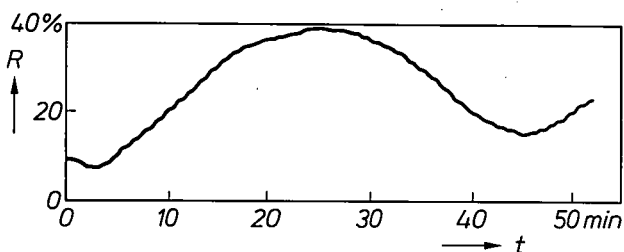


Fig. 6. Reflectance R of an Na_2KSb film at 500 nm, measured as a function of the deposition time t .

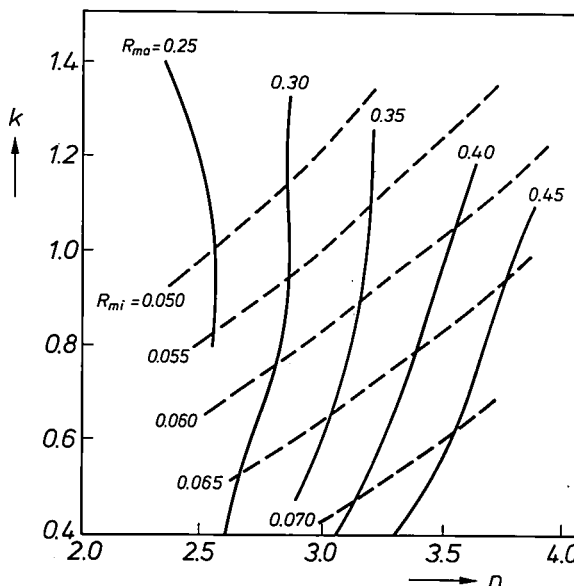


Fig. 7. Calculated isoreflectance curves in the plane of the extinction coefficient k and the refractive index n , for the first minimum R_{mi} and the first maximum R_{ma} in the $R(t)$ -curve (see for example *fig. 6*). The values of k and n can be derived by comparing these isoreflectance curves with the measured values of R_{mi} and R_{ma} .

[8] J. J. Scheer and P. Zalm, Philips Res. Repts 14, 143, 1959.
 W. H. McCarroll, Phys. Chem. Solids 16, 30, 1960.
 [9] S. E. Webber and S. R. Scharber Jr., Appl. Optics 10, 338, 1971.
 V. E. Kondrashov and A. S. Shefov, Engl. Transl. Bull. Acad. Sci. USSR, phys. Ser. 28, 1349, 1964.

these extreme values can be expressed in terms of k and n . This permits the 'isorefectance' curves to be drawn, which give the theoretical reflectances at these extreme values in a plane of k - and n -coordinates. These curves are shown in *fig. 7* for the first minimum and the first maximum on the $R(t)$ -curve. The k - and n -values can then be deduced by comparing the curves with the extreme values of the reflectance observed during the growth (*fig. 6*).

If the k - and n -values are known, the optical reflectance at the wavelength λ can be calculated as a function of d ; see *fig. 8*. A comparison of the calculated $R(d)$ -curve with the measured $R(t)$ -curve (*fig. 6*) gives the film thickness $d(t)$ and hence the growth rate $r(t)$ as well. A more elaborate way of determining $d(t)$ and $r(t)$ depends on the calculation of R as a function of λ for a large number of d -values. *Fig. 9* shows how the theoretical curve for $R(\lambda)$ varies with d . Comparison with the experimental curves for $R(\lambda)$ measured during the growth allows the values of $d(t)$ and $r(t)$ to be obtained.

In determining the escape depth L and the escape probability $P(W,0)$ we have to remember that L can be of the same order as d . This implies that some of the photoelectrons will travel to the interface with the substrate, so that eq. (1) is no longer true. If there is no recombination at the interface these electrons can

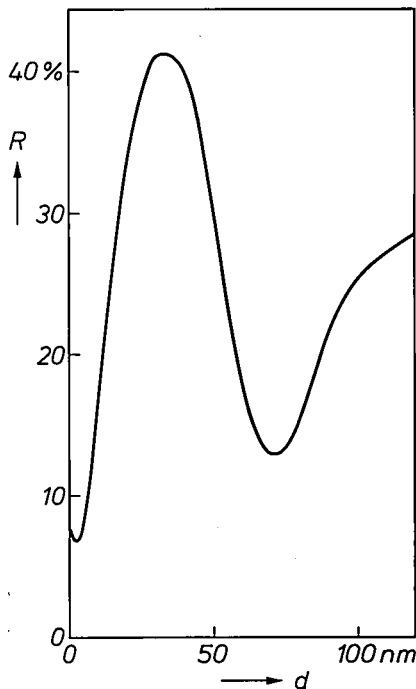


Fig. 8. Calculated reflectance R at 500 nm for a film with $k = 0.8$ and $n = 3.5$, as a function of the film thickness d . A comparison of the $R(d)$ -curve, calculated from the known k - and n -values of a film, and the measured $R(t)$ -curve (see for example *fig. 6*) gives the film thickness and the growth rate as a function of time.

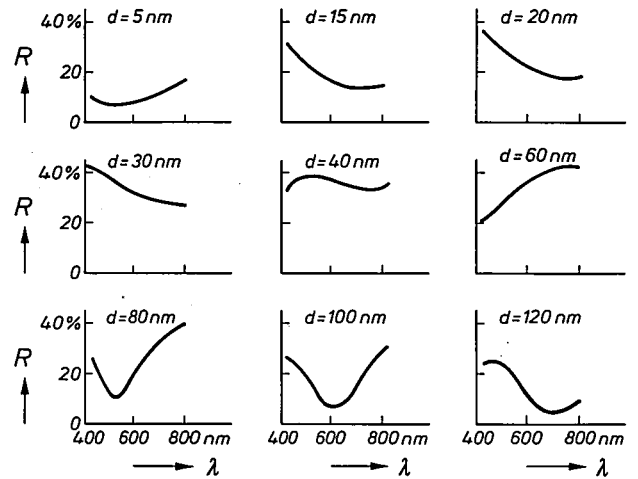


Fig. 9. Calculated optical reflectance R as a function of the wavelength λ , for different thicknesses of an Na_2KSb film. Comparison of these curves with the $R(\lambda)$ -curves measured during the growth provides another method for determining the film thickness and the growth rate as a function of time.

still escape on the vacuum side. The escape probability for a given photon energy W is then given by

$$P(W,x) = P(W,0) \frac{\cosh(d/L - x/L)}{\cosh(d/L)}. \quad (7)$$

However, if the recombination velocity at the substrate interface is very high, the electrons arriving at this interface can no longer escape on the vacuum side. The escape probability is then given by

$$P(W,x) = P(W,0) \frac{\sinh(d/L - x/L)}{\sinh(d/L)}. \quad (8)$$

In practice the recombination velocity will have a finite value, so that $P(W,x)$ is a very complicated function, which will not be discussed here^[10].

The values of L and $P(W,0)$ during the growth of the film on a substrate with known optical properties can be determined from measurements of the photo-sensitivity in the reflection mode (ϱ_R) and in the transmission mode (ϱ_T). If A_R and A_T are the distribution functions, dependent on k , n and λ , of the light absorbed in the film in the reflection and transmission modes, ϱ_R and ϱ_T can be expressed as

$$\begin{aligned} \varrho_R &= \int_0^d A_R(k,n,\lambda,x) P(W,x) dx \\ \varrho_T &= \int_0^d A_T(k,n,\lambda,x) P(W,x) dx. \end{aligned} \quad (9)$$

If k , n and d are known, the value of L at a given wavelength λ can be derived from the ratio ϱ of ϱ_R to ϱ_T . In *fig. 10* the calculated ratio ϱ is plotted against L for various combinations of k , n , λ and d . The value of L is obtained by comparing the experimental value of ϱ , determined by relative measurements of ϱ_R

and ϱ_T , with the theoretical curves. A necessary condition for the combination of parameters is that ϱ should vary sufficiently with L . In the situations of fig. 10 a good determination is possible with blue light if the thickness is larger than about 40 nm. If L is known, the escape probability $P(W,0)$ can be obtained from an absolute measurement of ϱ_R or ϱ_T .

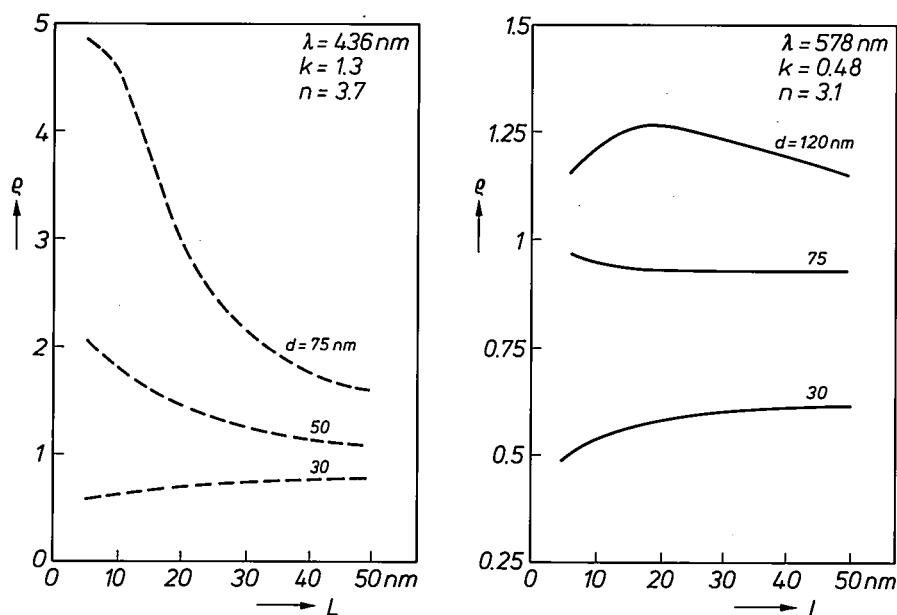


Fig. 10. Calculated ratio ϱ of the photosensitivity in the reflection mode to the photosensitivity in the transmission mode, as a function of the escape depth L of the photoelectrons, for various values of the film thickness d , wavelength λ , extinction coefficient k and refractive index n . The escape depth L can be derived from the $\varrho(L)$ -curve calculated from the known values of λ , d , k and n , and the measured value of ϱ , provided that the calculated curve varies sufficiently with L .

Growth monitoring

The evaluation of the properties during the growth of a photocathode provides useful information for monitoring this growth. Let us first consider coating a glass substrate with antimony, and examining the film by reflection measurements. At a substrate temperature above 140 °C the sticking coefficient of antimony on a very clean glass substrate is close to zero if the antimony flux is no larger than 10^{13} at/cm²s. In the presence of an alkali vapour, however, this coefficient approaches unity for substrate temperatures between 140 and 250 °C. When the growth of a photocathode film is started at 200 °C, for example, antimony must therefore be evaporated in the presence of an alkali vapour, to permit the antimony to be satisfactorily deposited on the glass and the chemical reaction between solid antimony and the alkali vapour to start.

During the chemical reaction the film thickness increases. To obtain a thicker film when this reaction has been completed, more antimony must be deposited and saturated again with alkali vapour. The growth of the film can easily be monitored by optical reflection or transmission measurements, as described. As an

example, fig. 11 shows the reflectance during the growth of an Na₃Sb film¹ on a glass substrate. In this case the growth occurs via a number of separate antimony evaporations, whereas sodium is continuously present in the vapour. The theoretical curve of the reflectance as a function of the film thickness (fig. 11a) is similar to that of fig. 8. If the reflectance at the

beginning of the growth is measured in detail, the effects of the various antimony depositions can be seen; see fig. 11b. The change of reflection during one deposition is shown in fig. 11c. After the antimony dispenser has been switched on, the film thickness first increases steeply, since antimony is deposited in the presence of sodium vapour. When the evaporation is temporarily stopped, the antimony atoms still present on the film react with the sodium atoms to form a compound close to Na₃Sb. This chemical reaction causes a further increase in the thickness. When the reaction has been completed, the thickness remains constant until the next antimony evaporation begins. Similar results have also been obtained in growing K₃Sb and Cs₃Sb layers.

The rate of the reaction depends upon the quantity of antimony and alkali metal at the substrate surface. If there is an excess of antimony, the reaction rate is controlled by the local amounts of alkali. The rate of growth of the film during the chemical reaction is then proportional to the partial alkali vapour pressure.

[10] See C. Piaget, Thèse No. 1892, Université de Paris Sud, 1977.

However, if more alkali atoms than antimony atoms arrive at the substrate, the growth rate is proportional to the deposition rate of the antimony.

The growth of Na_2KSb films can be studied in a similar way to that of the monoalkali-antimonide films. To obtain an Na_2KSb film with the highest photosensitivity the addition mechanism of eq. (6) is found to be the most satisfactory. In this mechanism, the number of antimony evaporations in the presence of sodium should be twice the number of evaporations in the presence of potassium. In practice, however, the situation is more complicated, since potassium in the deposited film is partly substituted by sodium, depending on the substrate temperature and the ratio of the sodium and potassium vapour pressures^[7]. In a tube these vapour pressures are never negligible, because of substitution reactions on the

in the presence of sodium are required than the addition mechanism would indicate.

Depending on the film thickness, the substitution of potassium by sodium can have a marked effect on the reflection. Fig. 12 shows how the reflectance at 520 nm varies with the film thickness for the compounds Na_2KSb and Na_3Sb . If the potassium of Na_2KSb is substituted by sodium, the reflectance at a thickness of about 30 nm (the first maximum) decreases; whereas at a thickness of 80 nm (the second minimum) the reflectance increases. Thus, to obtain a pure Na_2KSb film with no Na_3Sb , the reflectance at the first maximum should be as high as possible and the reflectance at the second minimum should be as low as possible.

The homogeneity of a film can be tested by repeating the determination of k and n at various stages of the growth. A film can be said to be homogeneous if

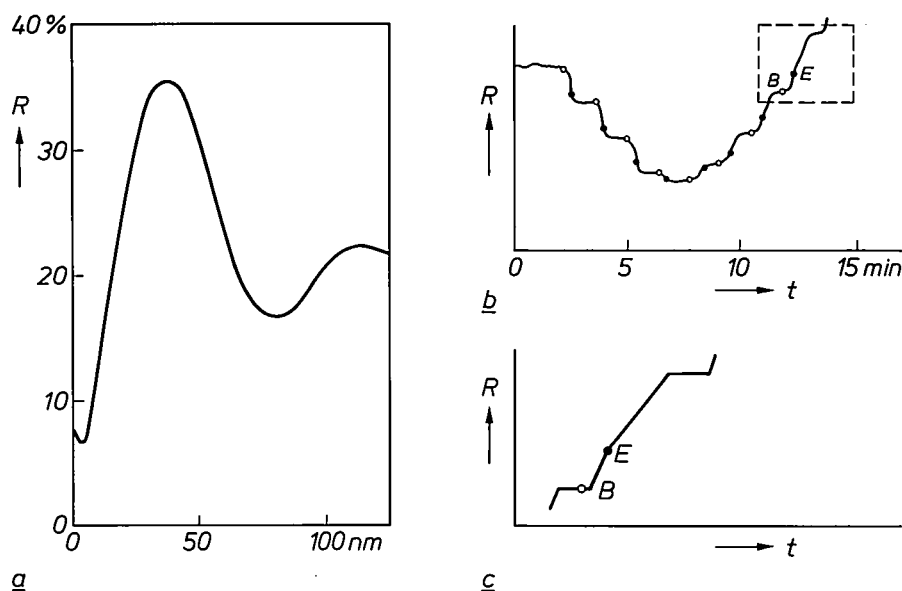


Fig. 11. Calculated and measured reflectance at 520 nm during the growth of an Na_3Sb film. a) Calculated reflectance R as a function of the thickness d of a film with $k = 1.1$ and $n = 3.3$. b) Detail of the reflectance R (in arbitrary units) measured at the start of the growth, as a function of the evaporation time t . The start and finish of the antimony evaporations in the presence of sodium are denoted by open and filled circles. The effect of these evaporations on the reflection is clearly demonstrated. c) Detail of (b), showing more clearly the effect of a single antimony evaporation, starting at B and finishing at E . During the evaporation the reflection increases with time because of the increase in film thickness as a result of the deposition of antimony in the presence of sodium vapour. After the antimony evaporation has been stopped the reflectance first continues to increase, because of the formation of a compound of similar composition to Na_3Sb , owing to the reaction of the residual antimony with sodium. As soon as all the antimony has reacted, the reflection remains virtually constant until the next evaporation.

walls. When both alkali metals are present on the tube walls, heating one of the two alkali dispensers produces a simultaneous pressure increase of the other alkali metal (fig. 5). The difference from the rate of reaction with antimony should therefore also be taken into account. We have found that at the same pressure sodium reacts 1.3 times as fast as potassium. As a result of the substitution reactions and the higher reaction rate of sodium, fewer antimony evaporations

the k - and n -values obtained from different pairs of extreme values in the reflectance-thickness curve are the same. On the other hand, a change in k or n indicates that there have been some changes in the film during the growth. We have found that a low and uniform growth rate and a constant film composition during the growth favour the formation of homogeneous Na_2KSb films. The homogeneity is also better as the growth due to addition reactions increases. In-

creasing the homogeneity gives higher photosensitivities because of the reduction in crystal imperfections, grain boundaries and stresses. In addition, the dimensions of the microcrystals may become comparable to the film thickness, so that the photoelectrons can have a large escape depth.

The presence of unwanted deposits on a grown film can be detected easily, since they usually have a marked effect on the reflection. Fig. 13 shows how drastically the reflectance at 520 nm for different thicknesses of Na_2KSb film can vary because of additional deposits of pure antimony and Na_3Sb .

To obtain the highest possible photosensitivity the growth of a film must be stopped at the optimum film

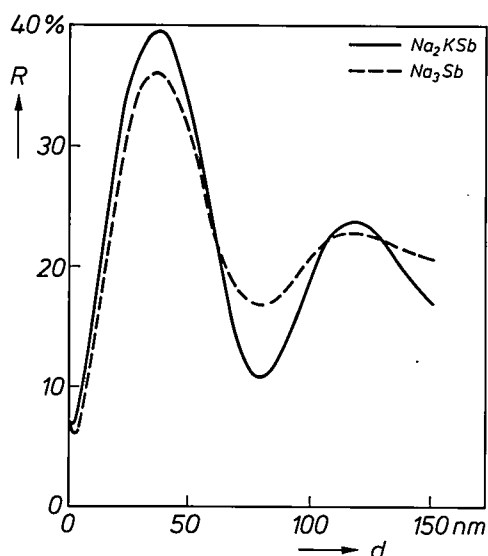


Fig. 12. Measured reflectance R at 520 nm for an Na_2KSb film and an Na_3Sb film, as a function of the thickness d . The reflectance changes drastically when Na_2KSb is converted into Na_3Sb .

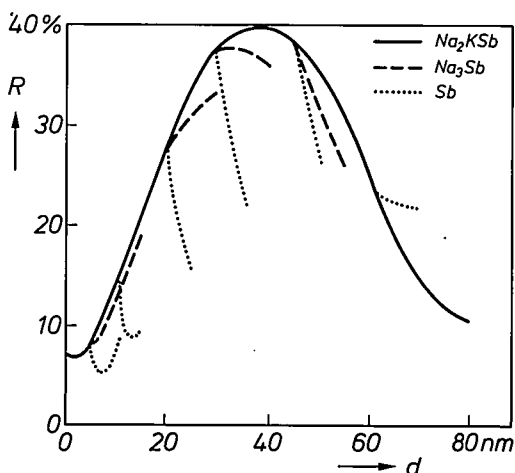


Fig. 13. Measured reflectance R at 520 nm for an Na_2KSb film as a function of the thickness d , and the effects of antimony and Na_3Sb depositions on the reflection. These depositions were made on Na_2KSb films of various thicknesses: 5, 10, 20, 30, 45 and 60 nm. The antimony deposition has a particularly marked effect on the reflection.

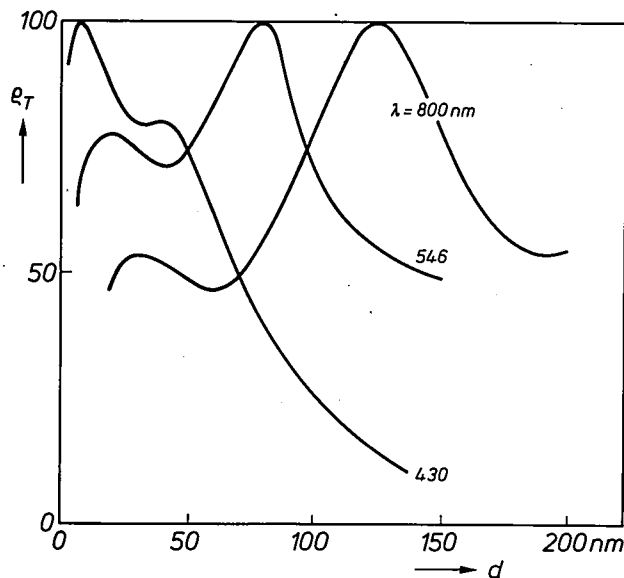


Fig. 14. Photosensitivity e_T in the transmission mode for different wavelengths, as a function of the thickness d of an Na_2KSb film. The photosensitivities are given in relative units; the maximum of each curve has been set equal to 100. The optimum thickness depends strongly on the wavelength λ of the incident light: about 10 nm for $\lambda = 430$ nm, about 80 nm for $\lambda = 546$ nm and about 120 nm for $\lambda = 800$ nm.

thickness. This optimum thickness depends strongly on the wavelength of the incident radiation. In fig. 14 the photosensitivity is shown as a function of the thickness of an Na_2KSb film for three different wavelengths. The optimum thickness increases with the wavelength: 10 nm for $\lambda = 430$ nm, 80 nm for $\lambda = 546$ nm and 120 nm for $\lambda = 800$ nm. If the light to be detected has a low intensity, as in night-vision applications, the optimum thickness of the photocathode film is 120 ± 10 nm.

Surface treatment of Na_2KSb films

The photosensitivity of a pure and homogeneous Na_2KSb film remains virtually constant for several hours at 200 °C. It increases by a factor of about 3 if the temperature falls to 20 °C, but nevertheless remains fairly low [11]. A considerably higher photosensitivity can be obtained by applying a surface treatment in which caesium is evaporated on to the Na_2KSb film. The highest efficiencies that we have measured for a test cell after the surface treatment are shown in fig. 15 as a function of the wavelength of the incident light. The curves relate to two Na_2KSb films, one made mainly by substitution reactions and the other mainly by addition reactions. The addition mechanism gives the highest efficiency; the maximum photosensitivity is $705 \mu\text{A}/\text{lm}$ for an Na_2KSb film of about 120 nm [12].

[11] D. E. Persyk, J. L. Ibaugh, A. F. McDonie and R. D. Faulker, IEEE Trans. NS-23, 186, 1976.

[12] The measurements were carried out with no amplification of the signal by an electric field.

A problem with this caesium treatment is that small quantities of sodium and potassium are also present in the evaporation chamber during the evaporation, as we have found from atomic absorption spectrometry. Most of the sodium and potassium atoms come from the chamber walls as a result of the alkali-displacement reactions mentioned earlier. The caesium treatment of an Na_2KSb film is therefore rather difficult to control and to study.

We have carried out special experiments to measure the thickness of this top layer and analyse its chemical composition. An Na_2KSb film was deposited on half of a glass substrate, with the other half covered by a mask during the evaporation. The mask was then removed and the top layer was deposited over the total substrate area. The evaporation was carried out at a temperature such that both parts were identically covered. We can therefore assume that the thickness and composition of the top layer were the same for both halves. The thickness of this top layer was estimated to be less than 0.8 nm from reflection measurements. Complementary Auger experiments [13] confirm this and also indicate that the top layer contains potassium as well as caesium and antimony. The photosensitivity of such a (K,Cs,Sb) top layer is of the order of 1 $\mu\text{A}/\text{lm}$, whereas an Na_2KSb film of about 100 nm coated with this top layer can have a photosensitivity of up to 350 $\mu\text{A}/\text{lm}$.

Investigators have often wondered whether the atoms of the top layer continue to remain on top of the Na_2KSb film or whether they partly diffuse into it. Although conclusive evidence is not yet available, we have sufficient indications to lead us to think that interdiffusion of the alkali atoms between the top layer and the Na_2KSb film occurs along the grain boundaries. The actual surface composition of a caesium-treated Na_2KSb film probably does not therefore correspond to the top layer originally deposited. The highest photosensitivities are obtained if the sodium and potassium pressures are kept as low as possible during the caesium treatment. Slow growth of the photocathode film is also beneficial: it limits the number of grain boundaries so that the alkali diffusion is less significant.

One of the problems encountered in practice is that undesirable substances such as sodium and potassium appear inside the tube during the caesium evaporation.

[13] The equipment used for Auger spectrometry of multi-alkali-antimonide layers has been described in [7].

[14] W. H. McCarroll, R. J. Paff and A. H. Sommer, *J. appl. Phys.* 42, 569, 1971.

A. A. Dowman, T. H. Jones and A. H. Beck, *J. Physics D* 8, 69, 1975.

C. Ghosh and B. P. Varma, *J. appl. Phys.* 49, 4549, 1978.

R. Holtom, G. P. Hopkins and P. M. Gundry, *J. Physics D* 12, 1169, 1979.

L. G. Estella, Thesis No. 8016823, Stanford University, 1980.

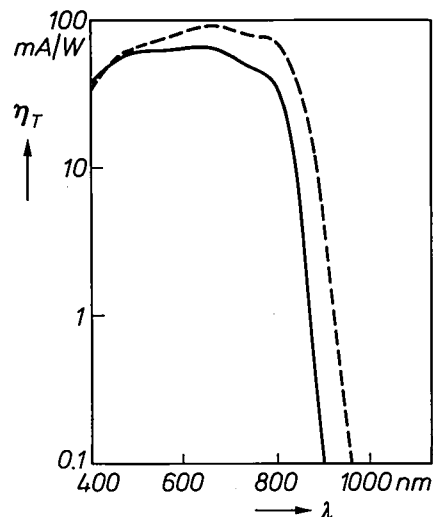


Fig. 15. Photoelectric spectral response η_T for two Na_2KSb films with a thin top layer containing caesium, measured in the transmission mode without electrical amplification. The solid curve was obtained with an Na_2KSb film of thickness about 110 nm, grown mainly by substitution reactions. The maximum photosensitivity ρ_T of this film was about 545 $\mu\text{A}/\text{lm}$. The dashed curve was obtained with an Na_2KSb film of thickness about 120 nm, grown mainly by addition reactions. This film had a significantly higher spectral response, with a maximum photosensitivity of 705 $\mu\text{A}/\text{lm}$.

There have been many investigations of these effects [13] [14], but no unambiguous conclusions can as yet be drawn for the best top-layer composition for high photosensitivity. We have therefore initiated special experiments using the molecular-beam technique to help us to study and improve the top-layer deposition. We have for example been able to demonstrate that a top layer with about one monolayer of caesium gives an appreciably higher photosensitivity than a top layer of the (Cs,Sb) or (K,Cs,Sb) type.

In the investigations described here useful discussions and technical support were given by J. Beltramelli, M. Decaesteker, M. A. Deloron, O. De Luca, F. Grolière, J. Houdard, F. Maniguet, G. Marie and C. Piaget, all of LEP.

Summary. The growth of alkali-antimonide films for photocathodes on an appropriate substrate (such as glass) takes place in a vacuum chamber by chemical reactions between antimony deposited on a substrate and the vapour of alkali metals. The growth can be studied by measurements of the atomic absorptions in the alkali vapour, the optical reflection and transmission of the film, and the photosensitivity in the reflection and transmission modes. These measurements allow the partial alkali-vapour pressures, the optical properties, thickness and growth rate of the film to be evaluated during the growth, as well as the escape depth and escape probability of the photoelectrons. As has been demonstrated for Na_2KSb films, growth monitoring gives useful information about the preferred chemical reaction mechanism, the optimum film composition and the thickness of the film. It provides a basis for improvement of the growth conditions, so as to increase the photosensitivity and reproducibility. To obtain an Na_2KSb film of high photosensitivity it must be coated with a very thin top layer of caesium. Molecular-beam experiments indicate that the best results are obtained with a monolayer of pure caesium.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- H. A. Algra, L. J. de Jongh** (University of Leiden) & **J. Reedijk** (Delft University of Technology): Specific heat near the critical concentration for the dilute simple-cubic magnet $\text{Co}_p\text{Zn}_{1-p}(\text{C}_5\text{H}_5\text{NO})_6(\text{ClO}_4)_2$. *Phys. Rev. Letters* **42**, 606-609, 1979 (No. 9). *E*
- P. M. Asbeck, D. A. Cammack, J. J. Daniele, D. Lou, J. P. J. Heemskerk, W. J. Kleuters & W. H. Ophrey**: High-density optical recording with (Ga,Al)As DH lasers. *Appl. Phys. Letters* **34**, 835-837, 1979 (No. 12). *N, E*
- C. Belouet**: Vapour growth in a microgravity environment. *Thin Solid Films* **58**, 1-8, 1979 (No. 1). *L*
- F. Berz, R. W. Cooper & S. Fagg**: Recombination in the end regions of pin diodes. *Solid-State Electronics* **22**, 293-301, 1979 (No. 3). *R*
- J. Bloem**: Problems in the melt- and vapor growth of silicon for integrated circuits and solar cells. *J. solid State Chem.* **27**, 19-27, 1979 (No. 1). *E*
- F. R. de Boer, W. H. Dijkman, W. C. M. Mattens** (all with University of Amsterdam) & **A. R. Miedema**: On the valence state of Yb and Ce in transition metal inter-metallic compounds. *J. less-common Met.* **64**, 241-253, 1979 (No. 2). *E*
- M. R. Boudry & J. P. Stagg**: The kinetic behavior of mobile ions in the Al-SiO₂-Si system. *J. appl. Phys.* **50**, 942-950, 1979 (No. 2). *R*
- M. Boulou, M. Furtado, G. Jacob & D. Bois**: Recombination mechanisms in GaN:Zn. *J. Luminescence* **18/19**, 767-770, 1979 (Part II). *L*
- J. C. Brice & A. M. Cole**: Infrared absorption in α -quartz. *J. Physics D* **12**, 459-463, 1979 (No. 3). *R*
- J. W. Broer**: Professional re-identification in the writer/scientist. *Proc. 26th Int. tech. Comm. Conf., Los Angeles 1979*, pp. W 27-31. *E*
- K. H. J. Buschow & N. M. Beekmans**: Magnetic and electrical properties of amorphous alloys of Gd and C, Al, Ga, Ni, Cu, Rh or Pd. Rapid quenched metals III, Proc. 3rd Int. Conf., Brighton 1978, Vol. 2, pp. 133-136. *E*
- K. L. Bye**: An X-ray topographic assessment of cadmium mercury telluride. *J. Mat. Sci.* **14**, 619-625, 1979 (No. 3). *R*
- S. R. Chinn** (M.I.T., Lexington, Mass.) & **W. K. Zwickler**: FM mode-locked Nd_{0.5}La_{0.5}P₅O₁₄ laser. *Appl. Phys. Letters* **34**, 847-849, 1979 (No. 12). *N*
- T. A. C. M. Claasen**: Comments on 'The absolute stability of high-order discrete-time systems utilizing the saturation nonlinearity'. *IEEE Trans. CAS-26*, 138-140, 1979 (No. 2). *E*
- T. A. C. M. Claasen & W. F. G. Mecklenbräuer**: Application of transposition to decimation and interpolation in digital signal processing systems. 1979 IEEE Int. Conf. on Acoustics, speech & signal processing, Washington D.C., pp. 832-835. *E*
- J. A. Clarke**: Aspheric mirror optical systems for the infra-red. 2nd Int. Conf. on Low light and thermal imaging, Nottingham 1979 (IEE Conf. Publ. No. 173), pp. 18-19. *R*
- D. J. Coe & H. E. Brockman**: Corner breakdown in MOS transistors with lightly-doped drains. *Solid-State Electronics* **22**, 444-446, 1979 (No. 4). *R*
- N. H. Dekkers**: Object wave reconstruction in STEM. *Optik* **53**, 131-142, 1979 (No. 2). *E*
- M. Delfino**: A comprehensive optical second harmonic generation study of the non-centrosymmetric character of biological structures. *Mol. Cryst. liq. Cryst.* **52**, 271-284, 1979 (No. 1-4). *N*
- P. Delsarte**: Bilinear forms over a finite field, with applications to coding theory. *J. combin. Theory A* **25**, 226-241, 1978 (No. 3). *B*

- P. Delsarte, Y. Genin & Y. Kamp:** A simple approach to spectral factorization. *IEEE Trans. CAS-25*, 943-946, 1978 (No. 11). *B*
- P. Delsarte, Y. Genin & Y. Kamp:** An equivalence relation in planar least squares inverse approximation. *Proc. IEEE* 66, 1662, 1978 (No. 12). *B*
- P. A. Devijver:** Nonparametric estimation of feature evaluation criteria. *Pattern recognition and signal processing*, ed. C. H. Chen, pp. 61-82; Sijthoff & Noordhoff, Alphen aan den Rijn 1978. *B*
- P. A. Devijver:** On the amount of information conveyed by nearest neighbors and its use in pattern recognition. *Théorie de l'information*, Coll. Int. CNRS No. 276, Paris 1978, pp. 353-362. *B*
- C. Z. van Doorn & J. J. M. J. de Klerk:** Two-frequency 100-line addressing of a reflective twisted-nematic liquid-crystal matrix display. *J. appl. Phys.* 50, 1066-1070, 1979 (No. 2). *E*
- J. W. F. Dorleijn** (Philips Lighting Division, Eindhoven) & **A. R. Miedema:** The anomalous Hall effect in nickel alloys: contributions from skew scattering and side displacement in the two-current model. *J. Magn. magn. Mat.* 12, 26-30, 1979 (No. 1). *E*
- G. Engelsma:** Effect of daylength on phenol metabolism in the leaves of *Salvia occidentalis*. *Plant Physiol.* 63, 765-768, 1979 (No. 4). *E*
- W. van Erk:** A solubility model for rare-earth iron garnets in a $\text{PbO/B}_2\text{O}_3$ solution. *J. Crystal Growth* 46, 539-550, 1979 (No. 4). *E*
- R. M. van Essen & K. H. J. Buschow:** Hydrogen absorption in various zirconium- and hafnium-based intermetallic compounds. *J. less-common Met.* 64, 277-284, 1979 (No. 2). *E*
- L. F. Feiner & R. P. van Staple:** Ontaarde elektronenbanen en structurele faseovergangen. *Ned. T. Natuurk. A* 44, 111-114, 1978 (No. 3). *E*
- W. E. Fischer:** PHIDAS — a database management system for CAD/CAM application software. *Computer-aided Design* 11, 146-150, 1979 (No. 3). *H*
- B. Fitzhenry:** Identification of a charging mechanism using infrared spectroscopy. *Appl. Spectrosc.* 33, 107-110, 1979 (No. 2). *N*
- N. Fleurot, M. Nail, R. Verrecchia** (all with CEA, Villeneuve-Saint-Georges) & **G. Clément:** Characterization of image converter tubes and photodiodes in the infrared region. *Proc. 13th Int. Congress on High speed photography and photonics*, Tokyo 1978, pp. 440-442; 1979. *L*
- J. A. Geurst:** Mutual friction in the laminar flow of superfluid helium II through capillary tubes. *Physics Letters* 71A, 78-82, 1979 (No. 1). *E*
- J. P. Gex, R. Sauneuf** (both with CEA, Villeneuve-Saint-Georges), **J. P. Boutot & J. C. Delmotte:** Some new possibilities in direct visible and X-ray measurements. *Proc. 13th Int. Congress on High speed photography and photonics*, Tokyo 1978, pp. 405-408; 1979. *L*
- R. W. Gibson & R. Wells:** The potential of SSB for land mobile radio. *29th IEEE Vehicular Technology Conf.*, Arlington Heights, Ill., 1979, pp. 90-94. *R*
- A. A. van der Giessen** (Philips Electro-acoustics Division): Audio recording tapes based on iron particles. *J. Audio Engng. Soc.* 26, 838-842, 1978 (No. 11). *R*
- J.-M. Goethals:** Combinatorial decoding methods for block codes. *Théorie de l'information*, Coll. Int. CNRS No. 276, Paris 1978, pp. 223-231. *B*
- W. J. A. Goossens:** Temperature dependence of the pitch in cholesteric liquid crystals: a molecular statistical theory. *J. Physique* 40, C3/158-163, 1979 (Colloque C3). *E*
- R. G. Gossink & T. P. A. Lommen:** Secondary-ion mass spectrometry (SIMS) analysis of electron-bombarded soda-lime-silica glass. *Appl. Phys. Letters* 34, 444-446, 1979 (No. 7). *E*
- L. H. Guildford:** Experiments in real-time image processing. *2nd Int. Conf. on Low light and thermal imaging*, Nottingham 1979 (IEE Conf. Publ. No. 173), pp. 55-56. *R*
- P. Guittard, P. Jarry, C. Piaget, J. C. Richard, E. Roaux & P. Saget:** GaAs photocathodes for low light level imaging. *2nd Int. Conf. on Low light and thermal imaging*, Nottingham 1979 (IEE Conf. Publ. No. 173), pp. 24-25. *L*
- S. Herman:** Nonlinear capacitors improve the performance of saturable lead ballasts. *J. Illum. Engng. Soc.* 8, 122-125, 1979 (No. 3). *N*
- W. J. van den Hoek & J. A. Visser** (Philips Lighting Division, Eindhoven): Arc oscillations in a horizontal rare-earth metal iodide/cesium iodide/mercury arc induced by an external magnetic field. *Appl. Phys. Letters* 34, 357-359, 1979 (No. 6). *R*
- M. H. H. Höfelt:** On the stability of a 1-bit-quantized feedback system. *1979 IEEE Int. Conf. on Acoustics, speech & signal processing*, Washington D.C., pp. 844-848. *E*
- K. Holford:** Front ends for Doppler with-sense radar. *Electronics Letters* 15, 75-76, 1979 (No. 3). *R*
- K. Holford:** Microwaves can control traffic lights for fire appliances. *Fire, March* 1979, pp. 507-508. *R*
- L. Honds & H. Meyer:** Auswirkungen von Oberwellen auf die Materialfunktion im Hysterese-motor. *etz-Archiv* 1, 187-190, 1979 (No. 6). *A*

- H. Ihrig:** Reply to 'Comment on "A systematic experimental and theoretical investigation of the grain-boundary resistivities of *n*-doped BaTiO₃ ceramics"'. *J. appl. Phys.* **50**, 1158-1159, 1979 (No. 2). *A*
- G. D. Khoe:** Practical machine for electric arc splicing of optical fibres in the field. *Electronics Letters* **15**, 152-153, 1979 (No. 5). *E*
- G. D. Khoe, H. G. Kock & L. J. Meuleman:** Fiberless hermetic packaged lens-coupled laser diode for wide-band optical-fiber transmission. *Optical fiber communication, Dig. tech. Papers Topical Meeting Washington D.C. 1979*, pp. 94-97. *E*
- L. C. Kimerling*, P. Blood & W. M. Gibson* (* Bell Laboratories, Murray Hill, N.J.):** Defect states in proton-bombarded silicon at $T < 300$ K. *Int. Conf. on Defects and radiation effects in semiconductors, Nice 1978 (Inst. Phys. Conf. Ser. No. 46)*, pp. 273-280; 1979. *R*
- M. Klinck:** Control of the surface-water purification plant for the Amsterdam Water-Supply Authority. *Journal A* **20**, 59-70, 1979 (No. 2). *H*
- G. Kowalski:** Multislice reconstruction from twin-cone beam scanning. *IEEE Trans. NS-26*, 2895-2903, 1979 (No. 2, Part 2). *H*
- G. Kowalski & W. Wagner:** Patient dose rate: An ultimate limit for spatial and density resolution of scanning systems. *Biomed. Technik* **24**, 38-42, 1979 (No. 3). *H*
- H. K. Kuiken:** The cooling of low-heat-resistance cylinders by radiation. *J. Engng Math.* **13**, 97-106, 1979 (No. 2). *E*
- G. Kurtze, C. Klingshirn** (both with Universität Karlsruhe), **B. Hönerlage** (Laboratoire de Spectroscopie et d'Optique du Corps Solide, Strasbourg), **E. Tomzig** (Universität Erlangen) & **H. Scholz:** Excitation spectroscopy, Raman scattering and the temperature dependence of the luminescence in highly excited red HgI₂. *J. Luminescence* **20**, 151-161, 1979 (No. 2). *A*
- J. van Laar:** Foto-emissie van vaste stoffen. *Ned. T. Natuurk. A* **44**, 91-95, 1978 (No. 3). *E*
- G. Laurence, F. Simonet & P. Saget:** Combined RHEED-AES study of the thermal treatment of (001) GaAs surface prior to MBE growth. *Appl. Phys.* **19**, 63-70, 1979 (No. 1). *L*
- F. H. de Leeuw, W. de Geus & P. Q. J. Nederpel:** High speed photography of moving domain walls in magnetic bubble materials. *Proc. 13th Int. Congress on High speed photography and photonics, Tokyo 1978*, pp. 317-320; 1979. *E*
- P. A. Lewis & M. J. Underhill:** Quiet tuning and matching of antennas for radio silence operation. *Proc. IEE Conf. on Recent advances in h.f. communication systems and techniques, London 1979 (IEE Colloq. Dig. 1979/48)*, pp. 31-44. *R*
- G. M. Martin, M. L. Verheijke, J. A. J. Jansen & G. Poiblaud** (RTC, Caen): Measurement of the chromium concentration in semi-insulating GaAs using optical absorption. *J. appl. Phys.* **50**, 467-471, 1979 (No. 1). *L, E*
- R. Memming & F. Schröppel:** Electron transfer reactions of excited ruthenium(II) complexes in monolayer assemblies at the SnO₂-water interface. *Chem. Phys. Letters* **62**, 207-210, 1979 (No. 2). *H*
- A. D. Mills, I. Mackenzie & R. J. Dolphin:** The use of a microcomputer for flexible automation of a liquid chromatograph. *J. autom. Chem.* **1**, 134-140, 1979 (No. 3). *R*
- A. Mircea & D. Bois** (INSA, Villeurbanne): A review of deep-level defects in III-V semiconductors. *Int. Conf. on Defects and radiation effects in semiconductors, Nice 1978 (Inst. Phys. Conf. Ser. No. 46)*, pp. 82-99; 1979. *L*
- B. J. Mulder:** Unbacked ultra-thin films of beryllium and other metals. *J. Physics E* **12**, 268-269, 1979 (No. 4). *E*
- A. van Oostrom, L. Augustus, A. Steinmetz* & G. v.d. Berg*** (* Philips Telecommunications Industries, Hilversum): Analysis of surface resistance and elemental composition of a reed contact. *Electrical contacts 1978, Proc. 9th Int. Conf. & 24th Ann. Holm Conf., Chicago*, pp. 521-526. *E*
- R. Orłowski & E. Krätzig:** Holographische Speicherung in elektrooptischen Kristallen. *Laserspektroskopie*, ed. F. Aussenegg (*Acta Physica Austriaca, Suppl. XX*), pp. 241-255; Springer, Wien 1979. *H*
- J. A. Pals & J. Wolter:** Measurement of the order-parameter relaxation in superconducting Al-strips. *Physics Letters* **70A**, 150-152, 1979 (No. 2). *E*
- J. B. H. Peek & J. M. Schmidt:** Een 'Station Programma Identificatie' systeem voor F.M.-radioomroep. *T. Ned. Elektronica- en Radiogen.* **44**, 25-29, 1979 (No. 1). *E*
- P. van Pelt & E. E. Havinga:** Electrochromism of organic dyes at high fields. *Nonlinear behaviour of molecules, atoms and ions in electric, magnetic or electromagnetic fields*, ed. L. Néel, pp. 291-300; Elsevier, Amsterdam 1979. *E*
- D. Polder, M. F. H. Schuurmans & Q. H. F. Vrehen:** Superfluorescence: Quantum-mechanical derivation of Maxwell-Bloch description with fluctuating field source. *Phys. Rev. A* **19**, 1192-1203, 1979 (No. 3). *E*
- D. Pons, A. Mircea, A. Mitonneau & G. M. Martin:** Electron traps in irradiated GaAs: comparison with native defects. *Int. Conf. on Defects and radiation effects in semiconductors, Nice 1978 (Inst. Phys. Conf. Ser. No. 46)*, pp. 352-359; 1979. *L*

- G. Prast:** A thermal tracking solar collector. Proc. DFVLR Int. Symp. on Solar thermal power stations, Cologne 1978, section 10, pp. 1-8. *E*
- J. L. Robert, B. Pistoulet, A. Raymond, J. M. Dussseau** (all with Centre d'Etudes d'Electronique des Solides, Montpellier) & **G. M. Martin:** New model of conduction mechanism in semi-insulating GaAs. *J. appl. Phys.* **50**, 349-351, 1979 (No. 1). *L*
- P. C. Scholten:** Colloid chemistry of magnetic fluids. Thermomechanics of magnetic fluids, ed. B. Berkovsky, pp. 1-25; Hemisphere, Washington D.C. 1978. *E*
- H. Schomberg:** Reconstruction of spatial resistivity distribution of conducting objects from external resistance measurements. *Z. angew. Math. Mech.* **59**, T 41-42, 1979 (No. 3). *H*
- H. Schomberg:** Monotonically convergent iterative methods for nonlinear systems of equations. *Numer. Math.* **32**, 97-104, 1979 (No. 1). *H*
- P. J. Severin:** Calorimetric measurement of the absorption coefficient of fibre-quality compound glass rods. *Óptica hoy y mañana, Proc. ICO-11, Madrid 1978*, pp. 499-502. *E*
- P. J. Severin:** Isotope separation by chemical vapour deposition and related processes. *J. Crystal Growth* **46**, 630-636, 1979 (No. 5). *E*
- M. Sintzoff:** Properties, feasibility and usefulness of a language for programming programs. *ALGOL Bull. No. 43*, 84-90, 1978. *B*
- M. J. Sparnaay:** Transfert d'électrons et d'ions dans et sur les semi-conducteurs. *Le Vide* **33**, 236-237, 1978 (No. 194). *E*
- W. Spiesberger:** Mammogram inspection by computer. *IEEE Trans. BME-26*, 213-219, 1979 (No. 4). *H*
- B. Steinmüller & R. Bruno:** The energy requirements of buildings. *Energy and Buildings* **2**, 225-235, 1979 (No. 3). *A*
- A. L. N. Stevels:** Red Mn^{2+} -luminescence in hexagonal aluminates. *J. Luminescence* **20**, 99-109, 1979 (No. 2). *E*
- F. L. H. M. Stumpers:** Some notes on the scientific work of H. Bremmer. *Radio Sci.* **14**, 179-188, 1979 (No. 2). *E*
- J. B. Theeten & F. Hottier:** *In situ* surface analysis of the vapor phase epitaxy of GaAs. *J. Electrochem. Soc.* **126**, 450-460, 1979 (No. 3). *L*
- J. B. Theeten, F. Hottier & J. Hallais:** Ellipsometric assessment of (Ga,Al)As/GaAs epitaxial layers during their growth in an organometallic VPE system. *J. Crystal Growth* **46**, 245-252, 1979 (No. 2). *L*
- J. J. P. Valetton** (Philips Electro-acoustics Division): Electronische signaalbehandeling in televisiecamera's. *T. Ned. Elektronica- en Radiogen.* **44**, 9-16, 1979 (No. 1). *L*
- H. Vantilborgh:** Exact aggregation in exponential queueing networks. *J. Ass. Computing Machinery* **25**, 620-629, 1978 (No. 4). *B*
- H. Verweij:** Raman study of the structure of alkali germanosilicate glasses: I. Sodium and potassium metagermanosilicate glasses, II. Lithium, sodium and potassium digermanosilicate glasses. *J. non-cryst. Solids* **33**, 41-53, 55-69, 1979 (No. 1). *E*
- H. Verweij, J. H. J. M. Buster & G. F. Remmers** (Twente University of Technology, Enschede): Refractive index and density of Li-, Na- and K-germanosilicate glasses. *J. Mat. Sci.* **14**, 931-940, 1979 (No. 4). *E*
- A. T. Vink, C. J. Werkhoven & C. van Opdorp:** Large defects: observation and influence on minority carrier recombination. *Semiconductor characterization techniques, Electrochem. Soc. Proc. Vol. 78-3*, pp. 259-288, 1978. *E*
- W. Wagner:** Reconstructions from restricted region scan data — new means to reduce the patient dose. *IEEE Trans. NS-26*, 2866-2869, 1979 (No. 2, Part 2). *H*
- H. W. Werner & A. E. Morgan:** Secondaire-ionen-massaspectrometrie. *Ned. T. Natuurk. A* **44**, 122-125, 1978 (No. 3). *E*
- J. S. van Wieringen:** Exact and approximate solution of the regenerator equation for the case of high heat exchange and moderate heat capacity. *Appl. sci. Res.* **34**, 145-158, 1978 (No. 2/3). *E*
- J. P. Woerdman:** Laser-excited broadband violet emission from sodium vapor. *J. Opt. Soc. Amer.* **68**, 714, 1978 (No. 5). *E*
- M. de Zwart & C. Z. van Doorn:** The field-induced square grid perturbation in the planar texture of cholesteric liquid crystals. *J. Physique* **40**, C3/278-284, 1979 (Colloque C3). *E*
- Published in Gallium arsenide and related compounds, 1978 (Proc. 7th Int. Symp., St. Louis; Inst. Phys. Conf. Ser. No. 45, 1979):*
- C. E. C. Wood, J. Woodcock & J. J. Harris:** Low-compensation n-type and flat-surface p-type Ge-doped GaAs by molecular beam epitaxy (pp. 28-37). *R*
- G. B. Scott & J. S. Roberts:** Photoluminescence in III-V compounds grown by MBE (pp. 181-189). *R*
- W. J. Bartels & H. Veenliet:** X-ray study of $Ga_{1-x}Al_xAs$ epitaxial layers grown with the metallorganic VPE technique on GaAs substrates (pp. 229-238). *E*
- Ch. Hurtes, L. Hollan & M. Boulou:** Impurity characterization of GaAs high-resistivity VPE layers for FET devices (pp. 342-352). *L*
- J. Hallais, J. P. André, P. Baudet & D. Boccon-Gibod:** New MESFET devices based on GaAs-(Ga, Al)As heterostructures grown by metallorganic VPE (pp. 361-370). *L*

An experimental assembly robot

J. G. van den Hanenberg and J. Vredenburg

When Karel Čapek coined the word 'robot' (from the Czech robota, meaning heavy work) in 1920 to denote a machine in the form of a man, he can hardly have dreamed that it would acquire such great significance in technology some sixty years later. In Čapek's play the robots become a threat to mankind. Robots in our day and age, however, are meant to release man from heavy or tedious work, or to perform work in inaccessible places or under conditions in which human beings cannot function effectively. The robot described here has been developed at Philips Research Laboratories to gain experience with a new aid that can dramatically change the nature of industrial mechanization.

Introduction

The traditional Philips products, such as incandescent lamps and electronic components, can only be mass-produced at sufficiently low prices and with satisfactory quality if specially designed machines are used. Philips have therefore built up extensive experience in industrial mechanization.

A common feature to many of the machines that have been developed for mass production in almost every industry is that they were designed for performing one or more specific tasks: making components, setting processes in motion, assembling components and carrying out monitoring and measuring operations, possibly in combination. These operations are built into the machine permanently, for example as movement patterns determined by cams. The processes or measurements to be carried out are controlled by electronic circuits or — less rigidly — by specially written computer programs. If the product has to be modified or redesigned, the production machines have to be modified or replaced. This type of mechanization is therefore called 'hard automation'.

Not all products are suited to this method of manufacture. An important factor is the size of the production run, of course, but much also depends on the com-

plexity of the operations required for assembling the products. Some relatively complicated products, such as electric shavers and vacuum cleaners, still require a good deal of manual assembly in production.

Nowadays compact and versatile machine control systems can be made by using computers. When the production process is modified it is a simple matter to reprogram the control computer of the machine. However, adjusting the hardware from the control system is still something of a problem. In numerically controlled machine tools this problem has been more or less satisfactorily solved. A recent development in assembly operations that varies is the use of 'industrial robots'. The introduction of these robots in the actual manufacturing process is referred to as 'flexible automation' [1][2].

An industrial 'robot' is not like the humanoids of science fiction, with arms and legs. It is usually a mechanical system in a fixed location, which consists of an 'arm' — controlled by an electronic 'brain' — that can independently perform a number of operations within a confined space. The arm consists of a number of movable interconnected parts and has a

[1] J. L. Nevins and D. E. Whitney, Computer-controlled assembly, *Sci. Amer.* **238**, Feb. 1978, p. 62-74.

[2] P. Saraga and J. A. Weaver, *Philips tech. Rev.* **38**, 329, 1978/79.

'hand', called a gripper. The arm is an all-purpose mechanism, whereas the gripper may have to be adapted to the objects to be picked up. Since the commands to the robot can easily be changed, it is able to perform a variety of short or long tasks. It can carry out short-cycle operations for a long time without becoming tired or it can work in adverse conditions. Since a robot is an all-purpose tool, it can remain a factory fixture, even when the production programme is subject to changes. The robot's job assignments can be changed by reprogramming the control computer, or by manually 'demonstrating' to the robot a new pattern of movements.

The robots in current use and in development are divided into three generations. Those now on the market and in operation belong to the first generation. The arm is controlled so that it moves from point to point or along a particular path in a coordinate system that is fixed in relation to the environment. The result of each movement is not measured or verified and cannot therefore be fed back into the control system. However, there is some uncertainty in the position of the gripper — of the order of 0.1 mm or more — and this is mainly due to friction, play and temperature effects in the arm mechanism. When the arm is given more degrees of freedom its range of applications is increased but at the same time the influence of these undesirable factors increases. Assembly with first-generation robots may therefore be unsuccessful if the positioning error is larger than the clearance between the parts to be assembled.

Because they have no 'senses', robots of this first generation do not respond to changes in their environment. This means that the components for assembly must always be presented at exactly the same place and with the same orientation, and they must also have the same shape and dimensions. There must be sufficient clearance between the objects to be assembled, and the path described by the arm must be free of foreign objects. Because of these constraints, the first-generation robots can only be used for relatively inaccurate operations, such as removing products from injection-moulding machines and performing process movements that do not require great accuracy, as in spot-welding and paint spraying in the automobile industry.

Deviations in the position of the arm mechanism and in the component feed can only be corrected by enabling the robot to observe its environment via 'sensors', that can 'touch', 'feel' and 'see'. Robots thus equipped belong to the second generation. They are less dependent on the inaccuracy of the mechanical positioning. The system has become an 'adaptive' one through feedback of the information obtained

from the sensors and its immediate processing in the control system in 'real time'; in other words, the system parameters are adapted to changing conditions. As a consequence the positioning and the geometry of the components do not have to be highly accurate. Without any improvement in the measurement of the absolute position, the robot can perform more accurate work, such as the assembly of components with very small clearances.

So as to identify the fundamental problems in the development and application of robots, and at the same time to facilitate the application of robots in the company, work was started at Philips Research Laboratories some years ago on designing and making an experimental industrial robot. It was called EPAAS, for 'Experimental Programmable Adaptive Assembly System'. This robot, which will be described here, is capable, for example, of assembling components with clearances of no more than 10 μm between them. Features of particular interest in the robot include a gripper whose fingers have compliance and sensors for force, a miniature camera in the gripper, movement with controlled force in difficult assembly operations, and the straightforward manner in which the user can compile his application program. Extending the system to include, for example, the pattern-recognition method described earlier in this journal ^[3] and introducing certain forms of self-organization ^[4] will make the robot even more versatile in the near future. It will then have grown into a fully-fledged robot of the third generation, capable of processing previous experience in its control system, in other words equipped with artificial intelligence.

We shall now deal first with the mechanical configuration of the robot arm (we sometimes call it an 'anthropomorphic' arm) and the gripper, and then we shall consider the sensors incorporated in the robot. Next, we shall look at the control system and its software, showing how the information from the sensors is used. Finally, we shall give an example of an assembly program for inserting a pin or 'peg' in a hole.

The robot arm

The arm has the function of conveying the gripper from place to place in the working space, and in a specified direction for each position. There are very many possible variations of the arm mechanism, with the kinematic joints interconnected by means of pivoting or telescopic joints. There are technical disadvantages with telescopic joints, however, such as a considerable amount of play and friction, and in many cases the coordinate transformation is difficult ^[5] — this will be mentioned later. There are also

variations in the mass-to-stiffness ratio when telescopic joints are used, causing marked differences in the dynamic characteristics of the mechanism in its various positions. Mainly for these reasons we decided to use an arm with pivoting joints only, which

determining the position of the point W of the wrist joint in the working space. The coordinate system (x, y, z) is fixed in relation to the environment. The control is designed in such a way that the points through which the wrist joint has to move are calculated in

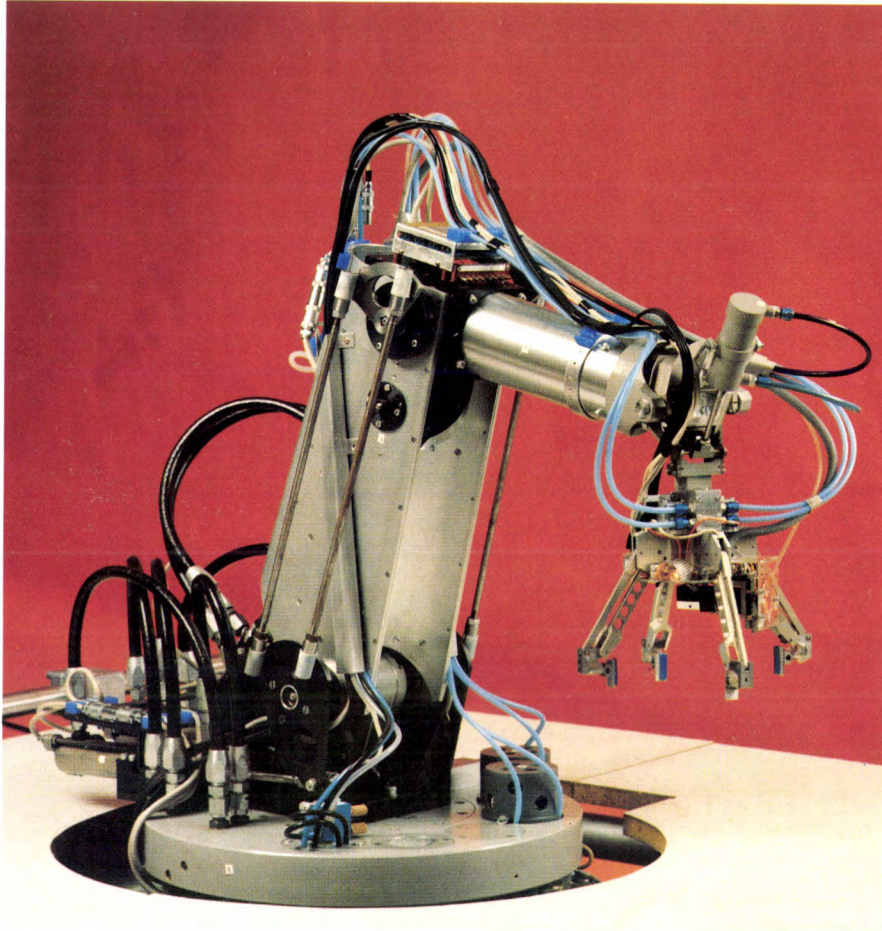


Fig. 1. Photograph of the arm mechanism of EPAAS, the experimental assembly robot developed at Philips Research Laboratories. To combine great stiffness for bending and torsion with low mass, the arm sections are made from thin-walled aluminium tubes.

resembles a human arm in its kinematic design; see *fig. 1*. The arm has a 'shoulder joint', an 'elbow joint' and a 'wrist joint'. The shoulder has two degrees of freedom, the elbow one and the wrist two. The hand attached to the wrist, the gripper, can move in a number of ways, which will be described in the next section.

Fig. 2 is a diagram showing the axes of rotation that form the arm's five degrees of freedom. The axes of rotation ZZ and AA are those of the shoulder joint, BB is the axis of rotation of the elbow joint, while CC and DD (corresponding to angles ψ and γ) are the axes of rotation of the wrist joint. The figure also indicates the coordinate systems that are used for

cylindrical coordinates (r, z, ϕ) . The z -axes of the two systems coincide and the r -axis corresponds to the x -axis if $\phi = 0$. The arm mechanism thus lies in an (r, z) -plane and the axes of rotation AA , BB and DD are perpendicular to this plane. The rotation about the axis ZZ corresponds to the ϕ -coordinate. The linear displacements l_3 and l_4 (see *fig. 3b* and *c*) of the pistons in the drive units that produce the rotations about the axes AA and BB have to be determined by the control system, however, from the coordinates r and z . The same applies to the displacements l_1 and l_2

[3] E. H. J. Persoon, Philips tech. Rev. 38, 356, 1978/79.

[4] T. J. B. Swanenburg, Philips tech. Rev. 38, 364, 1978/79.

[5] See H. Rankers, Metaalbewerking 45, 451, 1979 (in Dutch).

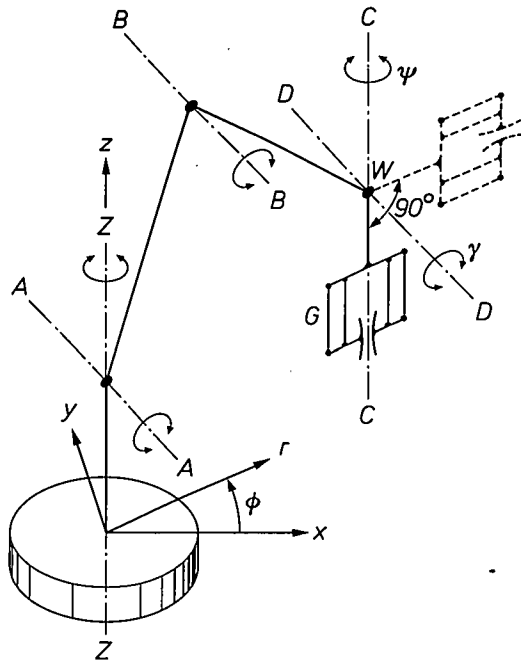


Fig. 2. Diagram of the arm mechanism. ZZ, AA, BB, CC and DD are the axes of rotation of the mechanism, which together give the five degrees of freedom. By rotation about the axes ZZ, AA and BB, the point W, the 'wrist joint', can be moved to any point in the working space. The axes AA and ZZ together form the 'shoulder joint', the axis BB can be considered as the 'elbow joint'. The direction of the 'hand' or gripper G can be varied by rotation about the axes CC and DD (with associated angles ψ and γ). The axis CC is fixed in relation to the axis DD, and intersects it at the point W. A rotation about DD thus also changes the direction of CC. An additional rotation through 90° — effected by means of a separate drive — alters the mean position of the gripper from vertical to horizontal. The dashed lines indicate the mean horizontal position. The coordinate system (x, y, z) is fixed in relation to the environment. The control system is designed in such a way that the points through which the point W has to move are calculated in cylindrical coordinates (r, z, ϕ) . The z-axes of both systems coincide; the r-axis corresponds to the x-axis when $\phi = 0$.

(see fig. 3d and e) of the pistons in the drive units for the rotations about the axes DD and CC (the angles γ and ψ). These coordinate transformations must be carried out 'in real time'. To make the calculations as simple as possible the arm mechanism is therefore designed in such a way that a rotation about the axis AA does not change the angle of the second arm-element in relation to the (r, z) -system, i.e. there is no change in the rotation about the axis BB. The angles γ and ψ also remain unchanged when rotations take place about the AA or BB axes.

The diagrams a to e in fig. 3 illustrate how the drives for the five rotations are produced. The axis ZZ corresponds to the axis of the turntable on which the actual arm mechanism is mounted. The rotation l_5 about the axis ZZ is measured with a rotary precision potentiometer. The components that have a relatively large mass, such as the linear servomotors for producing the other four rotations, are mounted close to this axis to minimize the moment of inertia about ZZ. These linear servomotors and the vane-rotor type servomotor for ZZ are driven hydraulically. Since the oil pressure can be raised to a relatively high value — in our case to 7×10^6 Pa (70 bars) — the hydraulic servomotors are much more compact and lighter for the same power than pneumatic or electrical devices. The frequencies of the mechanical resonances of the

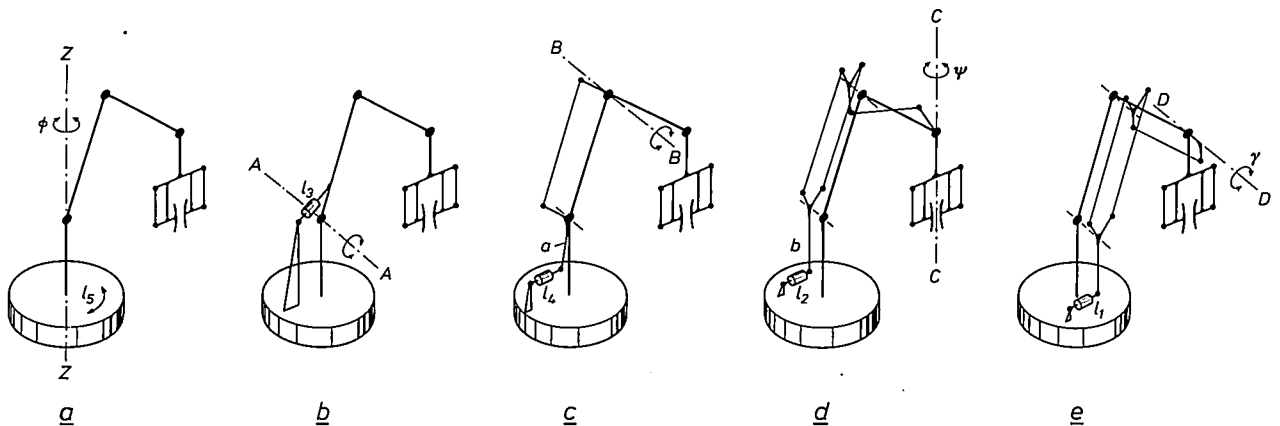


Fig. 3. a) The axis of rotation ZZ. This axis corresponds to the axis of rotation of the turntable on which the arm mechanism is mounted. The turntable is driven by a hydraulic servomotor with vane rotor. The angle l_5 through which the rotation takes place is measured with a rotary precision potentiometer and corresponds to the coordinate ϕ . b) The axis of rotation AA. A rotation about this axis is obtained by moving the first arm element by means of a hydraulic servomotor with integrated piston-position measurement, known as a bar-type actuator (BTA, see fig. 4). l_3 indicates the position of the piston in the cylinder. c) The axis of rotation BB. Rotation about this axis — like that about the axes CC and DD — is effected with the aid of a BTA mounted on the turntable and a system of rods in a parallelogram. This ensures that the angle

between the second arm element and the horizontal plane does not change when there is a rotation about AA. l_4 indicates the position of the piston in the cylinder. d) The axis of rotation CC. The rotation about CC is not affected by a rotation about AA. Independence of the rotation about BB is achieved by electronic coupling between the control circuits for both axes: when the servomotor for BB is actuated, the servomotor for CC is given a similar compensatory displacement, since the lengths of the levers a and b in fig. 3c and 3d are equal. l_2 indicates the position of the piston in the cylinder. e) The axis of rotation DD. The drive with the system of rods illustrated ensures that the position of the gripper is not changed by rotations about AA or BB. l_1 indicates the position of the piston in the cylinder.

two sections of the arm are high (about 90 Hz and higher) and are well above those of the control systems for the hydraulic servomotors. This is achieved by designing the arm sections in the form of cylinders made from thin-walled aluminium tubing.

of about 10 kHz. As soon as the arm approaches an object, the capacitance of the capacitor changes, so that the frequency and voltage of the oscillator change as well. The detected voltage change is used to stop the arm movement in time.

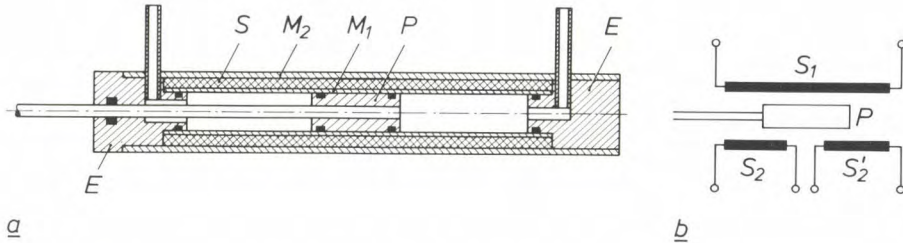


Fig. 4. a) The ‘bar-type actuator’ BTA [6]. The hydraulic drive and the position measurement for the piston are integrated in this design. M_1 and M_2 cylinder walls. The inner cylinder wall M_1 , made of non-magnetic material, forms one of the walls of the oil-filled space; the outer cylinder wall M_2 , made of steel, forms part of the iron circuit for the set of coils S . E cylinder ends, also forming part of the iron circuit. P steel piston. b) The coils from which S is built up. The primary coil S_1 is supplied with alternating current, which induces alternating voltages in the two secondary coils S_2 and S_2' . The difference between these two secondary voltages is a measure of the position of the piston; in the mean position the difference is zero.

The linear hydraulic servomotors are referred to as BTAs (bar-type actuators) [6], in which the hydraulic drive and the system for measuring the position of the piston are integrated; see fig. 4a. Between the two cylinder walls M_1 and M_2 there are three coils that together form a differential transformer; see fig. 4b. The primary coil S_1 is supplied with an alternating current. The position of the steel piston determines the magnitude of the alternating voltages induced in the secondary coils S_2 and S_2' . When the piston is in its central position they are equal; as the piston moves towards one of its extreme positions the voltage difference increases, so that this difference is a measure of the piston displacement.

Fig. 5 shows the angles that the arm sections can cover. The hatched area, which is bounded by circular arcs, is the area that the point W of the wrist joint in the (r, z) -plane can reach as determined by geometrical considerations. To give the r - and z -coordinates a constant limiting value, the useful region is confined to the grey shaded area. The control software prevents the gripper from colliding with the second arm-element. The angular displacement is $\pm 160^\circ$ for rotation about the axis ZZ and $\pm 45^\circ$ for rotation about the axes CC and DD .

To ensure that the arm does not make an undesired contact with an object, insulated metal strips are attached to both arm sections. These strips form a capacitor with the metal table-top. This capacitor in turn forms part of an RC oscillator with a frequency

The gripper and the force sensors

Since a robot should be capable of performing a wide variety of tasks, the gripper should be designed as an all-purpose instrument for handling objects of diverse shapes, mass and ruggedness. Our gripper (see fig. 6) has four ‘fingers’; diametrically opposite pairs can be independently controlled. Since most of the

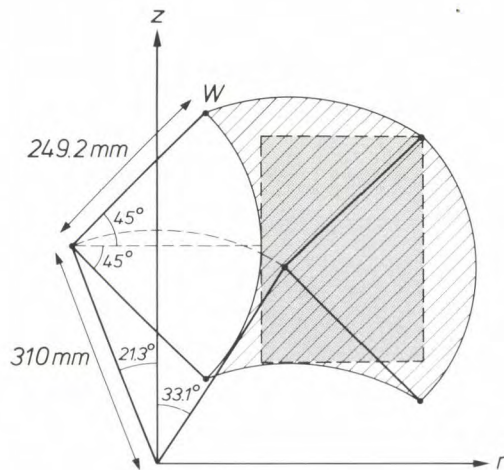


Fig. 5. The area in the (r, z) -plane that the arm can cover by rotation about the axes AA and BB . The diagram shows the maximum angles of the arm sections in relation to the coordinate system and the lengths of the arm sections. The hatched area can theoretically be reached by the point W of the wrist joint. To give the coordinates r and z a constant limit value, the useful area is limited to the region shaded grey.

[6] The BTA was designed by the Philips Centre for Technology.

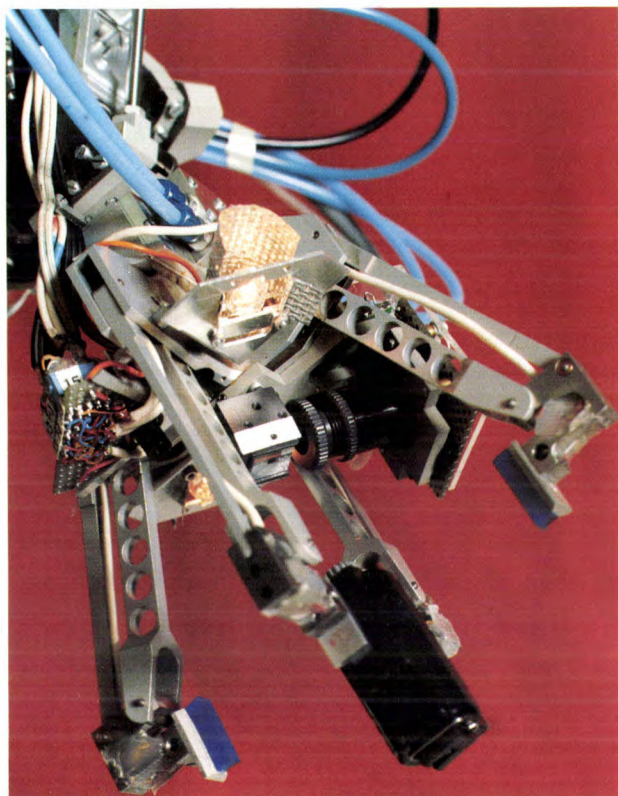


Fig. 6. Photograph of the gripper mechanism. The gripper consists of two pairs of 'fingers', and each pair can be separately controlled. When the control for one pair of fingers is fixed at a certain value, the gripper can be used as a two-finger type. The semiconductor camera mounted in the gripper (see fig. 11a) can be seen just to the right of centre.

objects to be picked up are round, square or rectangular a gripper with four fingers has a wider range of application than the frequently used types with three or two fingers. Furthermore, since the fingers are controlled in pairs, the gripping forces of a four-finger type are more exactly defined. It can also be used as a two-finger gripper if the control for one pair of fingers is rendered inoperative. *Fig. 7* shows the mechanism of one pair of diametrically positioned fingers. The two fingers are centred in relation to the longitudinal axis of the gripper. The pairs of fingers are operated by compressed-air cylinders Cy , which are designed as an integrated part of the gripper to minimize the mass of the gripper system. The gripping force between finger and object can be linearly varied for each finger pair in 256 steps between -45 and $+45$ N. This is done by means of a control system in which the pressure in the compressed-air cylinders is accurately measured by medical cerebral-pressure indicators (Honeywell & Philips, type 9822810 30001). The fingers can close around objects with dimensions between 2 and 110 mm, ranging from heavy (4 kg) to light and from rugged to fragile.

Two strain gauges G ^[7] are fitted to each finger as close as possible to the point where finger and object

touch. The strain gauges for each finger are connected into a Wheatstone bridge in such a way that only the torque acting in a cross-section of the finger is measured, and not the gripping force along the direction of its length. This enables sensory information to be obtained about the forces exerted on the object: the vertical force F_z and the horizontal forces F_r (for one pair) and F_ϕ (for the other pair); see *fig. 8a* and *b*. Adding the output signals from the Wheatstone bridges gives a signal that is a measure of F_z , and taking the difference of these signals results in a signal that represents F_r or F_ϕ . The bandwidth of the complete force-measuring system is about 100 Hz; the resolution of the F_z measurement is 0.01 N.

When an object is gripped by two diametrically opposed fingers, one of the fingers will touch it first, owing to the unavoidable inaccuracy in positioning. This causes the object to tilt until it comes into contact with the second finger. After the fingers have pushed the object back to approximately its original position, the object exerts forces on the finger sections, and the opposing frictional components of these forces set up torques, which are measured by the strain gauges. This takes place without any external forces acting on the object, other than the force of gravity. Since the force signal is used as a reference

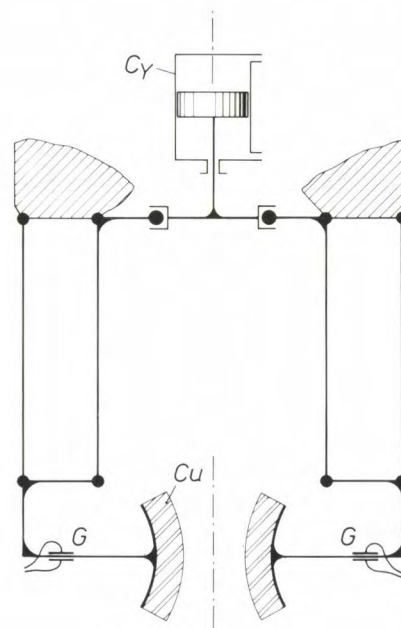


Fig. 7. Diagram of the mechanism for one pair of fingers. Each pair is operated by its own pneumatic cylinder Cy , which is integrated with the gripper. Each finger is fitted with a 'cushion' Cu of a rubbery material (the function of this is illustrated in *fig. 9b*). The parts of the fingers closest to the point of contact with the object are fitted with strain gauges G . These provide sensory information about the forces acting on the object (see *fig. 8a* and *b*). The strain gauges are connected into a Wheatstone bridge in such a way that only the torque acting in the cross-section of the finger where the strain gauges are fitted is measured.

during assembly, these initial torques (whose magnitude and direction are unknown) may cause difficulties. The force signal is therefore reduced to zero by subtracting an appropriate amount from it in the control electronics just before assembly starts.

Fig. 9a illustrates the situation when a force is exerted on the object in the gripper during a sideways displacement and the maximum frictional force between finger and object is reached. If the gripper is displaced further, the object shifts between the fingers. It becomes more and more tilted, while the output signal from the force sensors remains virtually constant; the signal is 'saturated'. When the gripper is moved back again, the position of the object in relation to the fingers becomes different but the value of the sensor signal is the same. Since the sensor signal can then no longer be used for controlling the gripper, this situation must be avoided. We therefore gave the gripper some 'compliance' by covering the fingers with an elastic, rubbery material (*Cu* in fig. 7) that possesses a high coefficient of friction. In fig. 9b the torque detected by the strain gauges is plotted against the displacement *x* of the gripper, for both 'hard' (1) and 'soft' (2) fingers. It can clearly be seen that the maximum permissible lateral displacement of the gripper before the output signal 'saturates' is much larger with soft fingers. This is partly due to the greater flexibility — smaller Young's modulus — and partly to the higher coefficient of friction of the material on the fingers.

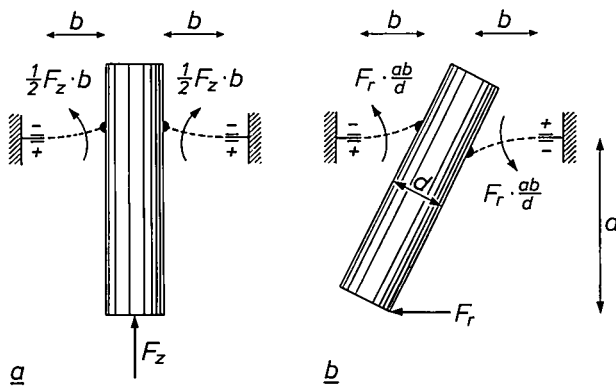


Fig. 8. a) The torques acting in the cross-sections of the fingers at the location of the strain gauges when a vertical force F_z is exerted on the base of the object. The dimension b is the distance between the object and the strain gauges. The lower strain gauges are lengthened (+), the upper ones shortened (-). Since the fingers are covered with soft material (*Cu* in fig. 7), the contact surfaces of finger and object may be regarded as pivot joints each transmitting a force of magnitude $\frac{1}{2}F_z$. Adding the signals from the two pairs of strain gauges yields a signal that is a measure of F_z . b) The torques acting in the cross-sections of the fingers at the location of the strain gauges when a horizontal force F_r is exerted on the base of the object. The contact surfaces transmit vertical forces $F_r \cdot a/d$, which have opposite directions. The tilt of the object is exaggerated (as in fig. 9a) for clarity. When the difference of the signals from the two pairs of strain gauges is taken, the resulting signal is a measure of F_r (or F_ϕ for the other pair of fingers).

We shall now show how feedback of the sensory information from the fingers to the servo-control system enables the robot to place a peg in a hole independently. The hole does not have to be chamfered, and the clearance between the peg and the hole can be much smaller than the error in the position of the gripper. The various positioning stages are shown in fig. 10a. The first stage is reached after coarse positioning with the aid of the position transducers in the arm drive. The shaft is tilted in relation to the hole, thus apparently producing a greater clearance. Next the control circuit — to be dealt with later — is switched on, and this processes the information received from the force sensors in the fingers. Once the peg has reached the required depth in the hole, the movement stops and the operation of placing the peg in the hole has been completed. Fig. 10b shows the measured forces F_H and F_V as a function of time.

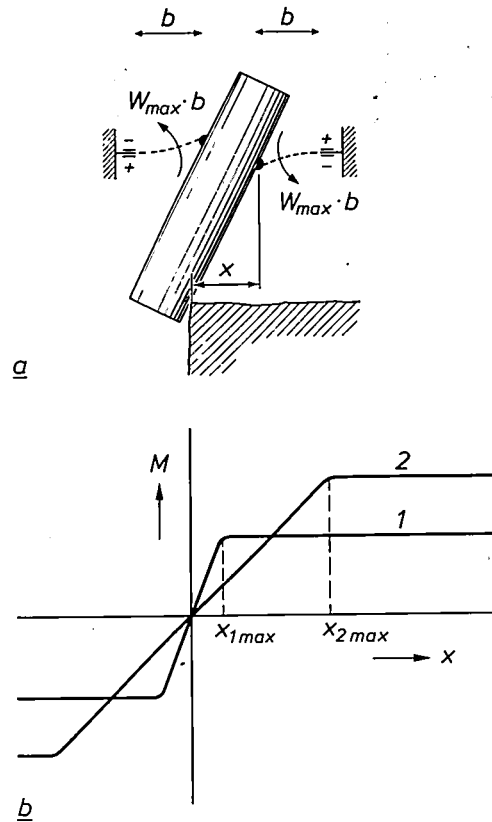


Fig. 9. a) Situation of fig. 8b when the maximum frictional force W_{max} between finger and object is reached. The object begins to 'slip through the fingers'. The torque of magnitude $W_{max} \cdot b$ detected by the strain gauges no longer increases when the gripper is moved further sideways. b) The torque M detected by the strain gauges, as a function of the lateral displacement x of the gripper after the object in the gripper has collided with some other object. Curve 2 applies to the case where the fingers are covered with soft rubbery material (*Cu* in fig. 7); curve 1 relates to hard uncushioned fingers. It can clearly be seen that the lateral displacement that can take place before the signal becomes saturated is much greater for curve 2 than for curve 1 ($x_{2max} > x_{1max}$). This is attributable to the low Young's modulus and the high coefficient of friction of the material covering the fingertips.

[7] See K. Bethe and D. Schön, Philips tech. Rev. 39, 94, 1980.

The visual sensors

The visual sensors serve to recognize and localize the assembly positions and the objects to be picked up, so that the requirements for the absolute accuracy of the arm mechanism and the component feed are less exacting. The visual sensors most commonly used are cameras, mounted in a fixed relation to their environment and equipped with a camera tube. The disadvantage of cameras is that the picture tends to drift in relation to the fixed coordinate system (x, y, z), partly as a result of mechanical and thermal effects. What is more, the camera image is not a true representation of the environment because nonlinear distortion is introduced into the image. In robot systems described elsewhere these effects have been minimized by means of calibration procedures, but these are a nuisance because they take a relatively long time. To get around these problems we use a camera that

moves with the arm, as well as conventional fixed TV cameras. The moving camera is incorporated in the upper part of the gripper; see fig. 6.

The camera in the gripper should be small and light (see fig. 11a). Instead of a camera tube, with its heavy deflection coils, we therefore use a semiconductor detector. This consists of a matrix of 100×100 light-sensitive diodes with associated charge-coupled devices (CCDs) for transmitting the charges in the capacitors formed by the diodes. A simple achromatic lens is also used; this can be set in one of two focusing positions by the control system. The electronic circuits — except for the amplifier stage — are not contained in the arm but in the control unit. The camera always 'looks' in the direction of the axis of symmetry of the gripper and gives an accurate image of the region of the environment that is of interest for the robot. The information obtained from this picture enables the system to adapt itself to the changing conditions in the working space.

The gripper camera is shown diagrammatically in fig. 11b. The distance from the object O to the lens Le varies very little during use, so that with the small aperture angle, the availability of only two lens positions still gives a sufficiently sharp focus. The object is illuminated by two lamps La with reflectors M (one set is shown in the figure). The amount of light reflected is measured with a photodiode D , which has approximately the same spectral sensitivity as the detector S . The lamp current is adjusted by means of a control circuit to ensure a sufficiently constant illumination of the detector. The optical axis between the object and prism P coincides approximately with the axis of symmetry of the gripper; the prism deflects the light beam through an angle of 90° . The incident illumination of the object adopted in our case is preferable to the transmitted illumination (back-lighting) often used, in spite of the addition of moving mass, because it also permits observation of details such as blind holes.

The processing of the picture observed by the semiconductor camera is at present carried out with a separate fast microprocessor, the Signetics 8X300. Owing to the limited frame frequency of the cameras (25 Hz for the fixed cameras and 100 Hz maximum for the semiconductor camera), it is not possible to recognize and determine the centre of gravity and angular orientation of objects in every 6 ms sampling period of the system (the sampling will be discussed later). The image processing is therefore generally done at times determined by the control program. Collecting and processing the image takes at least 50 ms. When a moving object is tracked by the camera in the gripper, the camera images are continuously

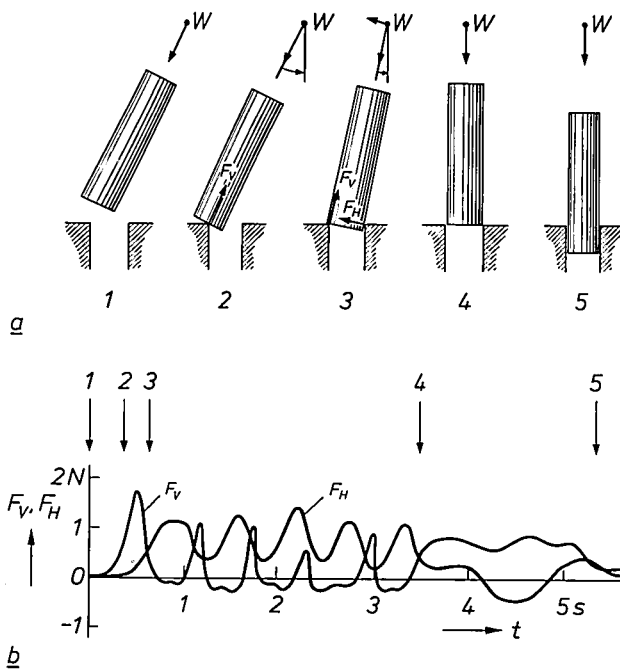


Fig. 10. a) The various stages in the assembly of a peg in a hole, using the information received from the force sensors in the fingers. Stage 1 is reached after coarse positioning with the aid of the position transducers in the arm drive. The peg is rotated about the axis DD (see fig. 2) to produce an apparently larger clearance. W point of intersection of the axis DD with the plane of the drawing. Stage 2: after a downward movement of the gripper in the direction of the centre-line of the peg, a force F_V is exerted in this direction that exceeds a critical limit value. The gripper is now slowly rotated back about the axis DD . Stage 3: the right-hand side of the wall is touched with a force F_H . Lateral corrections are also applied now to bring F_H back to a specified value. Meanwhile the peg is moved far enough towards the hole to maintain at all times a specified low value for F_V . Stage 4: the peg is in its vertical position and the rotary movement about DD stops. The point W is displaced further to try to bring F_V to a small target value, and F_H is reduced to an appropriate value if necessary. The result is that the peg moves down into the hole. Stage 5: the peg has now been inserted into the hole to the required depth and the movement stops. b) Recorded plots of the forces F_H and F_V as a function of time. Stages 1 to 5 from fig. 10a are indicated by arrows. The apparently temporarily negative value of F_V is due to uncertainty of the zero level as a result of noise.

processed in the control system at a frequency of about 10 Hz (a higher frequency will be used in the future). Since observation of the position where an object has to be assembled can be obstructed when the object is held in the gripper, the robot arm will probably be fitted with two semiconductor cameras in the future.

The fixed cameras are responsible for checking the result of the actions performed and for observing the position where the components are fed in (e.g. by conveyor belts). The images from the fixed TV cameras can be processed sequentially by the same electronic

circuit, because these cameras generally observe static objects. The complete duration of one robot action — one second or more — is available for processing these images.

The control system

The hardware

The heart of the control system is a Philips P851M minicomputer, which has a main memory with a capacity of 32×10^3 words of 16 bits. Peripherals con-

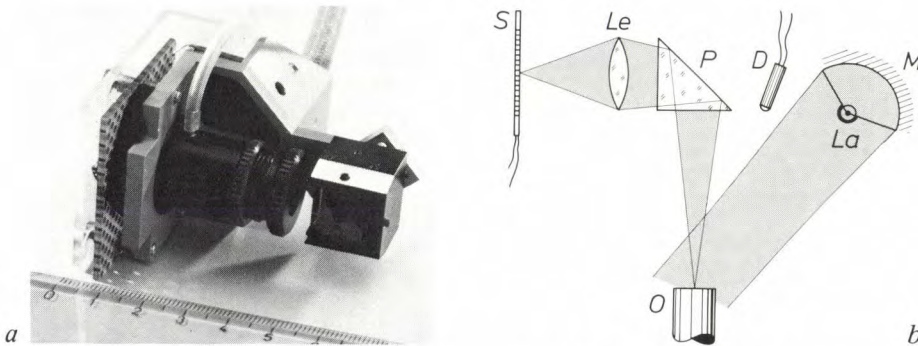
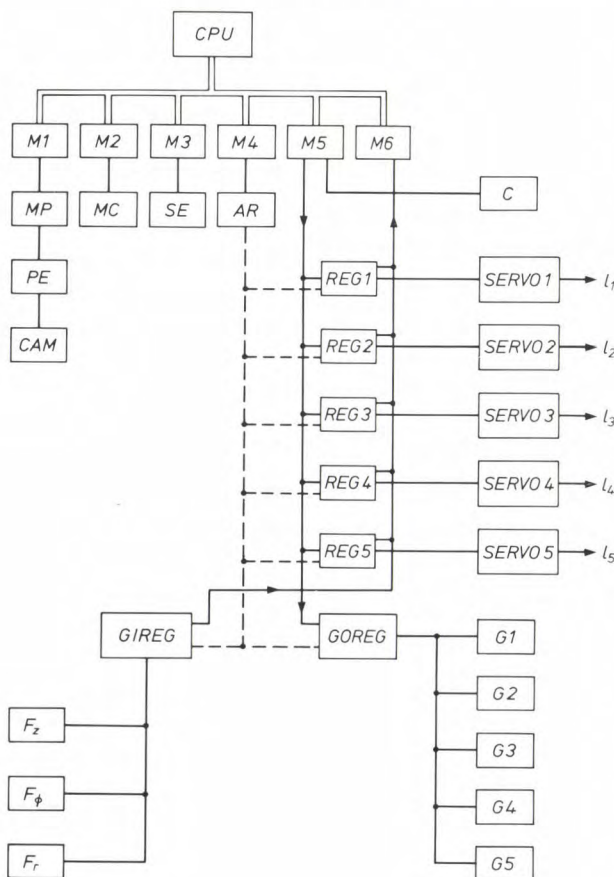


Fig. 11. a) Photograph of the semiconductor camera. b) The path of the rays in the semiconductor camera in the gripper. *S* detector, consisting of a matrix of 100×100 light-sensitive diodes. *Le* lens. *P* prism that deflects the optical axis through 90° . *O* object. *D* light-sensitive diode, which measures the reflected amount of light. *La* one of two incandescent lamps; a control circuit including the diode *D* keeps the intensity of the illumination of detector *S* constant. *M* one of the two reflectors. In practice the two lamps with reflectors are mounted in a plane perpendicular to that of the drawing.



nected to the computer are a display with keyboard and a floppy-disc unit that constitutes the backing memory with a storage capacity of $2 \times 250 \times 10^3$ bytes. Data exchange between the computer and the electronics of the robot control system is effected by electronic circuits belonging to the Philips MIOS6 interface system (Modular Input-Output System); see *fig. 12*. These circuits are connected directly to the

Fig. 12. Block diagram of the hardware. *CPU* central processing unit of the Philips P851M minicomputer, which is the heart of the control system. The double lines indicate the standardized transmission channel GPBS (General Purpose Bus Simplified) for digital information, the single lines represent the transmission of information, and the dashed lines indicate the allocation of the addresses. *M1* to *M6* electronic circuits of the Philips MIOS6 interface system for the input and output of the data. *MP* Signetics 8X300 microprocessor, which processes the information from the camera mounted in the gripper. *PE* the associated electronic circuits and the image store. *CAM* camera. *MC* panel for manual control of the robot. *SE* circuits for start and stop procedures and for monitoring the status of the system. *AR* 16 bit address register that indicates one of the registers *REG 1* to *5*, *GIREG* or *GOREG*. *C* clock generator. *REG 1* and *2* are 12 bit registers that receive the information for the required positions l_1 and l_2 of the servomotors *SERVO 1* and *2* for the wrist movement, via the interface circuit *M5*. *REG 3* to *5* are registers for the required positions l_3 to l_5 of the servomotors *SERVO 3* to *5* for the arm movement. The positions l_1 to l_5 of the pistons in the servomotors can be determined by the computer via the interface circuit *M6*. *GIREG* general 16 bit input register that passes on the information from the force-measuring systems for F_z , F_ϕ and F_r via *M6* to the computer. *GOREG* general 16 bit output register that controls the various gripper functions *G1* to *G5*. The choice from F_z , F_ϕ and F_r and *G1* to *G5* is made with the aid of an ancillary address register (not shown).

GPBS transmission channel (General-Purpose Bus Simplified), which handles the internal data flow in the minicomputer. In the version of the MIOS6 circuits that we have adopted the wiring of the computer and that of the peripherals are electrically isolated by optical couplers, which consist of a combination of an LED (light-emitting diode) and a light-sensitive transistor. This suppresses the greater part of any interference due to 'earth loops' in the hardware.

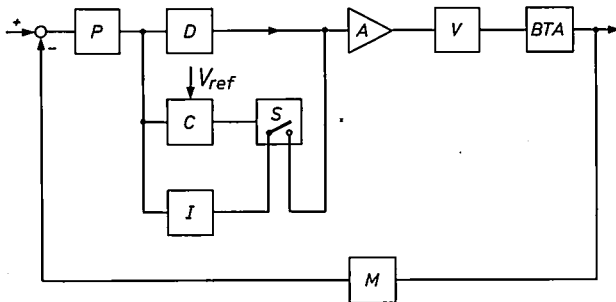


Fig. 13. One of the control systems associated with the servomotors for the arm and wrist movements. *P*, *D* and *I* are sections of the controller with proportional, differentiating and integrating actions, respectively. *C* comparator circuit, which compares the residual deviation with a preset limit value V_{ref} . *S* electronic circuit that switches on the integrating section *I* of the controller when the residual deviation reaches the limit value. *A* electronic amplifier. *V* electrically operated hydraulic servovalve. *BTA* servocylinder. *M* measuring system integrated with the servocylinder.

In addition to the computer peripherals, there is a manual control panel *MC* connected to the computer via an interface circuit *M2*. The user can give the robot single movement commands from this panel. The manual control panel is particularly useful for 'teaching' the robot how to perform one or more actions.

In the preceding section we mentioned the term 'sampling period'. At the start of each sampling period the control circuits of the servomotors that actuate the robot arm and gripper are provided with fresh information about the desired position. During each sampling period the new positions l_1 to l_5 of the pistons in the servomotors are calculated, using coordinate transformation. The duration of the sampling period is still 6 ms at present, but will be reduced in the future by providing each control circuit with its own microcomputer. The clock generator *C* signals the start of each sampling period. The computer then uses an address register *AR* to designate successively the registers *REG1* to *5* associated with the servomotor-control circuits. These registers are then filled with the calculated new values for l_1 to l_5 . This digital data is converted into an analog value by DACs (digital/analog converters) and fed to the appropriate

control circuits. The position of the robot at any given moment can be determined by reading out the contents of the registers *REG1* to *5* via the interface circuit *M6*.

Incorporated in the control system of the servomotors for the arm and wrist movements are hydraulic servovalves that control the oil supply to the cylinders. The controllers in these systems have PID (Proportional Integrating and Differentiating) charac-

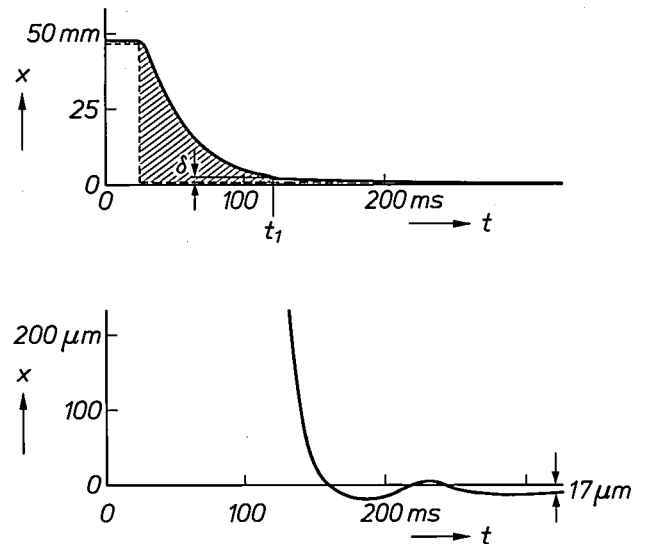


Fig. 14. The measured displacement x of the piston in the servocylinder as a function of time t . The lower diagram is a $250\times$ enlargement of the upper one; the first part of the curve is outside the figure. The dashed line indicates the desired displacement of the piston, the solid line the actual displacement (step-function response). At the time t_1 the integrating section *I* (see fig. 13) of the controller is switched on; the residual error has then reached the limit value δ . Since the hatched area is negated by the integrator, over-reaction from the controller is avoided. This means that there is a negligible overshoot and the resulting static deviation is only $17\ \mu\text{m}$.

teristics. In assembly robots, movement overshoot is not generally permitted, since workpieces or tools could be damaged. When the PID controller is used without any modification, some overshoot always occurs when the system responds to a step-function input signal. We have avoided this by not switching on the integrating part of the controller — which has the function of reducing the static deviation to nearly zero — until the residual error reaches a particular value. Fig. 13 shows a block diagram of the control system designed for this purpose. The switch *S* for the integrating action is operated by an electronic circuit *C* that compares the residual error with a preset limit value V_{ref} . Fig. 14 shows (as a result of measurements) that after application of a step corresponding to a very large piston displacement of 47 mm — which

does not occur in the actual assembly operation — the static deviation is only 17 μm and that overshoot is almost entirely avoided.

The gripper system of the robot is controlled by the computer via a general output register *GOREG*; see fig. 12. This register enables functions *G1* to *G5* to be performed, such as moving the gripper from the vertical to the horizontal position, opening or closing the fingers or adjusting the gripping force. The information about the magnitude of forces F_z , F_ϕ and F_r , originating from the circuits associated with the strain gauges in the fingers, reaches a general input register *GIREG* via an analog-to-digital converter, which processes the three signals sequentially.

The Signetics microprocessor *MP* mentioned earlier, with memory, controls the camera in the gripper and the storage and processing of the camera images. The positional information resulting from the image processing is passed to the computer via the interface circuit *MI*.

The hardware also contains a number of circuits *SE*, which are responsible for switching the system on, bringing the hydraulics up to pressure, switching the system off and monitoring it. In addition these circuits supply the computer with information about the status of the total system.

The software

The essential principle in the design of the software for EPAAS was that the user need have no special knowledge of computers to be able to compile working instructions for the robot. Such instructions, formulated in an 'application program', consist of a number of elementary operations such as gripping, displacing, inserting, releasing and so on. The user communicates to the system the operations — in chronological order — that make up his application program, and provides each operation with the appropriate parameters such as position coordinates, maximum velocity and force limits. The system then puts the operations together to form a complete program. A program can be stored in the computer's backing memory; programs already stored can be recalled and transferred to the main memory.

The set of 'statements', which in our case each contain a routine for performing an operation, is called the library. The statements are compiled to form a complete program by means of an all-purpose supervisory program, for which we have chosen the *FRAME* program. *FRAME*, also developed by Philips Research Laboratories, is capable of composing process-control programs consisting of routines that may each contain a maximum of 12 parameters. The routines that make up the library have to be entered in

FRAME by the system designer once for each field of application. *FRAME* is also responsible for checking the validity of the statement entered by the user, and offers facilities for correcting and adapting the application programs [8].

The commands that the user can give via *FRAME* include:

SAVE, for storing an application program in the external memory;

LOAD, for transferring an application program from the external memory to the main memory;

MANUAL CONTROL, for switching the system to manual control during the composition of an application program;

START, for initiating an application program;

WHERE, calling for the position and orientation of the gripper;

MACHINE STATUS, asking for the status of the system, including the data relating to 'WHERE';

SLAVE, an instruction, based on feedback of the force information, to adjust the position of the arm by manually moving an object in the gripper;

and the following 'EDIT' operations:

— 'hot editing', changing statements and parameters or both by interrupting an application program being run; after the change, the program can be continued at a point indicated by the user;

— 'immediate execution', stopping a program in progress, so that a statement with parameters can be input and immediately executed without it becoming part of the program;

— 'interactive control', stopping a program in progress that can then be continued (possibly in stages) at another point to be indicated and

— 'fill in the blanks', that is to say filling in during the execution of the program the values of parameters that have been left open.

The use of these commands leads to an application program in a continuous process of interaction between the user and the system, and the result of the elementary operations can at all times be monitored. The user thus 'teaches' the robot to carry out the program until a satisfactory result is achieved.

As already described, control circuits as illustrated in fig. 13 are used to produce the displacements calculated by the computer. The statements that make use of these circuits alone are *MOVE*, which enables the wrist joint *W* to describe a linear path, and *TRANSP*, which enables *W* to describe an elliptical path. The control circuits form the lowest hierarchical level in

[8] In many respects *FRAME* resembles the *INDA* program, which will be described by D. J. Burnett in a forthcoming article in this journal. *INDA*, however, is designed for larger computer systems than the P851M with peripherals, which we used, and can therefore offer more facilities.

the control of the arm; see *fig. 15*. Control circuits of a higher level use the force-measuring system in the gripper fingers or the camera in the gripper. The choice between the direct control C of the lowest control circuits or the two higher circuits with image or force information is made by means of the software 'switch' SW in the control program. The 'position' of SW is determined by the statement being dealt with and by the force information being processed in P_F . The statements that make use of force information

at the same time there is a rotation about the wrist joint W (see *fig. 10*). When the peg has reached its vertical position (stage 4), other limit values are passed to P_F via the dashed connecting line in *fig. 15*, and the downward movement starts in accordance with the statement *INSERT*. The parameter indicated for the insertion depth determines when the movement must stop.

The statements *CAMERA* and *FOLLOW*, which make use of image information, serve to position the gripper at the centre of gravity of the observed image or to follow this centre of gravity with the gripper.

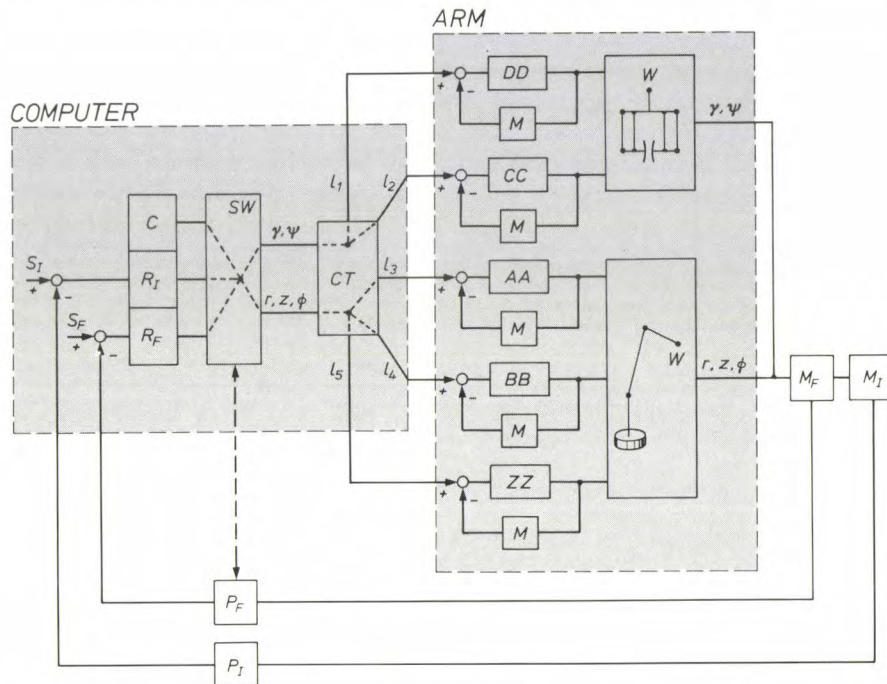


Fig. 15. Diagram of the control system. The main units, shaded grey, are the arm mechanism and the control program in the computer. SW 'switch' in the program, permitting a basic choice between direct control C , the controller R_F for processing the force signal or the controller R_I for processing the image signal. The information in the form of coordinates γ , ψ , r , z and ϕ are converted by coordinate transformation CT into positions l_1 to l_5 of the pistons in the servomotors. The arm sections are moved and the gripper rotated by means of the local control circuits for the rotations about the axes ZZ , AA , BB , CC and DD , which include the measuring systems M . M_F force-measuring system in the gripper fingers. M_I camera in the gripper. P_F processing of the signal from the force-measuring system. P_I processing of information from the camera, for determining the position of the centre of gravity of the image in relation to the centre of the camera. S_F and S_I set-points for the control circuits with, respectively, force information and image information. The values of S_F and S_I are input into the computer by the control program. SW can be altered via the connection represented by the dashed line if a threshold value in the force-measuring system is exceeded. The program also supplies P_F with the threshold values for the force measurement along this connection, but in the reverse direction.

are *RAISE* and *INSERT*. *RAISE* produces the movement from stage 1 to 4 in *fig. 10*, *INSERT* the movement from 4 to 5.

As the peg goes from stage 1 to stage 2 the calculated values received by the registers *REG 3* to 5 in *fig. 12* increase for each sampling period, so that the peg moves along its axis at a low velocity. In each period, however, a check on the magnitude of the force F_V in this direction is made via P_F . As soon as the specified limit value for this force is reached, SW is set to the required 'position' and the movement along the axis of the peg is stopped. Use is now made of the control loop with force information, while

To conclude, we give below a summary of the principal statements that the user can call from the existing library for compiling his application programs: *MOVE* and *TRANSP*; these are respectively linear and elliptical movements — built up from elementary circular arcs — of the wrist joint W to a point given in (r, z, ϕ) -coordinates. *TRANSP* requires a second point on the elliptical path. *CAMERA*; positioning the axis of symmetry of the gripper at the centre of gravity of an object observed in the gripper camera.


```

0000  XXXXXXXXXXXXXXXXXXXXXXXX
0001  % INSERT PEG IN HOLE %
0002  XXXXXXXXXXXXXXXXXXXXXXXX
0003
0004  INIPOS (1, 407,0,1075,12569,-356)
0005  INIPOS (2, 300,-614,12018,15461,-1517)
0006  INIPOS (3, 407,0,1075,12569,-356)
0007  GRIP (-40,40,-40,40)
0008  10 IXMOVE (1,20)
0009  WAIT (2000)
0010  CAMERA (0,0,0,20)
0011  SAVPOS (1)
0012  ITRANSP (2,3,50)
0013  CAMERA (0,0,200,30)
0014  GRIP (125,125,125,125)
0015  WAIT (200)
0016  ITRANSP (2,1,40)
0017  RESET ( )
0018  FORWIN (40,0)
0019  RAISE (10,0,1)
0020  FORWIN (100,0)
0021  INSERT (10,250)
0022  GRIP (-125,125,-125,125)
0023  GOTO (10)
:EOF
    
```

Fig. 16. Example of an application program for inserting a peg in a hole. The significance of the individual statements is explained in the text; the numbers between brackets represent parameters such as coordinates (γ , ψ , r , z and ϕ), velocities and forces. Position 1 is the location where the object with the hole is presented, position 2 is a point on the elliptical path between position 1 and position 3, the place where the peg is presented. The program proceeds as follows (the numbers refer to the statements): 4, 5, 6 indication of the initial coordinates of the positions 1, 2 and 3. 7 opening of the gripper. 8 linear movement to the assembly location. 9 waiting, during the period required for presenting the object containing the hole. 10, 11 observation of the hole and storage of the improved coordinates in the memory. 12 elliptical movement to the place where the peg is presented. 13 positioning at the centre of gravity of the peg with the aid of the camera. 14, 15, 16 the peg is picked up and moved to the assembly location; waiting is necessary for moving the fingers to the peg. 17, 18, 19 the gripper is raised from the tilted position, following from statement 4. 20, 21 the actual assembly (insertion). 22 the peg is released. 23 return to the beginning of the program for the next assembly operation; since statement 4 is now negated, the system possesses more accurate information concerning the position where the objects with the holes are presented. The minimum permissible clearance between peg and hole is 10 μm .

FOLLOW; following the centre of gravity of a moving object observed in the camera.

RAISE; a tilting and raising movement of an object in the gripper relative to a stationary object, with feedback of the force information; the angle at which the rotation stops and the areas for F_V and F_H (see fig. 10) within which there is no response to a change in force have to be specified.

INSERT; the assembly of the object in the gripper and a stationary object, with feedback of the force information. The areas for F_V and F_H within which there is no response to a change in force follow from the statement **FORWIN**.

FORWIN; specification of limit values for the forces.

RESET; resetting to zero the force-measurement signals (before **FORWIN**, **INSERT** and **RAISE** are used).

GRIPV and **GRIPH;** moving the gripper to its vertical or horizontal position (see fig. 2).

GRIP; opening or closing the gripper fingers with a maximum gripping force to be specified for each pair of fingers.

INIPOS; defining a point in the working space in (r, z, ϕ)-coordinates.

SAVPOS; measuring the instantaneous position of W in the working space; these coordinates then replace the coordinates that were given with **INIPOS**.

IXMOVE and **IXTRANSP;** movements resembling **MOVE** and **TRANSP**, but the points of destination are derived from the coordinates defined by **INIPOS** or corrected by **SAVPOS**.

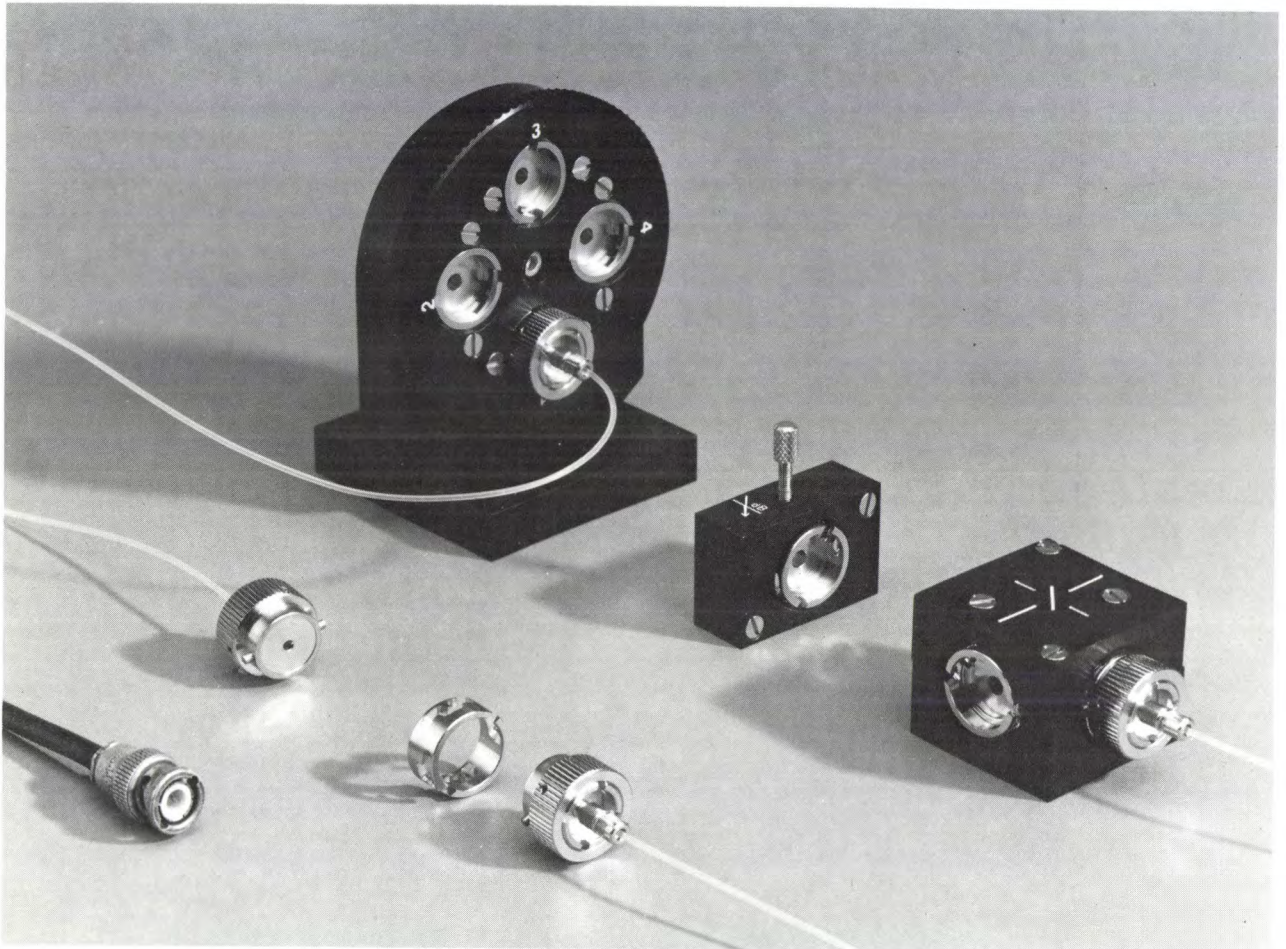
With the last four statements the system can be made to 'learn' from previous experience, thus achieving a kind of 'artificial intelligence'. The assembly location — this is the location where the devices in which the objects will be assembled are presented — can first be indicated roughly by **INIPOS**. Once **IXMOVE** or **IXTRANSP** has brought the gripper to this location it can be positioned more accurately by **CAMERA** at the centre of gravity of the assembly location. The improved coordinates are entered in the memory by means of **SAVPOS**. The actual assembly operation then takes place with the statements **RAISE** and **INSERT**. When an identical assembly has to be carried out a second time, the system possesses more accurate information about the position of the assembly location than it did for the first assembly operation.

The procedure described here is illustrated by the program given in *fig. 16* for the assembly of a peg in a hole.

R. J. Asjes, L. Leistra, F. J. de Munnik and G. J. Scholl also contributed to the work described in this article.

Summary. Hard automation, the established form of industrial mechanization, now has to compete with flexible automation, which makes use of industrial robots. To gain experience with the industrial robot as a new tool of factory mechanization, work started at Philips Research Laboratories some years ago on designing and building an experimental industrial robot. This robot, called EPAAS (for 'Experimental Programmable Adaptive Assembly System') has meanwhile developed into a second-generation robot and further additions will be made to turn it into a robot of the third generation. The arm mechanism possesses five degrees of freedom, and in its kinematic design it resembles a human arm. The fingers of the gripper are equipped with force sensors; to avoid collisions, the arm is fitted with proximity detectors; the gripper possesses 'vision' in the form of a semiconductor camera that moves with the arm. There are also fixed TV cameras that observe the working space of the robot. The heart of the control system is a Philips P851M minicomputer. The software, which includes the **FRAME** supervisory program, has been designed in such a way that the user can interact continuously with the system to 'teach' the robot to carry out the assembly program. The software contains a library of routines that can be compiled in a desired sequence to form a complete application program. This is illustrated with an example in the form of a program for inserting a peg in a hole.

Components for glass-fibre circuits



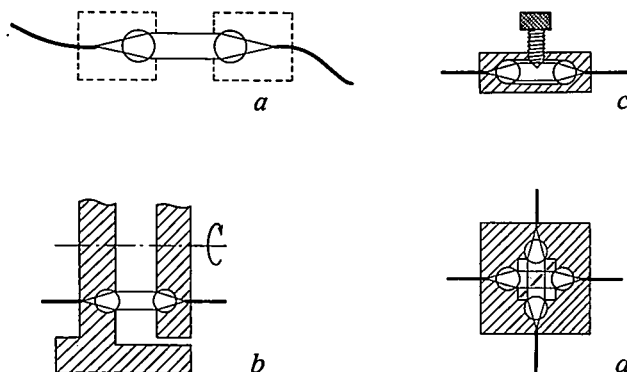
The use of glass fibres for optical telecommunication will require not only components and techniques for making permanent connections^[1] but also demountable couplings for fibres and subsidiary devices for various operations on the optical signals that travel along the glass fibres. Such components are required both for measurements in glass-fibre cable networks and for laboratory experiments. For microwave signals in waveguides or coaxial cables there are innumerable subsidiary devices of this type, and their use has been taken for granted for so long that sometimes people hardly realize that their design once involved a painful process of trial and error.

The title photograph shows a number of components that have been designed for glass-fibre circuits. They are all based on the principle devised for a dry coupling^[2], in which the divergent beam emerging from a fibre is converted into a parallel beam by a convex lens, and a second convex lens concentrates the beam on the end of the next glass fibre (*fig. 1a*). In this way, depending on the type of glass fibre, a coupling efficiency of 85-90% can be achieved, without requiring the use of an immersion liquid.

[1] A. J. J. Franken, G. D. Khoe, J. Renkens and C. J. G. Verwer, *Philips tech. Rev.* **38**, 158, 1978/79.

[2] A. J. A. Nicia, *Appl. Optics* **20**, 3136, 1981.

Fig. 1. *a*) Principle of a dry coupling (i.e. with no immersion liquid) for two glass fibres. The ends of the two fibres are situated at the focal points of two convex lenses. *b*) A switch for optical signals consists of two plates, one fixed and the other rotatable; in both plates there are four half-couplings. *c*) To attenuate an optical signal a screw is inserted in the expanded beam between the two convex lenses. *d*) A half-silvered mirror mounted as shown here between four half-couplings in a straight glass-fibre — connected to two opposite halves — can extract a fraction of the signal passing through the cable via one of the two other half-couplings. A signal can be fed on to the straight glass-fibre via the fourth half-coupling.



A conventional coaxial cable with connector is shown on the far left of the photograph to give an idea of the size of the components. Beside it there is a glass-fibre coupling that works on the principle just described; between the two halves of the coupling is the ring needed for connecting them together. To ensure good coupling efficiency the two halves must be designed in such a way that the emergent beam is accurately perpendicular to the front face.

Since the coupling efficiency is still good at a somewhat greater spacing between the two halves, they can be spaced apart to permit various operations to be performed on the beam in the space between them. This is the principle of the components shown in the back row of the photograph. There are four half-couplings in a rotary disc facing another four in a fixed disc, and the result is a switch (see fig. 1*b*).

A screw can be introduced into the beam between the two halves to provide variable attenuation (fig. 1*c*); a half-silvered mirror, placed diagonally between four half-couplings arranged in a square, allows signals to be coupled in and out in a straight line (fig. 1*d*).

A. J. A. Nicia
C. J. T. Potters

Ir A. J. A. Nicia and C. J. T. Potters are with Philips Research Laboratories, Eindhoven.

The pigmentation of phosphors for colour television

K. Carl, J. A. M. Dikhoff and W. Eckenbach

In colour television the reflection of the ambient light by the screen has an adverse effect on the contrast. To reduce this reflection use has recently been made of colour filters: the phosphor grains on the screen of the picture tube are coated with coloured inorganic pigment particles. The pigment for a phosphor is selected so that it transmits as much as possible of the light emitted by the phosphor and absorbs as much as possible of the ambient light. Scientists at Philips Forschungslaboratorium Aachen and the Philips Elcoma Division, Eindhoven, have carried out calculations and experiments on model systems to see how useful this filter concept can be for improving the contrast and brightness of a colour television receiver.

The filter concept for colour picture tubes

Brightness and contrast are important quantities in the picture quality of a colour television receiver. The brightness is determined by the intensity of the light that is emitted by the three phosphors for the primary colours red, green and blue when they are excited by electrons. The contrast is determined by the ratio of the intensity of the emitted light to that of the ambient light reflected by the picture-tube screen [1].

Both the brightness and the contrast are increased as more light is produced from the picture tube, e.g. by using phosphors of higher efficiency. A further improvement in the contrast can be obtained by decreasing the reflection of ambient light. This can be done by making the glass of the screen darker or filling the empty spaces between the phosphor areas of the screen with a black powder.

In a third method, which first came into use a short time ago [2][3], use is made of colour filters, normally in the form of coloured inorganic pigments. In principle, each of the three phosphors has a suitable filter material applied to it, which transmits as much as possible of the phosphor light and absorbs any other light. The use of such filters has the advantage that the reflection of the ambient light can be considerably reduced without appreciably affecting the brightness. A further advantage is that the choice of phosphors is

a little wider, since with the filters some correction can be made for a less suitable chromaticity of the phosphor emission.

Fig. 1 gives diagrams of two practical versions of the filter method for a conventional colour-tube screen, in which the phosphor light emitted isotropically by electron excitation is all directed towards the viewer because of the presence of a thin reflecting aluminium coating. In one version the filter consists of a layer between the phosphor layer and the glass of the screen; hence the name, filter sandwich. Since both the phosphor light and the ambient light pass through the filter layer it must of course be applied in exactly the same geometry as the dots or lines of the associated phosphor. In the other version the phosphor is mixed with the corresponding pigment. In practice this mixed version is preferable because of its simple technology and it is now widely used.

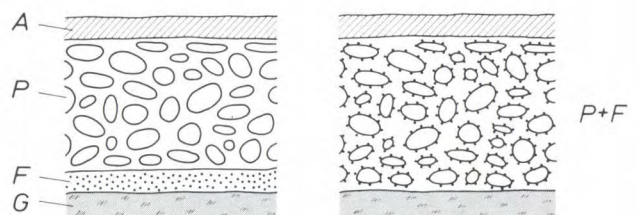


Fig. 1. Two practical versions of the filter method for the screen of a colour television receiver. *Left:* sandwich with screen glass G, filter layer F, phosphor layer P and aluminium layer A. *Right:* version with a mixture of phosphor and pigment.

Dr K. Carl and Dr W. Eckenbach are with Philips GmbH Forschungslaboratorium Aachen, Aachen, West Germany, and J. A. M. Dikhoff is with the Philips Electronic Components and Materials Division (Elcoma), Eindhoven.

Some years ago we started an investigation into the theoretical possibilities of the filter method for colour television. We were chiefly concerned with the filter sandwich, since it constitutes a 'model system' that will give a relatively simple description of the filter method. The brightness and contrast measured for filter sandwiches give good agreement with the results of the model calculations. This gives good support to the theoretical predictions relating to the choice of pigment, the optimum pigment concentration and the theoretical possibilities of the filter method. Some of the results are also applicable to the use of the theoretically less accessible mixtures of phosphor and pigment.

We shall now first consider the theory of the filter sandwich and then compare the results of the calculations with those of the measurements. Finally, we shall examine mixtures of phosphor and pigment as used in practice.

The filter sandwich

We shall refer to the four layers of the filter sandwich that we have investigated as *G* (glass), *F* (filter), *P* (phosphor) and *A* (aluminium); see fig. 1. The most important experimental variations relate to *F*. Since a *GFPA* sandwich is relatively difficult to make, we made most of our optical reflection and transmission measurements with combination *GF*. By using the model that we have developed we can derive the

brightness and contrast of a sandwich from these measurements. Our model relates the optical properties of single layers to those of the sandwich. In determining the properties of the sandwich the spectral luminous-efficiency function, the spectrum of the phosphor emission and the spectrum of the ambient light are all taken into account. The characteristics of a picture tube with a filter are normalized with respect to those of a tube without a filter. We assume that all effects are one-dimensional, with no lateral or angular effects. We also assume that there are no interference effects in the diffusely reflecting phosphors and pigments, so that the intensities of the light can be summed. We assume further that there is no optical contact between the glass and the filter; in practice this means that there is a small intervening space between the two.

If a homogeneous layer of reflectance ρ_1 and transmittance τ_1 is combined with a second homogeneous layer of reflectance ρ_2 and transmittance τ_2 , then incident light gives rise to a number of reflections between the adjacent surfaces of the two layers; see fig. 2. The transmittance τ_{12} for the two layers combined then becomes

$$\begin{aligned} \tau_{12} &= \tau_1\tau_2(1 + \rho_1\rho_2 + \rho_1^2\rho_2^2 + \dots) \\ &= \tau_1\tau_2/(1 - \rho_1\rho_2). \end{aligned} \tag{1}$$

This transmittance is independent of the side of incidence: $\tau_{12} = \tau_{21}$. The reflectance ρ_{12} of the combination of the two layers is given by

$$\begin{aligned} \rho_{12} &= \rho_1 + \tau_1^2\rho_2(1 + \rho_1\rho_2 + \rho_1^2\rho_2^2 + \dots) \\ &= \rho_1 + \tau_1^2\rho_2/(1 - \rho_1\rho_2). \end{aligned} \tag{2}$$

This reflectance does however depend on the side of incidence: $\rho_{12} \neq \rho_{21}$. The reflectance and transmittance of a single layer are obtained by solving equations (1) and (2):

$$\rho_1 = \frac{\tau_2^2\rho_{12} - \rho_2\tau_{12}^2}{\tau_2^2 - \rho_2^2\tau_{12}^2}, \tag{3}$$

$$\tau_1 = \tau_2\tau_{12} \frac{1 - \rho_2\rho_{12}}{\tau_2^2 - \rho_2^2\tau_{12}^2}. \tag{4}$$

If the first layer acts as a planar light source, the

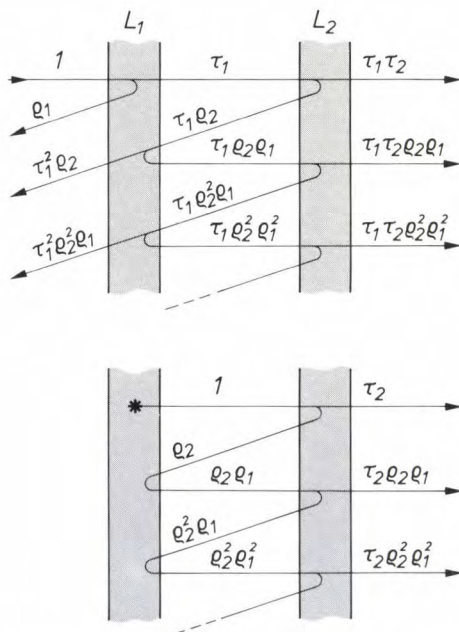


Fig. 2. Schematic representation of the light intensities for two layers L_1 and L_2 , with reflectances ρ_1 and ρ_2 and transmittances τ_1 and τ_2 . *Top*: for light incident on L_1 . *Bottom*: for light generated in L_1 .

[1] In this article we are only concerned with the diffuse reflection by the phosphor particles; the directional reflection from the glass of the screen is not considered.
 [2] S. S. Trond, Extended Abstracts Electrochem. Soc. **78-1** (Spring Meeting, Seattle 1978), p. 817.
 T. Nishimura, T. Takahara, S. Matsuura, M. Watanabe and T. Kawamata, *ibid.*, p. 820.
 [3] T. Takahara, T. Wakatsuki and T. Nishimura, Extended Abstracts Electrochem. Soc. **80-1** (Spring Meeting, St. Louis 1980), p. 564.
 W. Möller, W. de Rave and H. Widmann, *ibid.*, p. 576.
 K. Carl, J. A. M. Dikhoff, W. Eckenbach and H. G. Junginger, *ibid.*, p. 579, and J. Electrochem. Soc. **128**, 2395, 1981.

optical transmittance $\tau_{(1)2}$ through the second layer (fig. 2) is given by

$$\begin{aligned}\tau_{(1)2} &= \tau_2(1 + \varrho_1\varrho_2 + \varrho_1^2\varrho_2^2 + \dots) \\ &= \tau_2/(1 - \varrho_1\varrho_2).\end{aligned}\quad (5)$$

Equations (1), (2) and (5) can be extended for a sandwich consisting of three or more homogeneous layers, e.g.

$$\tau_{123} = \tau_{12/3} = \tau_{12}\tau_3/(1 - \varrho_{21}\varrho_3), \quad (6)$$

where τ_{12} and ϱ_{21} can be derived with the aid of equations (1) and (2). The location of the split is of no significance: $\tau_{12/3} = \tau_{1/23}$.

The reflectance ϱ_{GFPA} of a filter sandwich *GFPA* can be obtained by splitting it into *GF* and *PA*:

$$\varrho_{GFPA} = \varrho_{GF} + \tau_{GF}^2\varrho_{PA}/(1 - \varrho_{FG}\varrho_{PA}), \quad (7)$$

where ϱ_{GF} is not equal to ϱ_{FG} . By splitting in a similar way we can use equation (5) to obtain the transmittance $\tau_{(AP)FG}$ of the light from the planar source *AP* through *FG*:

$$\tau_{(AP)FG} = \tau_{FG}/(1 - \varrho_{FG}\varrho_{PA}). \quad (8)$$

We used equations (7) and (8) to find ϱ_{GFPA} and $\tau_{(AP)FG}$ from the measured values of ϱ_{FG} and τ_{FG} as a function of the wavelength of the incident light. In these calculations we took a constant value of 0.80 for ϱ_{PA} . We used equations (2), (3) and (4) to derive the values of ϱ_{GF} from the measured values of ϱ_{FG} and τ_{FG} and from the reflectance ϱ_G and transmittance τ_G of the glass. For a model with two interfaces, ϱ_G and τ_G can be calculated from the reflectance r of the glass surface and the internal glass transmittance θ :

$$\varrho_G = r + \frac{\theta^2 r(1-r)^2}{1 - \theta^2 r^2}, \quad (9)$$

$$\tau_G = \frac{\theta(1-r)^2}{1 - \theta^2 r^2}. \quad (10)$$

The measured reflectance ϱ_{FG} and transmittance τ_{FG} are plotted against the wavelength in fig. 3 for two red pigments, Cd(S,Se) and α -Fe₂O₃. As might be expected, the reflectance curves resemble the transmittance curves. The transition at about 600 nm from low to high reflectance and transmittance is considerably steeper for Cd(S,Se) than for α -Fe₂O₃, and this gives a better filter effect. The calculated reflectance ϱ_{GFPA} of the sandwich *GFPA* is also plotted against the wavelength in fig. 3, as is the calculated transmittance $\tau_{(AP)FG}$ of the light from *AP* through *FG*. These calculations were made for the case in which the internal glass transmittance θ is equal to unity, so that $\varrho_G + \tau_G = 1$; see equations (9) and (10). The transitions from low to high reflectance and transmittance

have become much steeper because of the reflection from *PA*. In the wavelength range of the red-emitting colour-television phosphor Y₂O₂S:Eu the transmittance $\tau_{(AP)FG}$ is high, especially with Cd(S,Se), so that the phosphor emission is only slightly impaired.

The measured reflectance ϱ_{FG} and transmittance τ_{FG} for the pigment cobalt blue are given in fig. 4, with the calculated reflectance ϱ_{GFPA} and transmittance $\tau_{(AP)FG}$. Here again we see that the transitions from low to high reflectance and transmittance, in this case at about 500 and 700 nm, are considerably steeper for the sandwich than for *FG*. The high value of $\tau_{(AP)FG}$ in the wavelength range of the blue-emitting colour-television phosphor ZnS:Ag shows that the phosphor emission is again only slightly impaired.

The average reflectance R for a sandwich *GFPA* of ambient light with a spectral distribution N_λ is determined not only by ϱ_{GFPA} and N_λ but also by the spectral luminous efficiency function V_λ :

$$R = \frac{\int \varrho_{GFPA} V_\lambda N_\lambda d\lambda}{\int V_\lambda N_\lambda d\lambda}. \quad (11)$$

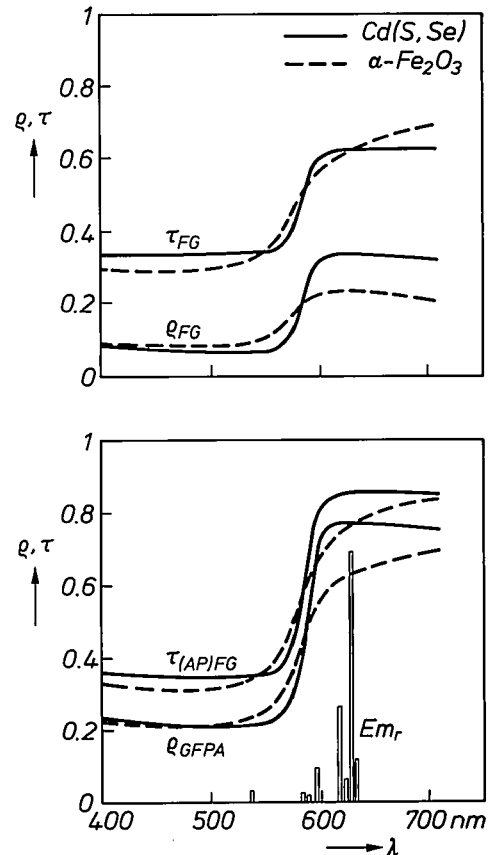


Fig. 3. Reflectances and transmittances as a function of the wavelength λ . *Top*: measured reflectance ϱ_{FG} and transmittance τ_{FG} for two red pigments on glass. The transitions to a high reflectance and transmittance at about 600 nm are considerably steeper for Cd(S,Se) than for α -Fe₂O₃. *Bottom*: the reflectance ϱ_{GFPA} of the filter sandwich in fig. 1, calculated from ϱ_{FG} and τ_{FG} , and the calculated transmittance $\tau_{(AP)FG}$ through *FG* of the light generated in *AP*. Here *P* is the red-emitting phosphor Y₂O₂S:4.25%Eu, whose emission spectrum Em_r is shown schematically.

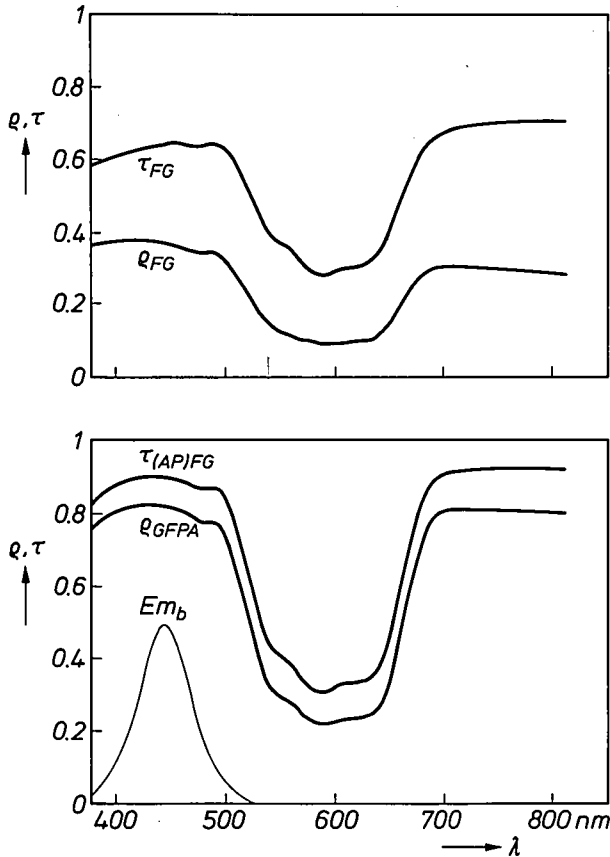


Fig. 4. Reflectances and transmittances as a function of the wavelength λ . Top: the measured R_{FG} and τ_{FG} for cobalt blue on glass. Bottom: the R_{GFPA} and $\tau_{(AP)FG}$ calculated from R_{FG} and τ_{FG} , where P is the blue-emitting phosphor ZnS:Ag with the emission spectrum Em_b .

The average transmittance T for a particular emission colour is determined by $\tau_{(AP)FG}$ and V_λ and by the spectral distribution E_λ of the emission from the phosphor:

$$T = \frac{\int \tau_{(AP)FG} V_\lambda E_\lambda d\lambda}{\int V_\lambda E_\lambda d\lambda} \quad (12)$$

The contrast for a particular emission colour is given by the ratio of the brightness B of the emitted light to the brightness of the reflected light, which is proportional to the intensity of the ambient light and the reflectance R . The contrast is generally expressed in terms of the quantity B/R , a ratio that does not depend on the ambient light. The improvement in quality resulting from colour filtering is usually expressed as the geometric mean of B and B/R . The relative value of this mean, the 'brightness/contrast performance', is defined as

$$BCP = \frac{(B/\sqrt{R}) \text{ with filter}}{(B/\sqrt{R}) \text{ without filter}} \quad (13)$$

Since the brightness is proportional to the transmittance T in equation (12), we can also write:

$$BCP = \frac{(T/\sqrt{R}) \text{ with filter}}{(T/\sqrt{R}) \text{ without filter}} \quad (14)$$

The calculated BCP -values are given in Table I for a number of combinations of well-known colour-television phosphors with suitable pigments. The use of a pigment yields an obvious improvement for red and blue. This is not the case for green, however: $BCP < 1$. We also calculated BCP -values for combinations of well-known phosphors with non-existent 'ideal' filters that pass everything in the emission range of the phosphor and nothing outside that wavelength range. The calculations for such filters naturally result in much higher BCP -values. Still higher values are obtained if the calculations are based on phosphors with a monochromatic emission in combination with filters that have 100% transmission over a wavelength range of only 10 nm and absorb all the incident light outside this range. For the three kinds of combinations we find that the BCP -value invariably increases in the same sequence — green, red, blue. This is because the sensitivity of the eye (the spectral luminous efficiency function) to the phosphor emission decreases in this sequence.

The absorption and emission spectra of existing pigments and phosphors can be moved along the wavelength axis by altering their compositions slightly. Our calculations, however, have shown that appreciable improvements cannot be obtained in this way. Varying the pigment composition improves the BCP -value by one per cent at most. The maximum increase that can be obtained by displacing the phosphor-emission spectra is five per cent, but this is accompanied by an unacceptable change in the emission colour.

Table I. Calculated BCP -values (equation (14)) for a number of existing and hypothetical combinations of phosphor and filter in a sandwich.

Colour	Phosphor	Filter	BCP
red	$Y_2O_2S:4.25\%Eu$	0.099 mg/cm ² Cd(S,Se)	1.25
	$Y_2O_2S:4.25\%Eu$	0.016 mg/cm ² $\alpha\text{-Fe}_2O_3$	1.17
	$Y_2O_2S:4.25\%Eu$	ideal, 605-635 nm	1.99
	monochromatic, 620 nm	ideal, 10 nm wide	3.39
green	(Zn,Cd)S:Cu	0.06 mg/cm ² cobalt green	0.94
	(Zn,Cd)S:Cu	ideal, 495-585 nm	1.09
	monochromatic, 540 nm	ideal, 10 nm wide	2.74
blue	ZnS:Ag	0.22 mg/cm ² cobalt blue	1.30
	ZnS:Ag	0.25 mg/cm ² ultramarine	1.45
	ZnS:Ag	ideal, 405-495 nm	2.18
	monochromatic, 450 nm	ideal, 10 nm wide	4.02

Experiments with filter sandwiches

We used various methods for making filter sandwiches on glass plates of dimensions 50×50 mm. In the simplest method the phosphor and pigment particles were precipitated from a suspension in an organic liquid. We also used electrophoresis, especially for the pigments; this is a method in which the particles in the suspension are given an electric charge, and the suspension is then passed through an electric field, so that the particles migrate. Another method is

the voltage, current density and excitation time corresponded as closely as possible to those of the Philips 30AX colour tube. The light intensity was measured with a detector whose sensitivity was very similar to that of the human eye. This detector was built into the PSEM, close to the specimen to be investigated.

We determined the *BCP*-values for red and blue for a number of sandwiches with different pigments and different surface concentrations of pigment. The experimental and calculated *BCP*-values for red are

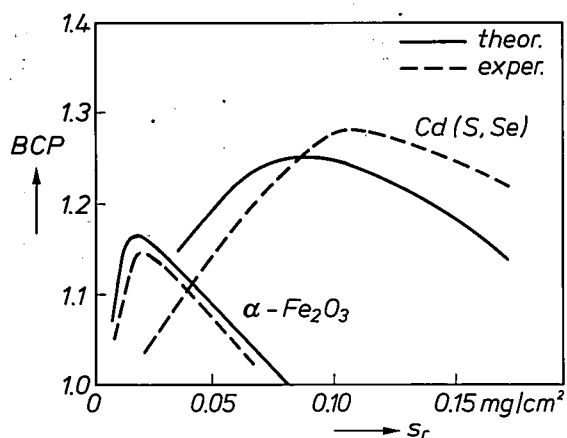


Fig. 5. *BCP*-values for a filter sandwich with the red-emitting phosphor $Y_2O_2S:4.25\%Eu$ as a function of s_r , the surface concentration of red $Cd(S,Se)$ or $\alpha-Fe_2O_3$ pigment. There is good agreement between the experimental curves (dashed lines) and the calculated curves (solid lines) both for the maximum *BCP*-values and for the optimum surface concentrations.

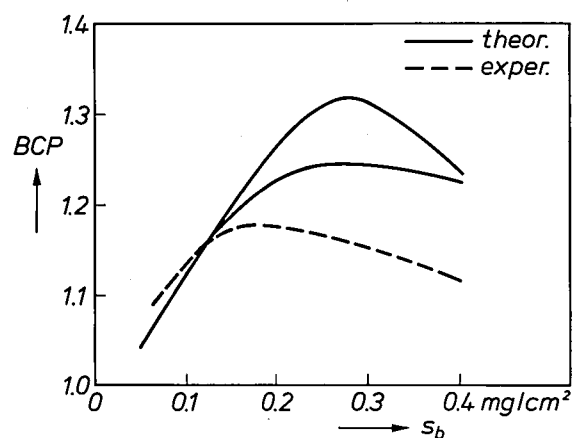


Fig. 6. *BCP*-values for a filter sandwich with the blue-emitting phosphor $ZnS:Ag$ as a function of s_b , the surface concentration of cobalt blue. The two calculated curves (solid lines) were obtained from measurements on differently prepared pigment layers. Partly because of the different preparation methods the agreement with the experimental curve (dashed line) is not as good as for the red.

Table II. Calculated and measured maximum *BCP*-values, BCP_{max}^{th} and BCP_{max}^{ex} , and the changes ΔX and ΔY produced by the filter in the chromaticity coordinates of the phosphor emission, for three combinations of phosphor and filter in a sandwich.

Colour	Phosphor	Filter	BCP_{max}^{th}	BCP_{max}^{ex}	ΔX	ΔY
red	$Y_2O_2S:4.25\%Eu$	$0.13 \text{ mg/cm}^2 \text{ Cd}(S,Se)$	1.27	1.28	0.010	-0.010
	$Y_2O_2S:4.25\%Eu$	$0.02 \text{ mg/cm}^2 \alpha-Fe_2O_3$	1.17	1.15	0.007	-0.007
blue	$ZnS:Ag$	$0.28 \text{ mg/cm}^2 \text{ cobalt blue}$	1.32	1.17	<0.001	-0.002

to centrifuge a suspension containing a lacquer, e.g. nitrocellulose or polyvinyl alcohol with ammonium dichromate. After the pigment and the phosphor layer have been deposited, a nitrocellulose layer is deposited to act as a temporary carrier for a $0.2 \mu m$ evaporated aluminium layer. Finally, all organic constituents are removed by heating the sandwich at $430^\circ C$.

The brightness of the filter sandwiches when subjected to cathode-ray excitation was measured with a Philips PSEM 500 scanning electron microscope. We arranged the excitation conditions in such a way that

plotted in *fig. 5* against the surface concentration of $Cd(S,Se)$ and $\alpha-Fe_2O_3$ pigments. In both cases there is good agreement between experiment and theory, both for the maximum *BCP*-value and for the optimum surface concentration.

When cobalt blue is used, the agreement is not so good; see *fig. 6*. This seems to be associated with the way in which the pigment is applied: results of model calculations for layers of pigments applied in different ways vary considerably. Table I indicates that ultramarine should produce better results than cobalt blue. In practice, however, this pigment is not suitable,

because it is not stable at the high temperatures encountered in tube manufacture.

Table II gives the maxima for the calculated and experimental *BCP*-values for red and blue. This table also indicates how the *X* and *Y* chromaticity coordinates of the phosphor emission are affected by the addition of a pigment. We determined these coordinates in the usual way by calculating the integrals of the phosphor-emission spectrum E_λ , multiplied by the transmittance $\tau_{(AP)FG}$ and the colour-matching functions, and normalizing the result^[4]. When a red pigment is used, *X* increases while *Y* decreases by about the same amount: the emission becomes redder and less green. These changes, however, are only three to four times as large as the 'just noticeable colour difference'^[5], so that they are not felt to be objectionable. When cobalt blue is used as a pigment, *X* is practically unaffected, while the decrease in *Y* is hardly noticeable.

It may even be possible to put the colour shifts to use sometimes, since they can help to compensate for a slight green shift in some of the efficient red- or blue-emitting phosphors. For example, if the amount of europium — an expensive material — in $Y_2O_2S:Eu$ is decreased from 4.25% to 3%, the phosphor emission becomes less red and more green: $\Delta X = -0.009$, $\Delta Y = +0.009$. If 0.13 mg/cm^2 of $Cd(S,Se)$ is used as a pigment, this shift is completely compensated (Table II).

Mixtures of phosphor and pigment

As we have already stated, a mixture of a phosphor and a pigment is preferable in practice to a filter sandwich because it is easier to make. In such a mixture the surface of the phosphor particles is partly covered by pigment particles. When the mixture is being prepared an organic binder is usually added to ensure that the pigment particles adhere well to the phosphor particles. The diameter of the pigment particles is of the order of $0.1 \mu\text{m}$ while that of the phosphor particles is 5 to $10 \mu\text{m}$. Fig. 7 shows an electron-microscope photograph of a red-emitting phosphor $Y_2O_2S:Eu$ pigmented with $\alpha\text{-Fe}_2O_3$.

The *BCP*-value for a mixture is lower than that of the corresponding sandwich: for the combination of $Y_2O_2S:Eu$ with $\alpha\text{-Fe}_2O_3$ the maximum value is 1.07 for a pigment concentration of about 0.3%; for the combination of $ZnS:Ag$ with cobalt blue *BCP* is at most 1.10 for a pigment concentration of about 4%. In both cases the improvement in *BCP* is approximately half of the improvement obtained with the filter sandwich (Table II). For the given surface concentrations of the phosphors — about 3.5 mg/cm^2 for

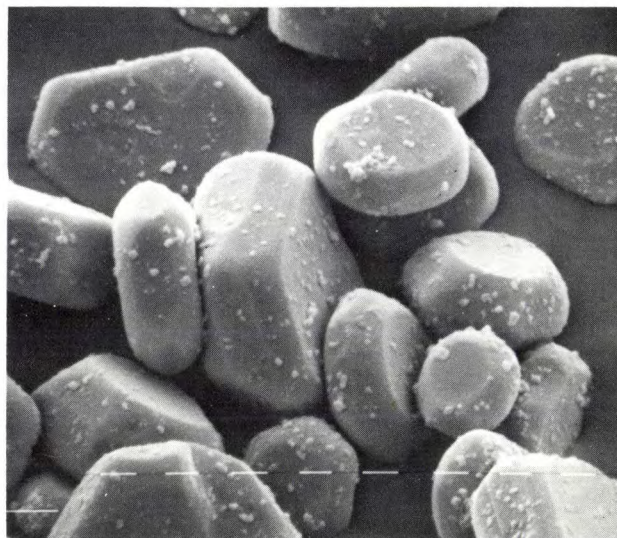


Fig. 7. Electron-microscope photograph (the dashes each represent $1 \mu\text{m}$) of a mixture of the red-emitting phosphor $Y_2O_2S:Eu$ and the red pigment $\alpha\text{-Fe}_2O_3$. The phosphor grains, with a diameter of $5\text{-}10 \mu\text{m}$, are coated by much smaller pigment particles, with a diameter of the order of $0.1 \mu\text{m}$.

red and about 3 mg/cm^2 for blue — the optimum surface concentrations of the pigments are about 0.01 mg/cm^2 for $\alpha\text{-Fe}_2O_3$ and about 0.12 mg/cm^2 for cobalt blue. These values are also about half of the optimum values when filter sandwiches are used. It would appear that for low pigment concentrations the *BCP*-value of a mixture increases at about the same rate as in a sandwich, but that the increase with pigment concentration continues only half as far as it does in a sandwich.

The main cause of the difference in behaviour is probably that in a mixture the ambient light and the phosphor light describe very similar paths, while in a sandwich the ambient light passes through the filter layer twice whereas the phosphor light passes through it only once. The light-scattering mechanism is also different. In a mixture with a low pigment concentration the optical interaction between the phosphor and pigment consists chiefly in scattering at the phosphor particles, which are much larger than the wavelength of the light ($0.4\text{-}0.7 \mu\text{m}$). In a sandwich, on the other hand, light is scattered at particles within the pigment layer, which are smaller than the wavelength.

For a quantitative description it may be assumed, in view of the small depth of penetration of the exciting electrons, that the phosphor light is only generated in a very small region in the neighbourhood of the aluminium layer. If we regard the layer with the phosphor and pigment as a single filter layer and assume that there is a thin light-emitting layer between this

[4] G. Wyszecki and W. S. Stiles, *Color science*, Wiley, New York 1967.

[5] D. L. MacAdam, *J. Opt. Soc. Amer.* 32, 247, 1942.

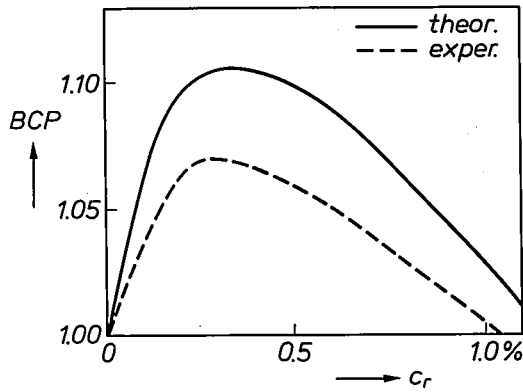


Fig. 8. *BCP*-values for a mixture of the red-emitting phosphor $Y_2O_2S:Eu$ and the red pigment $\alpha-Fe_2O_3$ as a function of the pigment concentration c_r . The experimental curve (*dashed line*) is lower than the curve calculated from the simple model (*solid line*). There is good agreement for the optimum pigment concentration, however.

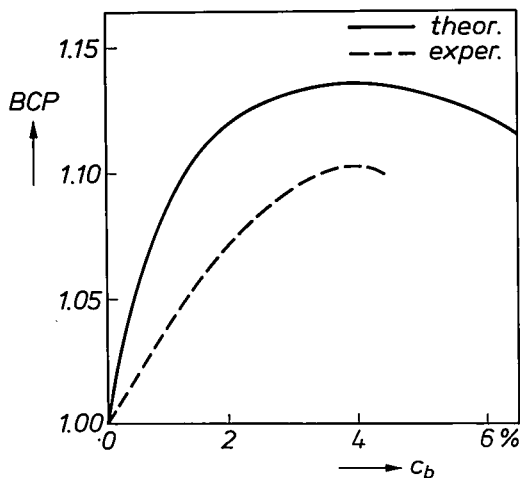


Fig. 9. *BCP*-values for a mixture of the blue-emitting phosphor $ZnS:Ag$ and cobalt blue as a function of the pigment concentration c_b . Here again the experimental curve (*dashed line*) is lower than the calculated curve (*solid line*).

filter layer and the aluminium layer, we can make the same type of calculations as for a filter sandwich.

According to the Kubelka-Munk theory for weakly absorbing scattering layers^[6], the transmittance and reflectance of this type of filter can be expressed in terms of the absorption coefficient K and the scattering coefficient S . By measuring the reflectance of an 'infinitely thick' layer of powder, we obtain the ratio K/S . The value of S can be derived, e.g. by measuring the reflectance for various layer thicknesses. This value is practically independent of the wavelength but it does depend on the phosphor material: about $500 \text{ cm}^2/\text{g}$ for $ZnS:Ag$ and about $400 \text{ cm}^2/\text{g}$ for $Y_2O_2S:Eu$.

The measured and calculated *BCP*-values for the mixture of $Y_2O_2S:Eu$ and $\alpha-Fe_2O_3$ are plotted against the pigment concentration in *fig. 8*. The measured values are all lower than the calculated ones, but if the approximations are taken into account the agreement is acceptable. The same applies to the mixture of $ZnS:Ag$ and cobalt blue; see *fig. 9*.

[6] P. Kubelka and F. Munk, *Z. techn. Physik* 12, 593, 1931.

Summary. Colour filters in the form of inorganic pigments can considerably improve the contrast of a colour television receiver without impairing its brightness. The effect is easy to describe if the filter is included as a layer in a sandwich formed by the glass of the screen and the light-emitting phosphor layer. The calculated improvement and the optimum surface concentrations for the red and blue pigments — the use of green pigment yields a loss — give good agreement with the experimental results. Mixtures of phosphor and pigment improve the contrast less than filter sandwiches, but are to be preferred for practical application because they are easier to make. Their action is more difficult to describe, but calculations based on a simple model still give reasonably good agreement with practice.

P²CCD in 60 MHz oscilloscope with digital image storage

H. Dollekamp, L. J. M. Esser and H. de Jong

'Time-axis conversion' is mentioned as an application in the title of one of the first publications about analog shift registers [1]. An example is given: in the PM 3310 storage oscilloscope signal samples recorded very rapidly are stored temporarily in an analog shift register — a 'profiled peristaltic charge-coupled device' (P²CCD) — and later read out more slowly to permit conversion to digital form at a lower rate. The subject is an example of the successful cooperation between Philips Research Laboratories, where P²CCD originated, and the Philips Scientific and Industrial Equipment Division [].*

The role of the P²CCD in the PM 3310 storage oscilloscope

In using an oscilloscope, it is sometimes desirable to be able to retain the picture after the input signal is no longer available. One example is in the observation of fast non-recurring effects. Until now this has been done by using the storage tube — a cathode-ray tube in which the electron beam records the measured effect in a pattern of electrical charges on a dielectric. This pattern is displayed on the screen by electron-optical methods and slowly fades away.

Digital solid-state memories — integrated circuits in which large numbers of binary numbers can be stored — have been commercially available for some years, so it would seem that these could replace the expensive storage tube. However, this requires the input signal to the oscilloscope to be digitized. This is a considerable undertaking in a broadband instrument. We shall describe an oscilloscope in which the signal is sampled up to 50 million times a second. The amplitude of each sample is recorded to an accuracy of 1 in 256, producing an information flow of 400 million bits per second. Circuits capable of analog-to-digital conversion at this rate without introducing large errors would be extremely expensive — if indeed they exist. This restricts the application of digital solid-state memories in broadband oscilloscopes.

In our oscilloscope, the PM 3310, which is shown in *fig. 1*, this difficulty has been overcome by using a new element known to electronic engineers as 'P²CCD'.

This is a shift register that stores signal samples in analog form and is fast enough for a sampling rate of 50 MHz [2]. The samples are then read out at a much lower rate — about 78 kHz — and converted individually to a digital code.

The P²CCD ('profiled peristaltic charge-coupled device') is a recent development of the analog charge-coupled device (CCD). In the analog CCD the individual elements are not restricted to one of the two conditions '0' or '1' as in the digital shift register, but can contain a continuously variable electric charge, whose magnitude is a measure of the magnitude of a sample taken from a signal at a particular moment. One element of an analog shift register or charge-coupled device can thus contain as much information as N elements of a digital shift register if the quantity being represented is expressed in numbers of N bits.

The analog charge-coupled device began its career as the 'bucket-brigade line' at Philips Research Laboratories [1]. The name 'bucket-brigade line' was later reserved for the version in which the 'packets' of charge are stored in separate diffusions in a single-crystal semiconductor. In the 'charge-coupled de-

[*] The Solid-State Special Products Group of the Electronic Components and Materials Division also put in considerable effort in making the production of the P²CCD possible. We are particularly grateful to Ir D. Daub and Ing. G. Gruitjes of the Group.

[1] 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.

[2] L. J. M. Esser and F. L. J. Sangster, Charge transfer devices, in: T. S. Moss (ed.), Handbook on semiconductors, Vol. 4, Ch. 3B; North-Holland, Amsterdam 1981.

Ir H. Dollekamp and Ing. H. de Jong are with the Philips Scientific and Industrial Equipment Division at Enschede, and Ir L. J. M. Esser is with Philips Research Laboratories, Eindhoven.

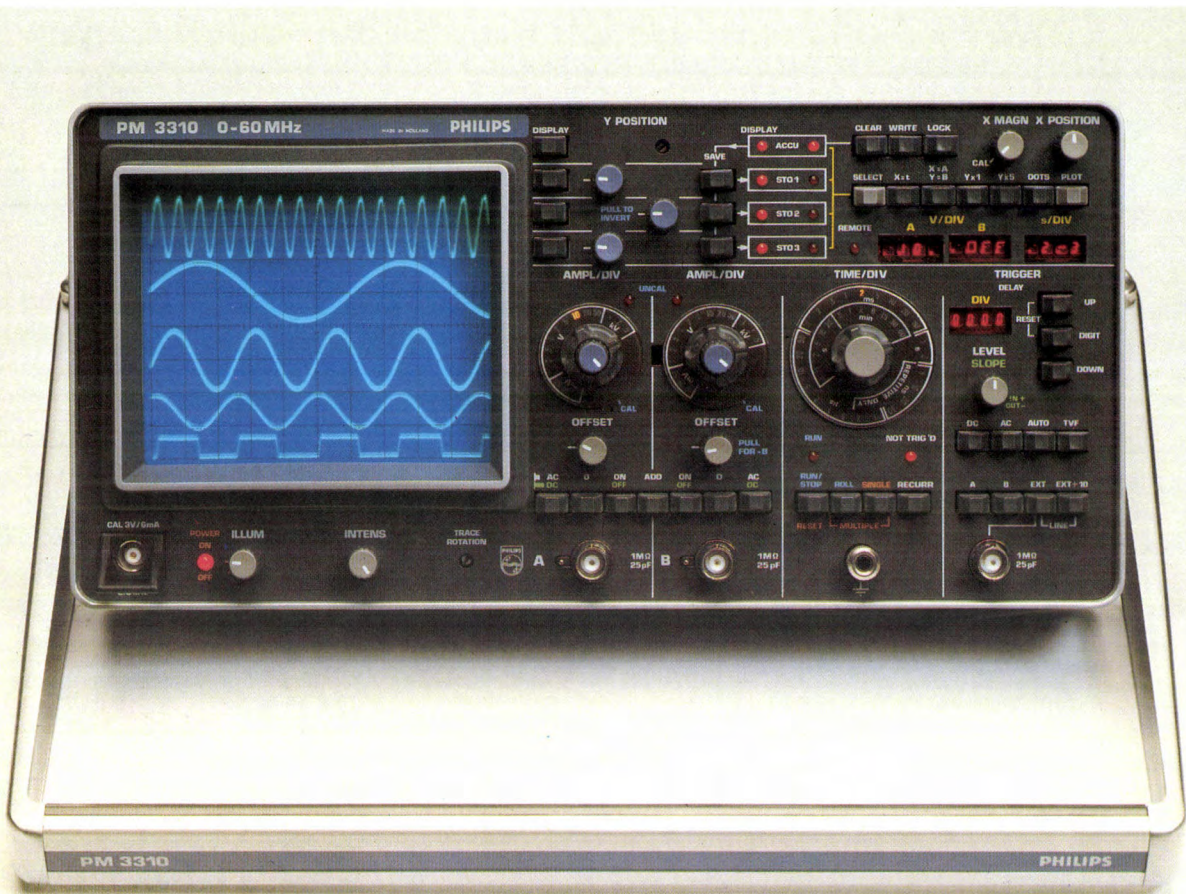


Fig. 1. The Philips PM 3310 two-channel storage oscilloscope. The contents of four semiconductor memories (*ACCU* and *STO 1-3*) can be displayed simultaneously, each with two traces if desired. Highest sampling rate of the input signals: 50 MHz.

vice', however, the packets of charge move through a homogeneously doped zone in the semiconductor.

Before examining the special features of the P²CCD we shall first consider the circuit of the PM 3310 oscilloscope, to show how the P²CCD is used.

The PM 3310, a digital storage oscilloscope

General

The PM 3310 digital storage oscilloscope is a portable two-channel oscilloscope for frequencies up to 60 MHz. It is not only the memory that is digital — all of the signal processing and display is controlled by a microprocessor, which is a vital part of the instrument. The images are displayed by read-out from the memory, so that even extremely fast 'single-shot' effects give a trace of the same brightness as with slower effects.

There are four digital memories. The events recorded in these four memories, each with its own time scale,

can be displayed simultaneously. If both channels are used, each memory is divided into two; eight simultaneous displays are then possible (*fig. 2*).

The contents of each of the four memories have their own delays relative to the trigger signal. These delays can be set anywhere between -9 and $+9999$ divisions with a digital delay circuit [3]. It is therefore possible to look at any of 1000 frames after the moment of triggering, and also at the last frame before triggering. For example, it is possible to see what happens before a fuse blows. *Fig. 3* shows another example. The image here has been triggered by a field pulse from a television signal. By selecting the appropriate delay any desired line in the field (i.e. raster) can be displayed. The bottom trace shows a number of the first 25 lines (which occur outside the television picture). The trace above it shows one of the lines carrying video information. The trace above that again shows the colour burst that is transmitted with each line for synchronization; the top trace shows this

again, but magnified. The sensitivities (volts/division) and the time scales (seconds/division), which are different for each trace, are held in the memories and can be read from an alphanumeric display panel formed from LEDs.

Another interesting possibility is the roll mode. In this mode the pattern displayed runs from right to left over the screen in four successive lines; the new signal appears at the top right while the oldest signal disappears from the screen at bottom left. With the time base at the slowest setting the screen can show an oscillogram corresponding to a total time of 40 hours.

When an extra module is added (the PM 3325 IEC bus interface), the PM 3310 oscilloscope can interact with other instruments via a standardized digital channel (the IEC bus^[4]). The control settings can be communicated to these instruments and modified by them; the contents of the memories (both data and settings) can also be communicated. The data, after having been processed by a computer, may be sent back to the oscilloscope in a suitable format for display on the screen, so that the oscilloscope acts as a display unit for the computer.

Signal housekeeping

One way in which the PM 3310 digital oscilloscope is essentially different from a conventional oscilloscope, which only works with analog signals, is in its 'signal housekeeping' — the management of the information-bearing and control signals circulating in the instrument, and the time relations between them. To illustrate this point we shall use the block diagram in *fig. 4* to follow the progress of a signal from the input of the oscilloscope until it appears on the screen.

The oscilloscope has two inputs. When two channels are used the two signals are processed in time-division multiplex. The setting of the input attenuators and amplifiers depends on the sensitivity selected on the front panel (*AMPL/DIV*). The signal is then input to the P²CCD, where it is sampled at regular intervals. The P²CCD consists of two parts, which take samples alternately; 128 analog signal values can be stored in each part.

The sampling rate depends on the finest detail in the signal, i.e. on the highest frequency in it that is to be resolved on the screen. The *TIME/DIV* control can be used to 'stretch' the picture to a greater or lesser extent in the horizontal direction. The screen always shows 256 successive samples; the more the picture is stretched, the closer the samples are to each other in time. The *TIME/DIV* control in fact adjusts the sampling rate, up to a maximum of 50 MHz. Here we see a basic difference from the conventional oscilloscope, in which the *TIME/DIV* control is used to

adjust the rate at which the electron beam passes over the screen in the horizontal direction.

What happens to the 256 analog signal samples stored in the P²CCD? As already stated, the read-out rate is invariable, at about 78 kHz. A different clock signal has therefore to be applied to the P²CCD for read-out. The P²CCD receives this signal, as it does the write-in clock signal, from the time-base generator (*TBU*) after an instruction from the signal-acquisition control logic (*ACL*), which only issues the instruction after it has received a trigger pulse from the trigger circuit *TRI*. This ensures that the processing of the 256 signal values starts at the right moment. The presence of 256 signal values at the moment when the trigger pulse arrives makes it possible for the oscilloscope to display one frame before the instant of triggering.

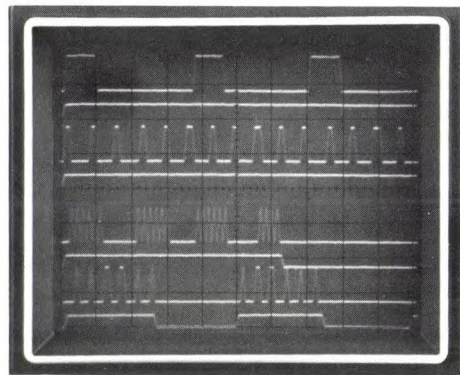


Fig. 2. The combination of two channels and four memories gives a maximum of eight traces on the screen.

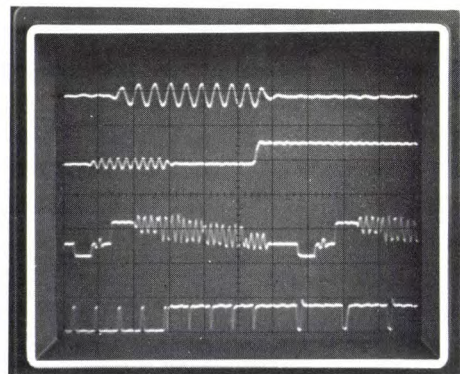


Fig. 3. Each memory has its contents displayed on its own time scale and with its own delay relative to the trigger signal, which is formed by the field pulse of a television signal. The bottom trace shows the TV lines that immediately follow this signal. By setting the delay appropriately each line can be selected for display; in the third trace this line is one of the lines carrying video information. By changing the time scale detailed display of the colour burst can be obtained on the first and second traces.

[3] By 'divisions' we are referring here to the coarse segmentation into squares.

[4] IEC Publication 625.

The first part of the processing is the conversion of every signal value to an eight-bit digital word. This is done in an analog-to-digital converter (*A/D*), which also receives the clock pulse necessary for this operation from the control logic. When the first conversion has been completed, the result is transferred via eight parallel lines to a buffer memory (*ShRe*) and stored there temporarily. The control logic then generates the next clock pulse, which triggers the next conver-

the display. This means that the 256 eight-bit words in the buffer memory are not written to *ACCU* until it has given a signal to indicate that it is ready to receive the information.

The information stored in *ACCU* can be copied into any of the other three memories: *STO1*, *STO2* and *STO3*. Each memory can accommodate 256 eight-bit words. Each word is stored together with its address, which determines the position on the horizontal

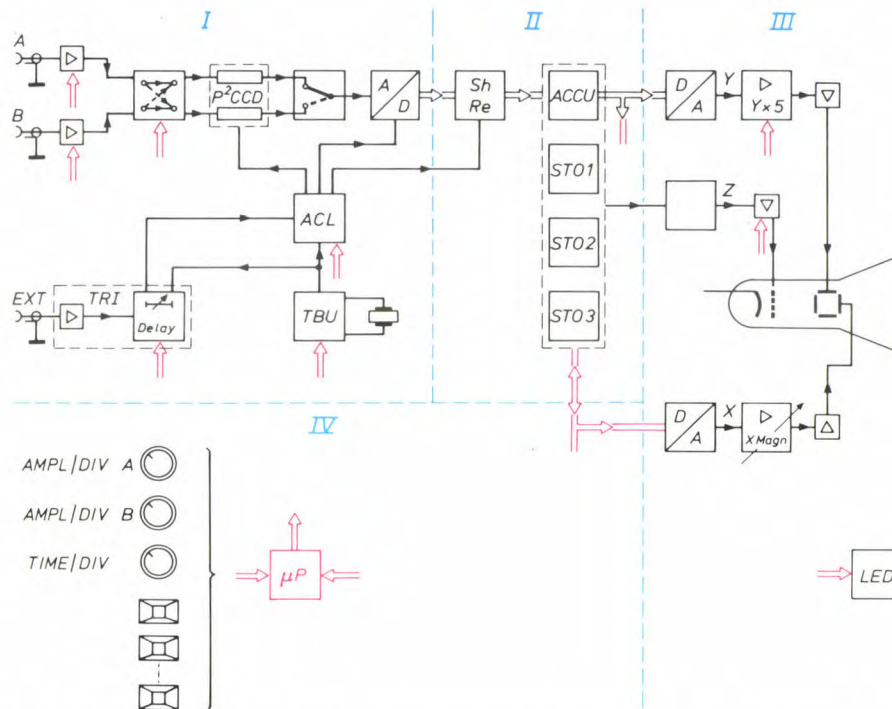


Fig. 4. Block diagram of the PM 3310 two-channel storage oscilloscope. *I*: signal processing. *A*, *B* signal inputs. *EXT* input for external trigger signal. *P²CCD* analog shift-register for 2×128 signal samples. *A/D* analog-to-digital converter. *ACL* control logic for signal processing. *TBG* time-base generator. *TRI* trigger unit. *Delay* variable delay line. *II*: store. *ShRe* shift register acting as a buffer memory. *ACCU* memory (256×8 bits) that receives the data. *STO 1*, *2* and *3* memories in which the data can be stored. *III*: display. *D/A* digital-to-analog converters for generating the *X* and *Y* deflection voltages. *Z* beam suppression. *LED* alphanumeric indication of the control settings with light-emitting diodes. *IV*: control. μP microprocessor that periodically checks the position of the controls and generates control signals.

sion. After 256 clock pulses the shift register is full and no more clock pulses are produced.

While the *P²CCD* now starts to write in analog signal samples again, a start is made with the transfer of the contents of the buffer memory to a memory that holds the information to be displayed. This is one of the four main memories; it is shown as *ACCU* in fig. 4. The buffer memory has to be brought into the procedure because signal acquisition and signal storage/display each have their own cycle and are not synchronous. The display cycle is based on a frame repetition rate of 50 Hz.

The transfer of the information from the buffer memory to *ACCU* is performed via a 'handshaking' procedure, so as not to interfere with the regularity of

X-axis where it must be displayed. The horizontal axis, which is ten divisions long, is divided into 256 discrete positions. The horizontal deflection voltage for the cathode-ray tube is generated by digital-to-analog conversion of the addresses.

At each horizontal position the contents of the associated eight-bit word determine the deflection of the luminous spot above or below the zero line, i.e. in the *Y*-direction. The eight bits permit 256 discrete *Y*-values to be distinguished. The vertical deflection voltage for the cathode-ray tube is generated by digital-to-analog conversion of the eight-bit words. The complete picture is thus formed in a matrix of 256×256 points. This matrix has a height of two divisions on the screen, or ten divisions in the *Y* \times 5 set-

ting (the complete screen is eight divisions high). If the user wishes, he can connect the points with straight-line sections to form a continuous curve (*fig. 5*).

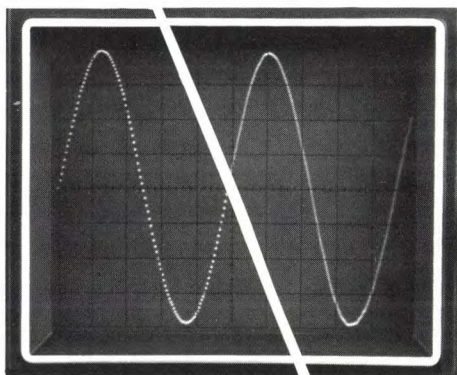


Fig. 5. Signal displayed at the setting $Y \times 5$. *Left:* the individual samples. *Right:* the smooth curve obtained by linear interpolation.

The P²CCD

CCD

In its simplest form a homogeneous shift register for analog signals (a charge-coupled device, or CCD) consists of a single-crystal semiconductor with a row of electrodes mounted on it; the electrodes are insulated from the semiconductor (*fig. 6*). If the semiconductor is silicon — the usual material — the insulation is provided by a thin film of SiO₂.

Under the electrodes there is a series of charge ‘packets’ that form a representation of the sampled signal. The charge can be injected electrically but may also be produced by incident light. If the charge is optically generated the CCD acts as a solid-state image sensor — an important application ^{[1][2][5]}, but we shall not discuss it further here. In our case the

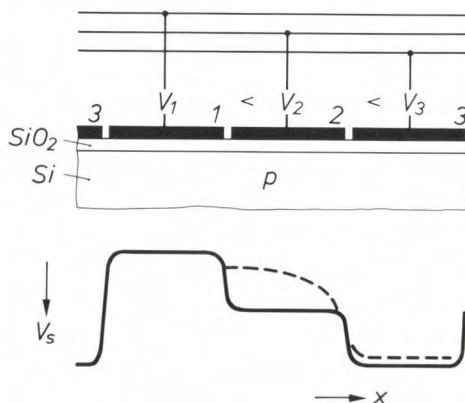


Fig. 6. Homogeneous shift register (CCD) in its simplest form. The single-crystal silicon carries a row of electrodes insulated from it by a thin layer of silicon dioxide. The electrodes are connected in turn to clock potentials V_1 , V_2 and V_3 . The potential V_s at the Si-SiO₂ interface is highest beneath electrode 3, since potential V_3 is the highest (*solid curve*). A potential well is formed here, in which free electrons collect. V_2 was previously the highest; the dashed line illustrates how the charge travels from electrode 2 to electrode 3. Electrode 1 separates the charge packets from each other.

charges are created electrically and it is the CCD’s task to store and transfer the charges and keep them separated. This requires at least three sets of electrodes. Three successive electrodes, one from each set, together form a memory element; an example is shown in *fig. 6*. Periodically varying clock voltages are applied to the electrodes; the clock signal has three phases, one for each set of electrodes.

This signal in the n-channel CCD changes the electrodes to a higher or lower positive potential in relation to the substrate. The free electrons in the silicon will travel to the electrode at the highest potential. If we assume that the silicon is of type p, the electrons in it will be minority carriers. They travel along the interface with the SiO₂, where inversion of the p-type silicon produces an n-type channel as in an n-type MOS transistor. In *fig. 6* electrode 3, which has the storage function, has the highest potential; the solid curve indicates the variation (in the absence of free electrons) of the potential at the Si-SiO₂ interface. Since the potential is plotted in the downward direction, the resultant picture is that of a potential well; the electrons flow into this well since they have their lowest potential energy there.

Shortly before, electrode 2 had the highest potential. Its potential, however, has been lowered, so that electrode 2 has acquired a transfer function. At the same time the potential of electrode 3 was increased, so that the situation shown in *fig. 6* was reached. The electrons collected beneath electrode 2 travel towards electrode 3, but they take some time to do so. The dashed curve in *fig. 6* is an instantaneous picture of the potential distribution at the Si-SiO₂ boundary during this transition; it is as if the charge from the well beneath electrode 2 is ‘transferred’ into the well under electrode 3. The potential under electrode 1 is even lower, preventing the charge from flowing away to the left. Electrode 1 has a separating function and ensures that there is a transfer of separate packets in the correct direction. The potential distribution over the successive electrodes moves one step at a time to the right, producing a kind of travelling potential wave that carries the charge packets separately. Transfer across a single element corresponds to three charge movements and takes one clock period.

Transfer efficiency, highest frequency

Not all of the charge is transferred to the next electrode. Perturbations of the crystal lattice occur at the Si-SiO₂ boundary, leading to local energy minima. These are generally referred to as ‘surface states’, and

^[5] A linear image sensor with 1728 elements, the P²CCD 1728 made by Valvo GmbH, is used in the Megadoc system for reading documents (Philips tech. Rev. 39, 329-343, 1980).

they can cause charges to be trapped. Charges that are trapped from a packet may be emitted later and added to a following packet. This reduces the transfer efficiency to less than 1. The effect of surface states can be reduced by ensuring that the size of the charge packet does not fall below a certain minimum (bias charge, 'fat zero'); most surface states are then permanently occupied. Efficiencies of 0.9999 are reported for CCDs used in this way, which means that 10^{-4} of the charge is left behind at each transfer step. In the 500 cells of a three-phase shift register, for example, a total of about 15% of the information is lost.

The cause of this limitation of the transfer efficiency may be considered as external; it is not determined by the transfer mechanism itself and is not very dependent on the clock frequency. On the other hand, the highest frequency at which the shift register will work is determined by the actual transfer process. During transfer there is a gradient in the charge density. This gradient $\partial Q/\partial x$ (where Q is the surface charge density and x the direction of transfer; see fig. 6) is accompanied by an electric field, which causes the charge to move in the x -direction. For this self-induced field E , we have:

$$E = -\frac{1}{C} \frac{\partial Q}{\partial x}, \quad (1)$$

where C is the oxide capacitance per unit area. The rate of charge transfer (the current density J) is given by

$$\frac{\partial Q}{\partial t} = -\frac{\partial J}{\partial x} = -\frac{\partial(\mu QE)}{\partial x} + D \frac{\partial^2 Q}{\partial x^2}, \quad (2)$$

where μ is the mobility of the electrons and D the thermal-diffusion coefficient. The rate decreases as the charge gradient $\partial Q/\partial x$ and the field it induces decrease; the factors Q and E in (2) then both become smaller. The charge finally remaining is transferred more by thermal diffusion than by the self-induced field. Thermal diffusion is a slow process; when the clock voltages change some residual charge will remain. If the initial charge is large, the self-induced field will also be strong and a relatively larger portion of the charge will be transferred as a result of this field than with a low initial charge. With a large initial charge the residual charge at the end of the clock cycle is larger in absolute value than with the small starting charge, but it is not larger relatively. The transfer efficiency, which is a relative measure, is consequently better for larger charges; to take advantage of this effect a constant amount of charge is added to all the charge packets. We thus have an intrinsic reason for a preset or 'bias' charge; the surface states were more of an extrinsic reason. The bias charge, however, limits

the signal-amplitude range and must become larger as the clock frequency increases.

The transfer efficiency is therefore limited for both extrinsic and intrinsic reasons. The highest frequency is determined by the intrinsic transfer mechanism; for the structure described, with conduction in an n-channel at the surface, the frequency range is about 10 MHz. The charge storage is

$$Q \approx C V_{cl}, \quad (3)$$

where V_{cl} is the clock-voltage swing. This charge storage is the maximum that is physically possible and is larger than with other existing structures. It is not equal to the maximum dynamic range; this is also affected by the noise level, which depends considerably on the activity of surface states. The signal-amplitude range is determined by the stored charge minus the bias charge, which is about 20%.

The CCD described here is known as a surface charge-coupled device, or SCCD for short.

PCCD

The charge transfer to the next electrode, especially for the residual charge, proceeds faster if the self-induced field is not the only electric field causing the charge carriers to move. Now the field of the electrode with the highest potential — towards which the electrons are moving — spreads out sideways at a greater depth below the surface. When the self-induced field has become considerably weaker, the last part of the charge packet therefore has to be moved to a greater depth, so that it comes under the influence of this electrode field. In the peristaltic CCD (PCCD) this is achieved by depositing a layer of n-type silicon on the substrate of p-type silicon; the thickness of this n-type layer should be, say, half an electrode length (fig. 7). The n-type layer is biased positively and is thus depleted of free electrons. The ionized donors, which remain behind, form a distributed positive space charge, which is fixed in position.

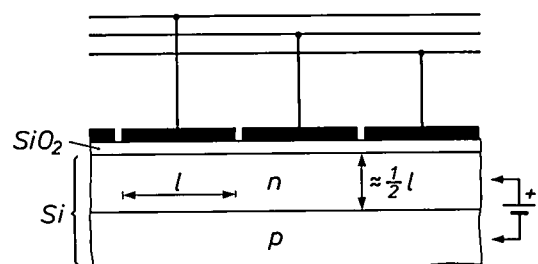


Fig. 7. Peristaltic CCD (PCCD). There is an n-type layer on the p-type silicon. Small packets of electrons collect in the bulk of this layer, while larger packets spread out from the bulk to the surface. The electric-field configuration in the bulk of the layer is more suitable for rapid charge transfer, particularly if the layer is about half the thickness of an electrode length l .

Because of this space charge the potential in the layer is higher than at the surface; see the solid curves in *fig. 8*, which represent the change of potential in the y -direction perpendicular to the surface [6]. The curves indicate potential wells some distance from the surface and it is in these regions that free electrons will initially collect. For example, there will be a small charge between the depths d' and d'' below the electrode connected to +15 V (see *fig. 8*). If the free charge increases, the charge packet will extend to the surface; the electric field in the SiO_2 layer then reverses its direction.

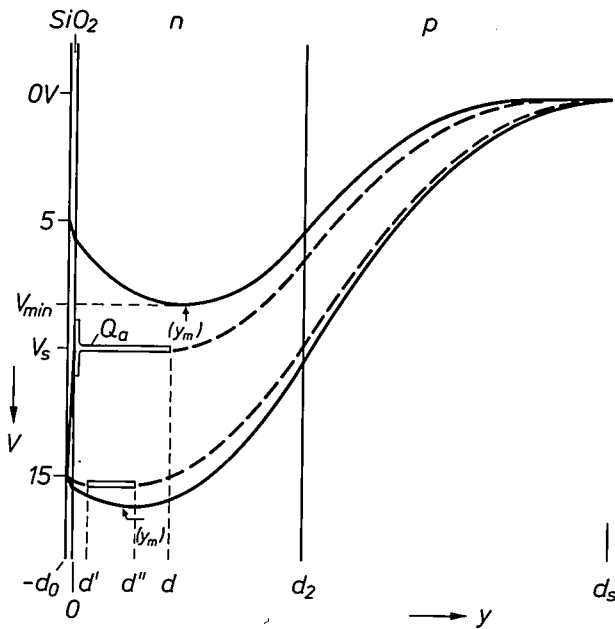


Fig. 8. The behaviour of potential V in the bulk of the PCCD as a function of the depth y . *Solid curves:* electrode potentials 15 V and 5 V, no free charge. The maximum of the potential is some distance from the surface. *Dashed curves:* electrode potential 15 V and small charge packet (between d' and d'') or a larger charge packet (between 0 and d). With the larger charge packet most of the charge collects at the surface (Q_a).

In the first case, when the charge is stored entirely in the bulk of the semiconductor, the charge per unit area is

$$Q_b = e N_2 (d'' - d'), \tag{4}$$

where e is the electronic charge and N_2 is the doping of the n-type layer. If the charge reaches the oxide layer the oxide capacitance receives a charge

$$Q_a = C_{ox} (V_s - 15), \tag{5}$$

where C_{ox} is the capacitance per unit area of the oxide film and V_s is the potential of the semiconductor at the interface with the oxide; it is assumed that the electrode potential is 15 V and hence that the voltage

[6] The method of calculating the potential curves is given in the Appendix.

across the oxide is $V_s - 15$. The charge in the bulk of the semiconductor is now

$$Q_b = e N_2 d, \tag{6}$$

which extends to a depth d .

The charge on the oxide is usually much larger than the charge in the bulk, which is given by (6); if N_2 is made small, e.g. $7 \times 10^{14} \text{ cm}^{-3}$, then in this example 90% of the charge will be on the oxide.

This means that the peristaltic CCD has a charge-storage capacity almost as large as that of the surface CCD. If the charge is considered to be concentrated at a 'centre of gravity', this centre of gravity will no longer be at the surface, but a small distance y_{cg} away from it. The storage capacitance C_{st} can now be regarded as the resultant of the oxide capacitance C_{ox} in series with the capacitance of the layer of thickness y_{cg} :

$$C_{st} = (d_0/\epsilon_{ox} + y_{cg}/\epsilon_{si})^{-1}, \tag{7}$$

where d_0 is the thickness, ϵ_{ox} the permittivity of the oxide and ϵ_{si} the permittivity of the silicon. In the case illustrated in *fig. 8*, y_{cg} is $0.15 \mu\text{m}$. If the oxide thickness d_0 is $0.1 \mu\text{m}$, the storage capacitance C_{st} is $0.7 C_{ox}$. The maximum charge storage is reached when $V_s = V_{min}$; when $C_{st} = 0.75 C_{ox}$. In this situation the separation function of the 5 V electrode is only just maintained.

The PCCD thus preserves the advantage of the surface CCD, high charge storage. It has however the added advantage of very rapid transfer of both large and small charge packets. Small packets are transferred in the bulk of the semiconductor and not at the surface.

At the start of transfer there is a combination of surface and bulk transfer. The transfer is then rapid because the self-induced fields are strong. When the packet has been reduced to some ten per cent of its initial size, it descends more deeply into the bulk of the semiconductor. The reduction of the capacitance between the remaining charge and the electrode has the effect of increasing the remaining self-induced field (eq. 1). In addition, the packet then comes under the influence of fields that are generated by the external potentials on adjacent electrodes; because of this the transfer remains fast to the very last electron. The PCCD thus exploits both the self-induced and the externally induced fields.

Charge transfer in the PCCD

Fig. 9 shows the cross-section of a four-phase PCCD. Four-phase structures are widely used in practice because they require only two layers of wiring, whereas three-phase structures need three layers. The

curve *S-S* indicates, for the particular electrode potentials, the variation of the potential at the interface between the semiconductor and oxide; curve *B-B* indicates the variation of the potential at a depth of about half an electrode length in the bulk. From curve *B-B*, which is a 'smoothed-out' version of curve *S-S*, we see that the field component E_x in the direction of transfer, which is given by

$$E_x = - \frac{\partial V}{\partial x}, \quad (8)$$

is not zero directly underneath the transfer electrode ($x = 0$). Any charge at that location in the bulk will therefore be transferred, unlike any charge at the surface.

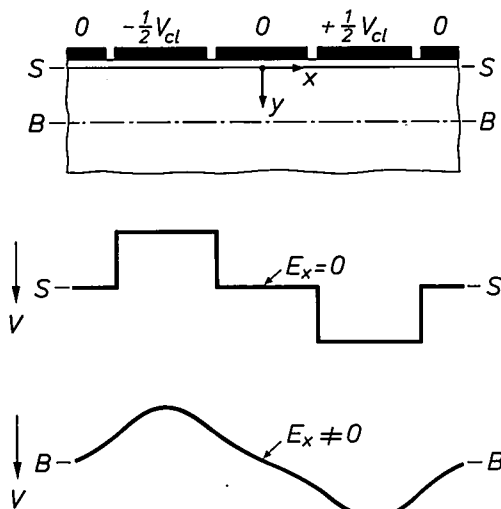


Fig. 9. Potential V in a four-phase PCCD at the interface with the oxide (section *S-S*) and at some depth in the semiconductor (section *B-B*). For the section *B-B* the steps in potential are 'spread out'; the electric field in the direction of transfer, E_x , has spread out sideways and removes the charges from underneath the neutral electrode, which is at a potential of 0 V. This is particularly important for the transfer of the last traces of charge and gives good transfer efficiency even at high clock frequencies.

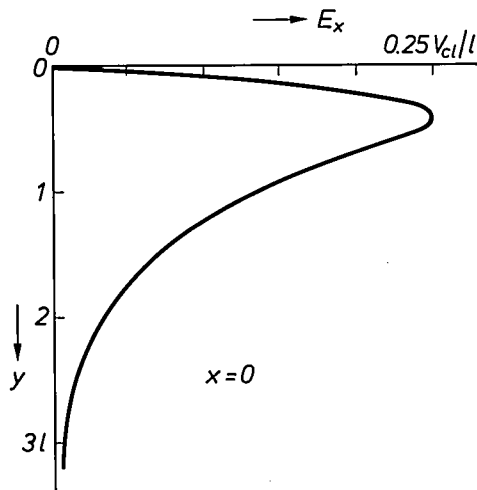


Fig. 10. The strength of the x -component E_x of the electric field as a function of the depth y in the semiconductor. The curve applies for the central position under the transfer electrode ($x = 0$).

The behaviour of the field-strength E_x as a function of the depth y is indicated in *fig. 10*^[7]. The figure shows the variation of this component below the centre of the transfer electrode ($x = 0$); the maximum is at a depth of about half an electrode length. The magnitude of the x -component there is about $\hat{E}_x = V_{cl}/4l$, where V_{cl} is the clock-voltage swing and l the electrode length. A minimum time for the transfer of the charge to the next memory element can be calculated from this information; in a four-phase structure this time is

$$T_{min} = 4l/\mu\hat{E}_x \approx 0.7 \text{ ns}, \quad (9)$$

where $l = 7.5 \mu\text{m}$, $V_{cl} = 10 \text{ V}$ and the mobility $\mu = 1300 \text{ cm}^2/\text{Vs}$. This relation corresponds to a maximum clock frequency of 1.5 GHz. The transfer time applies only to small charge packets that do not perturb the field to any great extent. They are transferred in the bulk of the semiconductor and not at the surface.

Since for the larger packets of charge the charge travels partly via the surface channel and partly via a channel in the bulk of the semiconductor, the term

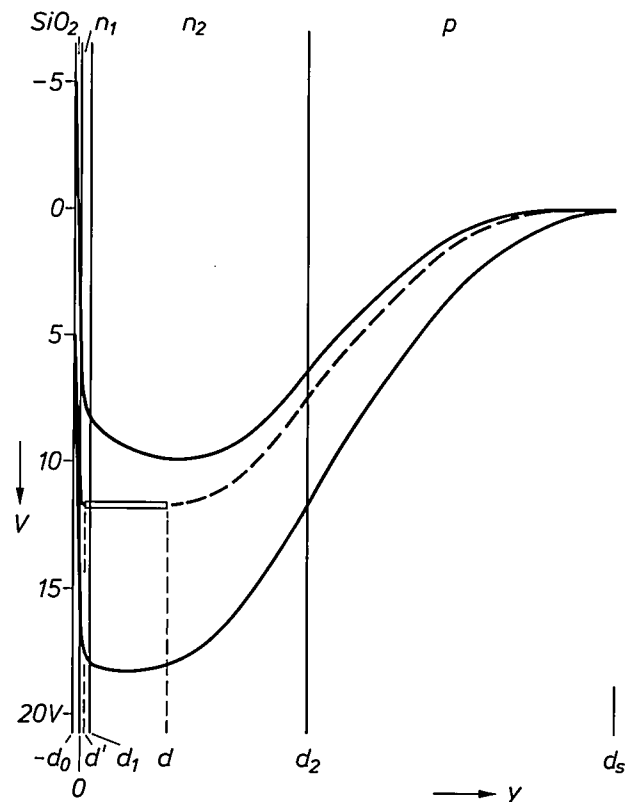


Fig. 11. The variation of the potential V in the bulk of the P²CCD as a function of the depth y . The P²CCD has a heavily doped top layer n_1 (about 10^{17} cm^{-3}) of thickness d_1 (about $0.2 \mu\text{m}$). *Solid curves*: electrode potentials of +5 V and -5 V; no free charge. *Dashed curve*: electrode potential +5 V, charge packet from $y = d'$ to $y = d$. Because of the high doping of n_1 , d' can be very small while the potential difference with respect to the surface is still large enough to prevent interaction between the charge and the surface states.

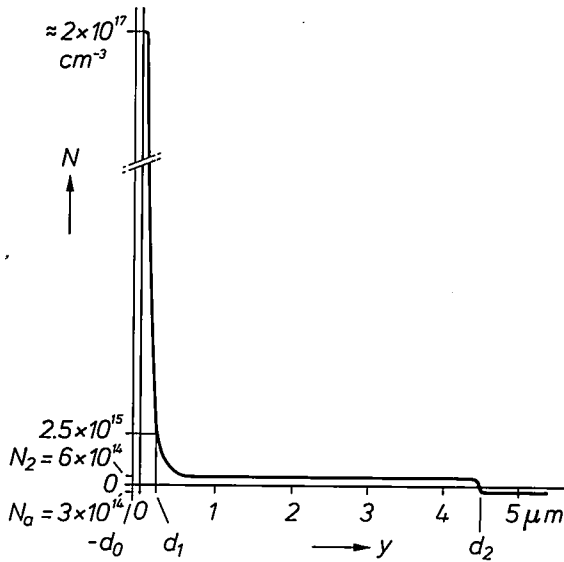


Fig. 12. The doping N of the various layers of the P^2CCD . The 'centre of gravity' of the highly doped layer n_1 is at a depth of less than $0.05 \mu m$.

BCCD

We should just briefly mention the 'buried-channel charge-coupled device' or BCCD. This has the same structure as the PCCD but a more heavily doped active layer (2×10^{15} - $5 \times 10^{16} \text{ cm}^{-3}$ instead of $5 \times 10^{14} \text{ cm}^{-3}$), so that the charge never reaches the surface. Since the transfer takes place in the bulk and there is no interaction with the surface states the BCCD has an excellent transfer efficiency (e.g. 0.99999), but it does not combine high speed with high charge storage.

P^2CCD

To overcome the restriction imposed on transfer efficiency by the surface states, we have modified the structure of the PCCD to form the profiled peristaltic CCD or P^2CCD . The modification is the addition of a heavily doped n-type top layer (about 10^{17} cm^{-3}), which is less than $0.2 \mu m$ thick and has been produced by low-energy implantation of 2×10^{12} arsenic atoms per cm^2 . The large number of ionized donors in this

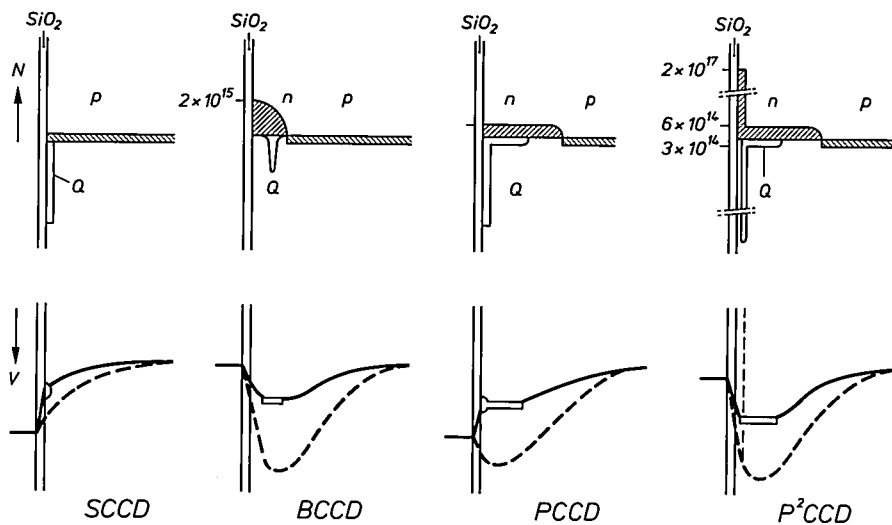


Fig. 13. Doping profile (N) and potential variation (V) in the four types of CCD discussed. n n-type silicon. p p-type silicon. The charge Q is entirely at the surface in the SCCD, and entirely in the bulk of the semiconductor in the BCCD; in the PCCD and P^2CCD there is a small charge in the bulk but the greater part of a larger charge is at the surface or in the surface layer respectively.

'twin channel' is applied to the PCCD. The contraction and expansion of the charge in its forward movement is reminiscent of the peristaltic movement of the alimentary tract, hence the name for this form of CCD.

Because of the twin channels the combination of a considerable charge storage and very fast and complete transfer is an intrinsic characteristic of the PCCD. The only threat to the transfer efficiency is presented, as in the surface CCD, by the surface states. We have observed transfer efficiencies of about 0.999.

top layer ensures that the charge that starts to collect on the top layer as the charge packet increases never reaches the interface with the oxide in practice; the surface states are not therefore exposed to any free charge. The potential curves are extremely steep in the heavily doped top layer (*fig. 11*); the potential step over the remainder of the top layer, which is not saturated with charge, need only be several tenths of a volt to make the interaction with the surface states negligibly small. In *fig. 11* this potential step (in the range $0 < y < d'$) is 1 V; for the potentials shown here 90% of the charge is in the top layer. The associated doping levels are shown in *fig. 12*. For easier comparison, *fig. 13* shows the doping profiles and poten-

[7] M. G. Collet and A. C. Vliegthart, Calculations on potential and charge distributions in the peristaltic charge-coupled device, Philips Res. Repts 29, 25-44, 1974.

tial distributions of all four types of CCD discussed above. The potential distribution in the P²CCD at successive stages of the charge-transfer process is illustrated in *fig. 14*.

In the P²CCD the charge transfer is also partly at the surface and partly in the bulk, and it can be said to have twin channels. Since there is no interaction with the surface states the transfer efficiency is 0.999999 without using a bias charge: for a clock-voltage swing of 10 V the charge storage is 1.2×10^{12} electrons per cm^2 . This means that the effective capacitance to the

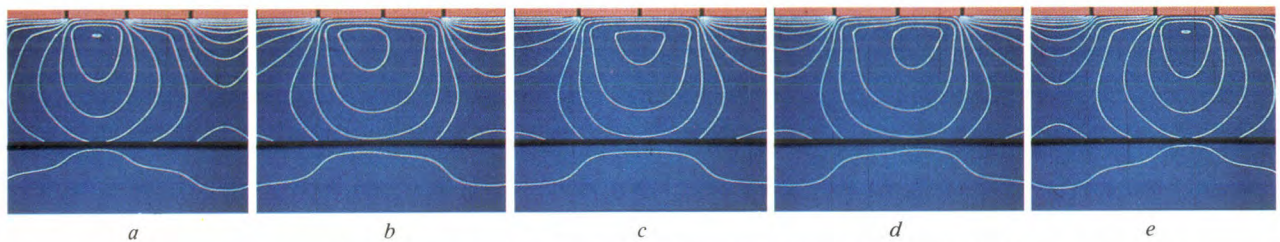


Fig. 14. Calculated potential distribution in the P²CCD at five successive times. The upper parts of the figures show equipotential lines beneath four successive electrodes (it is assumed that there are no free charge carriers); the lower parts show plots of the potential at a depth of half an electrode length [8]. *a*) Potentials of the four

Input circuit

At the start of the row of electrodes in the P²CCD an n⁺-type diffusion (*S* in *fig. 15a*) provides an ohmic contact to the semiconductor. The next electrode *1* is kept at a constant potential of about +6 V with respect to the substrate and acts as a barrier separating the successive charge packets from each other. The analog input signal, which has a maximum voltage swing of 1 V, is connected to electrode 2; a constant bias of about +7.5 V is added to the swing; the voltage at this electrode thus varies with the signal

electrodes: 0 V, +5 V, 0 V, -5 V. There is a potential maximum beneath the second electrode; free charge will begin to collect there. *b*) Intermediate stage. *c*) Potentials: -3.5 V, +3.5 V, +3.5 V, -3.5 V. *d*) Intermediate stage. *e*) Potentials: -5 V, 0 V, +5 V, 0 V. The distribution in (*a*) has progressed by one electrode.

electrode is $0.6 C_{\text{ox}}$. With this high transfer efficiency there is very little perturbation of the transfer, and since such perturbation is of a statistical nature and leads to transfer noise, the noise level is low and the dynamic range is large.

Practice has shown that the highest clock frequency that can be used is above 180 MHz for a four-phase structure; this means that there can be more than 720×10^6 transfers per second. This is not the absolute limit; measurements at higher frequencies are difficult, however, because of the capacitances between the electrodes, which partly overlap one another.

Input and output of the P²CCD

The series of travelling charge packets in the P²CCD must be a faithful representation of the sampled signal. It requires great care to ensure that a true representation is obtained when the charges are injected into the P²CCD; in the spatial separation during the formation of the charge packets precautions must be taken to ensure that the residual charges between them do not introduce misleading results. Distortion can also occur during read-out. The input and output circuits must therefore be designed with great care. There are various possibilities [2]; we shall now consider the input and output circuits of the P²CCD that have been developed for the PM 3310 oscilloscope.

between about +7 and +8 V. The square-wave voltage P_2 , one of the four phases of the clock signal, is applied to electrode 3; all four of these phases vary between a low potential of +2 V and a high potential of +11.4 V. This same square-wave voltage is also applied to source electrode *S*, but via a direct-voltage generator ($V_p \approx 14$ V) that gives a small internal potential difference in the n-layer of about 0.5 V.

Figs 15*b* to *d* show the potential distribution in the input section of the P²CCD at various times in the cycle of the clock signal P_2 . Between t_0 and t_1 the potential of the source electrode is lower than the barrier potential; negative charge flows across the barrier to the point of highest potential beneath electrode 2. How much charge is stored there depends on the signal voltage at this electrode; the higher it is, the deeper the potential well. An excess of charge also collects beneath electrodes 1 and 2. At time t_1 the potential of the source electrode starts to increase; the excess charge flows back to the source electrode. Because of the polarization voltage V_p the potential beneath electrode 3 is about 0.5 V more negative, so that the excess charge flows back towards the left.

When the potential of the source electrode reaches the barrier potential of electrode 1, flow-back ceases (*fig. 15c*). The potential of the source electrode continues to increase but this has no further effect because of the existence of the barrier. The potential of electrode 3 also increases simultaneously and the

well beneath it becomes deeper than the well beneath electrode 2. The charge beneath electrode 2 now moves to electrode 3 (fig. 15d).

The potential rise at time t_1 occurs rapidly. The charge packet that is separated by this increase is representative of the signal value at that moment. The input stages of the P²CCD thus carry out the sampling of the signal. The process in which an excess of charge is supplied and the surplus charge is returned is sometimes referred to as 'fill and spill' [9].

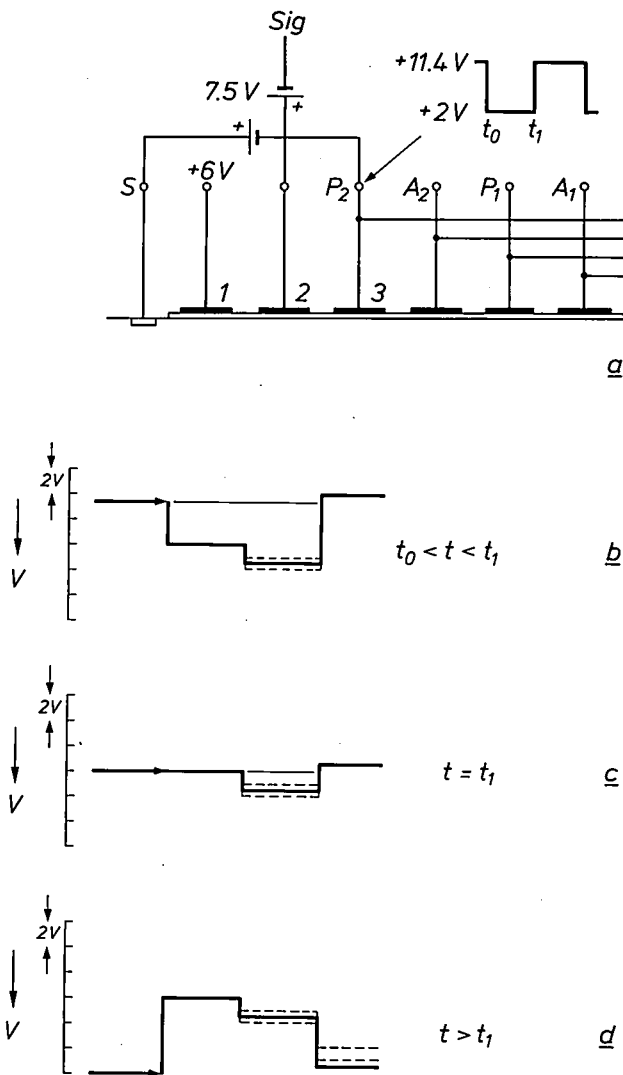


Fig. 15. Input of P²CCD. a) Circuit diagram. S source. 1, 2, 3 first three electrodes. Sig input for analog signal (maximum voltage swing 1 V). P₂, A₂, P₁, A₁ connections for four-phase clock signal. b) Potentials V in the bulk of the n-layer. P₂ is at +2 V, S at +16 V. So much charge flows in from S that the potential beneath electrodes 1 and 2 also drops to +2 V. c) The potential at P₂ and S increases rapidly. The charge flows back to the left to S until S has reached the same potential as that existing beneath electrode 1 (the time illustrated). The magnitude of the residual charge beneath electrode 2 depends on the instantaneous value of Sig. d) As the potential of P₂ continues to rise steeply, the residual charge flows into the potential well beneath electrode 3, so that a sample of the signal is recorded.

The return flow of the surplus charge takes place under the influence of the self-induced field. The last part of the charge therefore flows back slowly; it takes about 2 ns for the surplus charge to be reduced to 1%. The increase in the clock voltage proceeds faster; the situation shown in fig. 15c is reached about 0.5 ns after t₁. In fact, the situation at that moment is not so ideal as represented in fig. 15c: about 5% of the surplus charge is still left (fig. 16). Some of this flows into the charge well that is starting to form beneath electrode 3 and into which the signal charge also flows.

The magnitude of the surplus charge — and hence also of the fraction that is added to the signal charge — varies with the low level of the clock signal; a 1% variation of this level leads to a 10% variation in the signal charge. Stabilization of the low clock-signal level is therefore indicated.

Part B of the circuit in which the clock signal is generated is used for this purpose (see fig. 17). The clock signal is full-wave rectified and the resultant direct voltage is compared in a differential amplifier with a reference voltage that is held constant (+3.6 V). The difference signal drives the transistor T₁, which acts as a current generator for the push-pull stage T₂, T₄/T₃, T₆, which amplifies the clock signal.

Output circuit

At the end of the row of electrodes on the P²CCD another contact is formed by means of an n⁺-type diffusion; this acts as a drain and is connected to the gate of an MOS transistor integrated on the chip. This MOS transistor is connected as a source follower and gives a low output impedance (see fig. 18). After each charge packet the charge on the gate of this output transistor has to be equalized again. The equalization is carried out in another MOS transistor, which is also included on the chip; on receipt of a periodic external reset voltage it connects the drain to a reference voltage.

Zero correction

We mentioned earlier that light can be used to generate charge in a CCD; this facility is put to practical use when the CCD acts as an optical sensor. If



Fig. 16. Because of the slower rate at which the surplus charge flows back, the charge distribution at any particular instant is as shown here (and not as in the idealized representation in fig. 15c). The fraction of the charge left behind will flow to the right into the potential well that is forming and will introduce error by combining with the signal sample.

[8] H. W. Hanneman and L. J. M. Esser, Field and potential distributions in charge-transfer devices, Philips Res. Repts 30, 56-72, 1975.
 [9] J. E. Carnes, W. F. Kosonocky and P. A. Levine, Measurements of noise in charge-coupled devices, RCA Rev. 34, 553-565, 1973.
 M. F. Tompsett and E. J. Zimany, Jr., Use of charge-coupled devices for delaying analog signals, IEEE J. SC-8, 151-157, 1973.
 D. V. McCaughan and J. G. Harp, Phase-referred input: a simple new linear c.c.d. input method, Electronics Letters 12, 682-683, 1976.

the light is excluded, however, there is a small residual dark current, a leakage current caused by thermal excitation of pairs of charge carriers, e.g. at the site of crystal defects. Because of this leakage current a

It is obvious that the leakage current introduces an error into the picture, which is superimposed on a sloping base line. We correct this error by recording the zero level after each picture recording at the same

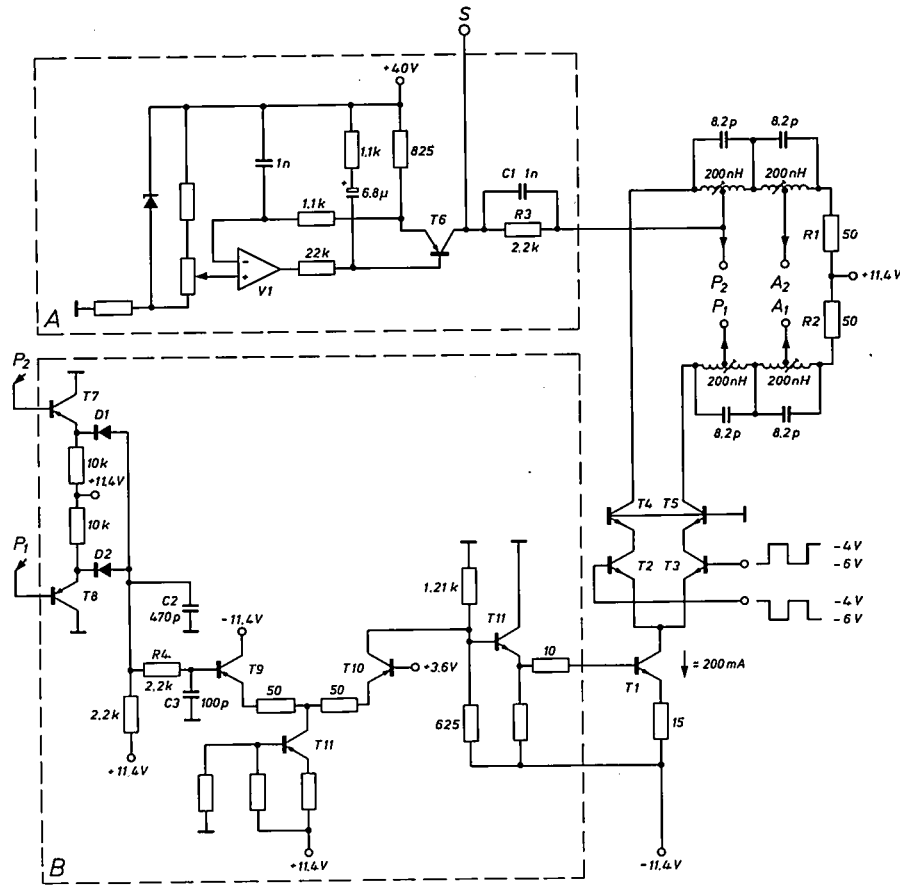


Fig. 17. Circuit for generating the clock signals P_1, A_1, P_2, A_2 . Part A is a stabilized current generator that generates, across resistor R_3 , the voltage to be applied between S and P_2 in fig. 15. Part B stabilizes the low level of the clock signal.

charge packet will become larger while it remains in the P²CCD. However, not all of the charge packets remain in the P²CCD for the same length of time. Each of the two parallel sections of our P²CCD has 128 elements and the n th signal sample remains present in the element for a time

$$\tau_n = \frac{n}{f_i} + \frac{128 - n}{f_o}$$

where f_i is the read-in frequency and f_o the read-out frequency; in our case f_o has a fixed value of about 78 kHz, while f_i varies through several powers of 10 depending on the oscilloscope setting. The charge Q_n added by the leakage current to the n th charge packet is plotted in fig. 19 against the number n ; this gives a pattern of sloping lines in which the slope depends on the ratio of f_i to f_o .

read-in frequency f_i and subtracting it from the picture content. This is done digitally: the buffer memory ($ShRe$ in fig. 4), which follows the analog-to-digital converter, first stores the uncorrected values and then applies them word by word to a subtraction circuit to

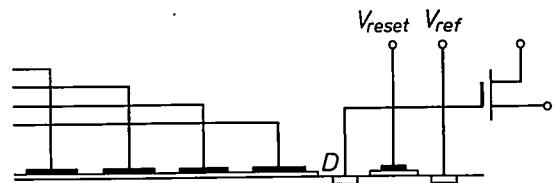


Fig. 18. Output circuit of the P²CCD. The drain diffusion D is connected to the gate of an MOS transistor, which is connected externally as a source follower. The second MOS transistor, operated by the reset voltage V_{reset} , returns the output to the reference potential V_{ref} between successive charge packets. The MOS transistors are included with the P²CCD on the same chip.

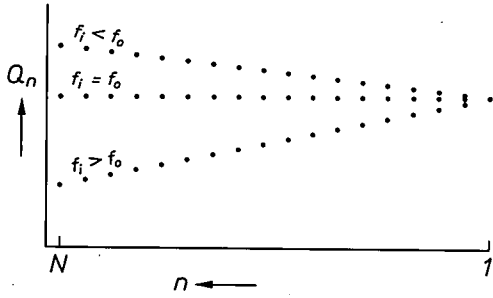


Fig. 19. Effects such as thermal generation of charge carriers cause a small parasitic charge Q_n to be added to the n th charge packet while it is in the P²CCD. The charge packets do not all remain in the P²CCD for the same length of time because the rate of read-in (clock frequency f_i) is generally not the same as the rate of read-out (clock frequency f_o), and Q_n is therefore generally not of the same magnitude for all the charge packets. This would give a sloping picture on the screen, and so zero correction is incorporated to prevent this.

which the associated value of the zero level is also applied. The difference between the two is again read into the buffer memory.

Appendix: Calculation of the potential distribution in a PCCD or a P²CCD

The potential V for every value of y in fig. 8 or fig. 11 can be calculated for any arbitrary variation of the doping $N(y)$ by integrating Poisson's equation. The one-dimensional form of the equation is sufficient here:

$$\frac{d^2V}{dy^2} = -\frac{dE}{dy} = -e\frac{N(y)}{\epsilon_{Si}}, \quad (A1)$$

where ϵ_{Si} represents the permittivity of the silicon. Taking a single integration between the limits y_1 and y_2 , which are still arbitrary at this stage, we have

$$E(y_2) - E(y_1) = \frac{e}{\epsilon_{Si}} \int_{y_1}^{y_2} N dy. \quad (A2)$$

Integrating again:

$$\begin{aligned} V(y_2) - V(y_1) &= \\ &= -\int_{y_1}^{y_2} E dy = -[Ey]_{y_1}^{y_2} + \int_{E(y_1)}^{E(y_2)} y dE = \\ &= y_1 E(y_1) - y_2 E(y_2) + \frac{e}{\epsilon_{Si}} \int_{y_1}^{y_2} y N(y) dy. \end{aligned} \quad (A3)$$

We now apply this result to the structure outlined in fig. 8 or fig. 11. It is assumed that no free charge carriers are present between $y = 0$ and $y = d_s$. The potential of the electrode is V_{el} . We shall first of all calculate the maximum potential at $y = y_m$; y_m is dif-

ferent for differing electrode potentials. We shall use the boundary condition $E(y_m) = 0$.

For the potential difference across the oxide layer we have

$$V_{el} - V(0) = \int_{-d_0}^0 E_{ox} dy. \quad (A4)$$

The electric flux density in the oxide layer is equal to that in the silicon at $y = 0$:

$$\epsilon_{ox} E_{ox} = \epsilon_{Si} E(0). \quad (A5)$$

It also follows from (A2) that

$$\begin{aligned} E(0) &= E(y_m) - \frac{e}{\epsilon_{Si}} \int_0^{y_m} N dy \\ &= -\frac{e}{\epsilon_{Si}} \int_0^{y_m} N dy. \end{aligned} \quad (A6)$$

From (A4) to (A6) the potential difference across the oxide layer is

$$V_{el} - V(0) = -\frac{e d_0}{\epsilon_{ox}} \int_0^{y_m} N dy. \quad (A7)$$

The potential difference across the silicon layer $0 \leq y \leq y_m$ is

$$V(y_m) - V(0) = \frac{e}{\epsilon_{Si}} \int_0^{y_m} y N dy. \quad (A8)$$

The potential difference between the electrode and y_m is therefore

$$V(y_m) - V_{el} = \frac{e d_0}{\epsilon_{ox}} \int_0^{y_m} N dy + \frac{e}{\epsilon_{Si}} \int_0^{y_m} y N dy. \quad (A9)$$

It can simplify the calculations somewhat if the entire space charge distributed between $y = 0$ and $y = y_m$ is considered to be concentrated at a 'centre of gravity' at a distance y_{cg} from the oxide-silicon interface. (A9) can then be rewritten in the form

$$V(y_m) - V_{el} = e \left(\frac{d_0}{\epsilon_{ox}} + \frac{y_{cg}}{\epsilon_{Si}} \right) \int_0^{y_m} N dy, \quad (A10)$$

where

$$y_{cg} = \frac{\int_0^{y_m} y N dy}{\int_0^{y_m} N dy}. \quad (A11)$$

We can then easily extend the equations so that they

give the potential at an arbitrary point $0 < y' < d_s$.

We then find:

$$V(y') - V_{e1} = e \left(\frac{d_0}{\epsilon_{ox}} + \frac{y'}{\epsilon_{Si}} \right) \int_{y'}^{d_s} N dy + e \left(\frac{d_0}{\epsilon_{ox}} + \frac{y'_{cg}}{\epsilon_{Si}} \right) \int_0^{y'} N dy, \quad (A12)$$

where y'_{cg} is the centre of gravity of the space charge in the range $0 < y < y'$.

Summary. The Philips PM 3310 storage oscilloscope has a digital semiconductor memory. A very high signal-sampling rate (up to 50 MHz) is possible because the signal samples obtained are stored temporarily in the form of charge packets in a profiled peristaltic charge-coupled device (P²CCD) and read out at a much slower rate (about 78 kHz); at this low frequency both analog-to-digital conversion and storage in a memory are more easily accomplished. The P²CCD is a development from Philips Research Laboratories and includes a double layer of n-type silicon on a p-type substrate. The double layer consists of a very thin, highly doped top layer and a less-highly doped bottom layer. The top layer is coated with SiO₂. Small charge packets are located in the bottom layer, near the top layer. Larger packets penetrate into the top layer, but will generally not reach the SiO₂ and thus will not be disturbed by the surface states at the Si-SiO₂ interface. The charge-handling capacity is nevertheless high.

Electrochemiluminescence in electrolyte-free solutions

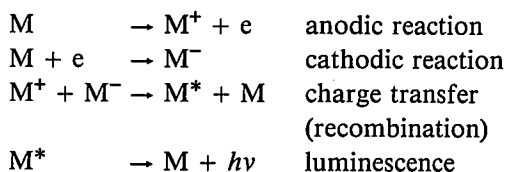
H. Schaper, H. Köstlin and E. Schnedler

In electrochemiluminescence (ECL) light is emitted as a result of charge transfer between positive and negative radical ions formed by electrochemical reactions. Until recently ECL cells seemed unsuitable for practical applications because of the problems involved in the use of polar solutions containing an electrolyte. Scientists at Philips Forschungslaboratorium Aachen (PFA) have now developed an improved type of ECL cell, consisting of a thin-layer cell in which the solution contains no additional electrolyte. This cell can be used in combination with a high-pressure liquid chromatograph to detect some of the aromatic hydrocarbons, or as background illumination in passive display devices.

Introduction

In the early sixties American scientists observed luminescence during the electrolysis of certain aromatic hydrocarbons in organic solvents^[1]. Investigations into the reaction mechanism revealed that the luminescence was produced by a charge transfer between positive and negative radical ions that had been generated electrochemically from the hydrocarbons. The effect was therefore called electrochemiluminescence (ECL).

In general terms, ECL is a regenerative process, in which the starting compounds are first electrochemically oxidized or reduced and then return to their original state via the charge transfer and the luminescence:



Owing to the reactivity of the M^+ and M^- ions, the solution must be free of water, oxygen, acidic or alkaline impurities. This imposes severe constraints on the solvents and reaction vessels.

Until recently the normal practice has been to add an electrolyte to ECL cells to make the solution sufficiently conductive. A polar organic solvent, such as acetonitrile or benzonitrile, is used for this purpose. The ECL cells are usually driven by an alternating voltage. As a result of the voltage reversals the positive and negative ions are generated at the same electrode, so that the charge transfer and luminescence are confined to a small area near this electrode^[2]. For practical applications, however, such cells are not very attractive. It is very difficult to make the polar solutions sufficiently pure, the use of an alternating voltage is not immediately compatible with battery-operated electronic circuits, the efficiency of the cells is at the most 1%, and their life is usually no more than a few hours.

ECL cells can also be driven by direct voltage^{[3][4]}. The problem with the types reported, however, is that after being generated at the anode and cathode the

-
- [1] D. M. Hercules, *Science* **145**, 808, 1964.
 R. E. Visco and E. A. Chandross, *J. Amer. Chem. Soc.* **86**, 5350, 1964.
 K. S. V. Santhanam and A. J. Bard, *J. Amer. Chem. Soc.* **87**, 139, 1965.
- [2] H. Köstlin and H. Schaper, in: *Philips – Unsere Forschung in Deutschland*, Vol. III, p. 59, 1980.
- [3] D. Laser and A. J. Bard, *J. Electrochem. Soc.* **122**, 632, 1975.
- [4] J. S. Dunnett and M. Voinov, *J. Chem. Soc., Faraday Trans. I*, **73**, 853, 1977.

Dr H. Schaper, Dr H. Köstlin and Dr E. Schnedler are with Philips GmbH Forschungslaboratorium Aachen, Aachen, West Germany.

positive and negative ions are relatively far apart and charge transfer between them is impaired because of their low stability. Another difficulty arises from electrochemical side reactions in a solution containing an electrolyte [6].

For these reasons we have investigated the possibilities of a less complicated and more stable type of ECL cell, which can be driven by direct voltage and is more suitable for practical applications. We therefore carried out ECL experiments with weak polar solvents, with no added electrolyte and which can be made thoroughly water-free and purified. Owing to the poor conductivity of such solvents the generation of ECL requires the use of a thin-layer cell, in which the distance between the electrodes is very small. We have obtained good results with thin-layer cells containing 5,6,11,12-tetraphenyltetracene (rubrene) as the luminescent compound and 1,2-dimethoxyethane (DME) as the solvent [6][7]. DME is a colourless liquid with a dielectric constant (relative permittivity) of 3.5. As far as we know, solvents with such a low polarity have seldom been used in ECL experiments.

Our investigations have led to a better understanding and hence to a better control of the ECL process. It is now possible to make thin-layer cells for specific applications: they can be used in combination with a high-pressure liquid chromatograph as detectors for certain aromatic compounds, or as background illumination in passive display devices.

In this article we shall first describe electrolyte-free thin-layer cells for ECL. We shall then discuss the processes occurring in these cells, paying particular attention to the electrostatic behaviour of the solutions and to the generation, charge transfer (recombination) and transport of the radical ions. Finally, we shall consider the applications mentioned.

Electrolyte-free thin-layer cells for ECL

Manufacture

Fig. 1 is a diagram of a thin-layer cell that we have developed for ECL. The cell contains two plates of soda-lime glass, both coated with an electrode film. The volume for the ECL solution is formed either by etching a cavity into one of the plates, or by separating the plates by strips of thin glass. To permit the insertion of a filling tube, also made of soda-lime glass, a semicircular groove may be ground in each plate on one side. The electrode films are deposited on the inner surface of the glass plates by well-known techniques such as spray pyrolysis, for $\text{In}_2\text{O}_3:\text{Sn}$ and $\text{SnO}_2:\text{Sb}$ [8], or vacuum evaporation, for metals like gold, silver and platinum. The electrode films are given the required configuration by means of photo-

lithography and chemical etching. The external electrical contacts are strengthened with fired silver paste. The cells are sealed with pyroceramics, as in television-tube manufacture.

An excess of solid rubrene, which has first been dried at 80 °C in a vacuum oven, is introduced into the vacuum-tight cell. The cell is then connected to a

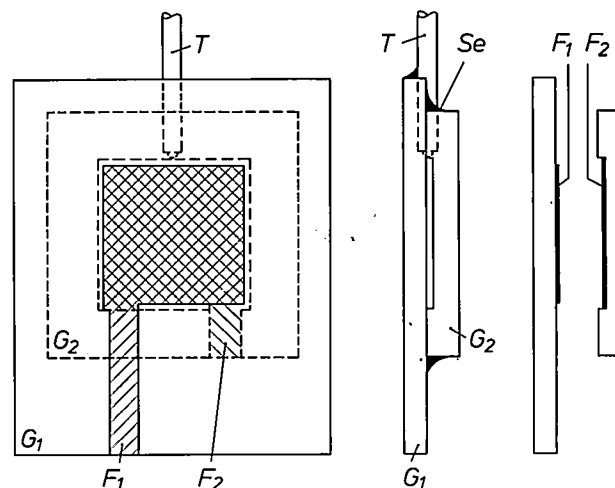


Fig. 1. Diagram of a thin-layer ECL cell. The glass plates G_1 and G_2 are coated with the electrode films F_1 and F_2 . T filling tube for introducing the ECL solution. Se pyroceramic seal.

vacuum line by its filling tube for the addition of the solvent, DME. The DME, to which a little benzophenone has been added, is usually first distilled several times over potassium until a blue colour reveals the presence of benzophenone radical anions, indicating that the solvent is sufficiently water-free. A reservoir containing such water-free DME, together with a small amount of potassium and benzophenone, is now connected to the vacuum line. The solvent is then repeatedly degassed and distilled under vacuum into the ECL cell. Next the solvent is redistilled into the reservoir, degassed and again distilled into the cell. This procedure is repeated a few times until the cell is free of volatile impurities. Finally, the filling tube is sealed off from the vacuum line. An ECL cell obtained in this way contains a saturated solution of rubrene in DME, corresponding to a concentration of about 4×10^{-3} mol/l at room temperature.

Properties

In fig. 2 the measured current density and luminescence intensity are plotted against the applied direct voltage for an ECL cell with an electrode spacing of 60 μm . At voltages up to about 2.3 V the solution behaves like an ohmic resistance of about $2 \times 10^5 \Omega$. At higher voltages an electrochemical current is gener-

ated, which is limited, as might be expected, by the diffusion rate of the rubrene molecules or by the distortion of the electric field. The current-voltage curve does not have a plateau, however, as in conventional polarography; instead the current increases linearly with increasing voltage above about 3 V. The luminescence starts at a voltage somewhat higher than the threshold voltage of the electrochemical current. The luminescence intensity then increases in about the same way with the applied voltage, so that the ECL efficiency — the ratio of the flux of photons emitted to the electrical current — remains constant at about 1%.

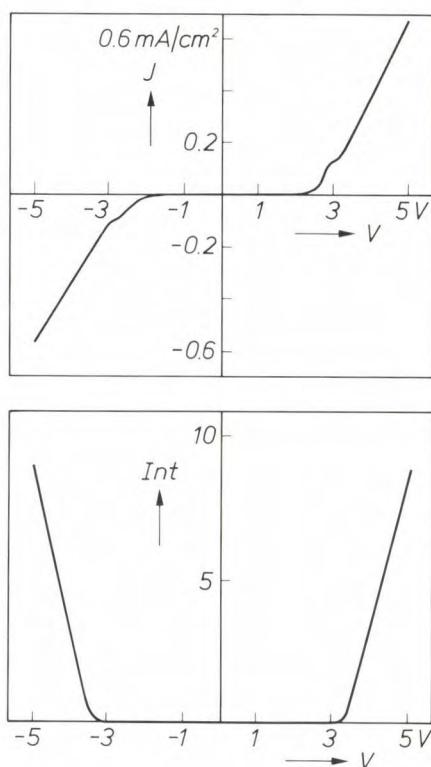


Fig. 2. Electrochemical current density J and luminescence intensity Int (in arbitrary units) as a function of the applied voltage V for an ECL cell in which the distance between the electrodes, both of $\text{In}_2\text{O}_3:\text{Sn}$, is $60 \mu\text{m}$. The voltage was increased at a rate of 5 mV/s . The cell contains a saturated solution of rubrene ($4 \times 10^{-3} \text{ mol/l}$) in DME. Above a threshold voltage both J and Int increase almost linearly with V . The luminescence starts at a somewhat higher voltage than the threshold voltage for the current.

In addition to the measurements with direct voltage we also measured the conductivity for alternating voltage under electrostatic conditions, i.e. at voltages lower than the threshold voltage for generating an electrochemical current. The conductivity is of the order of $10^{-7} \Omega^{-1}\text{cm}^{-1}$, and is virtually independent of the frequency, which we varied between 10 kHz and 1 MHz . Assuming that the molar conductivity is about $100 \Omega^{-1}\text{cm}^2\text{mol}^{-1}$, we can estimate that the ion concentration due to impurities in the solution is of the order of 10^{-9} mol/cm^3 .

Observation through a microscope reveals that the luminescence is not always uniformly distributed over the cell. An observer looking down at the plane of the thin layer can see a variety of regular microstructures, whose shape can be varied with the applied voltage and whose size depends on the electrode spacing. *Fig. 3* shows four luminescence photographs of an ECL cell with an electrode spacing of $75 \mu\text{m}$ at different volt-

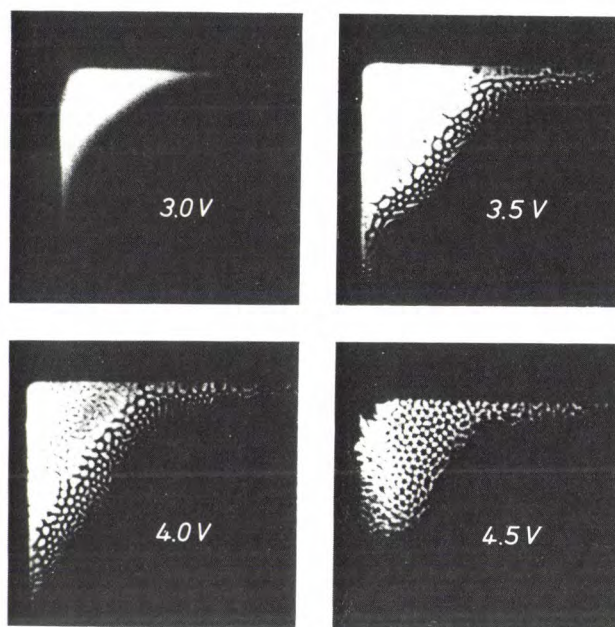


Fig. 3. Photographs showing the luminescence structure of an ECL cell with an electrode spacing of $75 \mu\text{m}$, at four different voltages. The structure of the luminescence becomes more evident as the voltage increases.

ages. Above 3.0 V structured patterns are observed that become sharper with increasing voltage. The patterns at about 4.5 V and above are regular hexagons resembling honeycombs. *Fig. 4* shows that the diameter of the hexagons is always larger than the electrode spacing by a factor of 1.5 to 2. The structures seem to stand motionless in the liquid, so that sharp photographs can easily be taken with exposure times up to a minute. If, however, the cell contains solid particles, e.g. undissolved rubrene crystals, very rapid movements of these particles can be seen, indicating a strong fluid flow.

Not only the luminescence structure but also the electrochemical current and the luminescence intensity are strongly dependent on the electrode spacing; see

- [5] J. S. Dunnett and M. Voïnov, *J. electroanal. Chem.* **89**, 181, 1978.
- [6] H. Köstlin and H. Schaper, *Phys. Lett.* **76A**, 455, 1980.
- [7] H. Schaper, H. Köstlin and E. Schnedler, *J. Electrochem. Soc.*, in print.
- [8] H. Köstlin, R. Jost and W. Lems, *Phys. Stat. sol. (a)* **29**, 87, 1975.

fig. 5. Both the luminescence intensity and the ECL efficiency increase with decreasing electrode spacing. The threshold voltages of the current and the luminescence, however, do not depend on the electrode spacing.

The current and the luminescence intensity are also affected by the temperature; see Table I. As the temperature rises from $-35\text{ }^{\circ}\text{C}$ they both first increase and then decrease. When the temperature rises above room temperature the ECL efficiency decreases. This

Table I. Electrochemical current I and luminescence intensity Int of an electrolyte-free thin-layer ECL cell at various temperatures, measured at a direct voltage of 5 V.

Temperature ($^{\circ}\text{C}$)	I (mA)	Int (relative units)
-35	3.50	2.20
-8	4.20	2.80
20	3.50	1
42	2.35	0.60
59	1.30	0.25

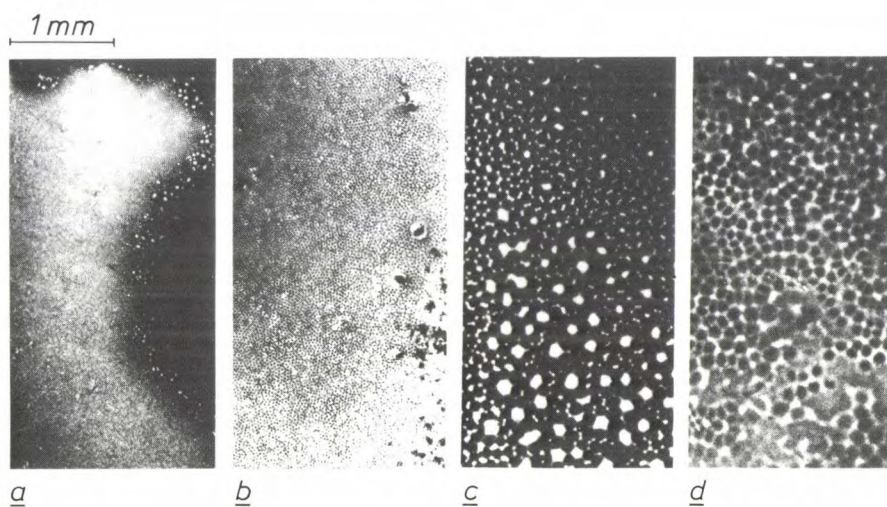


Fig. 4. Photographs of the hexagonal luminescence patterns observed in four ECL cells, driven at 5 V, with different electrode spacings: a) $12\text{ }\mu\text{m}$; b) $25\text{ }\mu\text{m}$; c) $50\text{ }\mu\text{m}$; d) $100\text{ }\mu\text{m}$. In all cases the diameter of the hexagons is larger than the electrode spacing by a factor of 1.5 to 2.

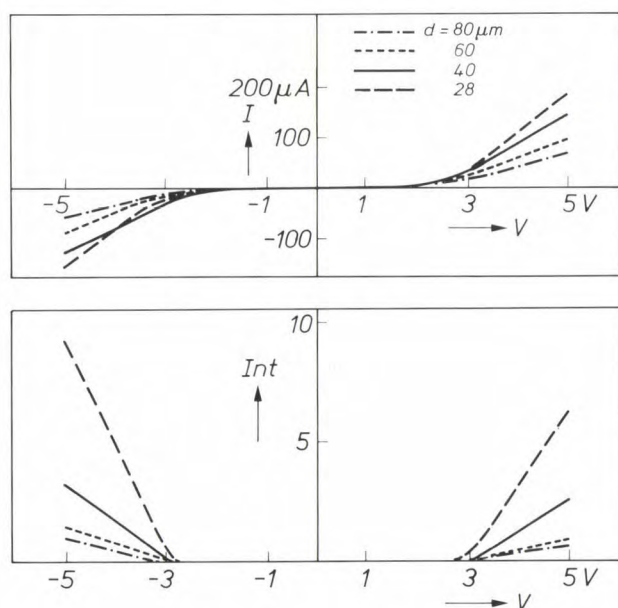


Fig. 5. Electrochemical current I and luminescence intensity Int (in arbitrary units) as a function of the applied voltage V , measured for ECL cells with a different electrode spacing d . A decrease in d causes a large increase in I and Int , but has very little effect on the threshold voltages.

is largely attributable to a decrease in the luminescence intensity, as is shown by measurements under ultraviolet excitation with a mercury discharge lamp. The decrease in the current, indicating a reduction in the ion transport, is however quite surprising in view of the lower viscosity at higher temperatures.

Compared with the direct-voltage ECL cells described earlier [3][4], the new electrolyte-free thin-layer cells have a much longer life. The chemical stability of rubrene radical ions in pure DME can be estimated if it is assumed that the electrochemical current density at constant voltage is determined by the rubrene concentration in the solution. The number of times N that a rubrene molecule on average is involved in the electrochemical reaction cycle, before being lost in side reactions, can be regarded as a measure of the reversibility of the ECL process. This number can be calculated from

$$N = \frac{J\tau}{acFd}, \quad (1)$$

where J is the initial current density, τ is the time in

which the current density has decreased by a factor α , c is the rubrene concentration, F is the Faraday constant and d is the electrode spacing. From fig. 2 it can be seen that $J = 0.35 \text{ mA/cm}^2$ at 4 V for a cell with $c = 4 \times 10^{-3} \text{ mol/l}$ and $d = 60 \text{ }\mu\text{m}$. The decay in the current density is about 1% per hour, so that N is of the order of 5×10^4 .

Processes in electrolyte-free thin-layer cells

In an attempt to explain the results obtained with electrolyte-free thin-layer cells for ECL, we have studied the electrostatic behaviour, the generation

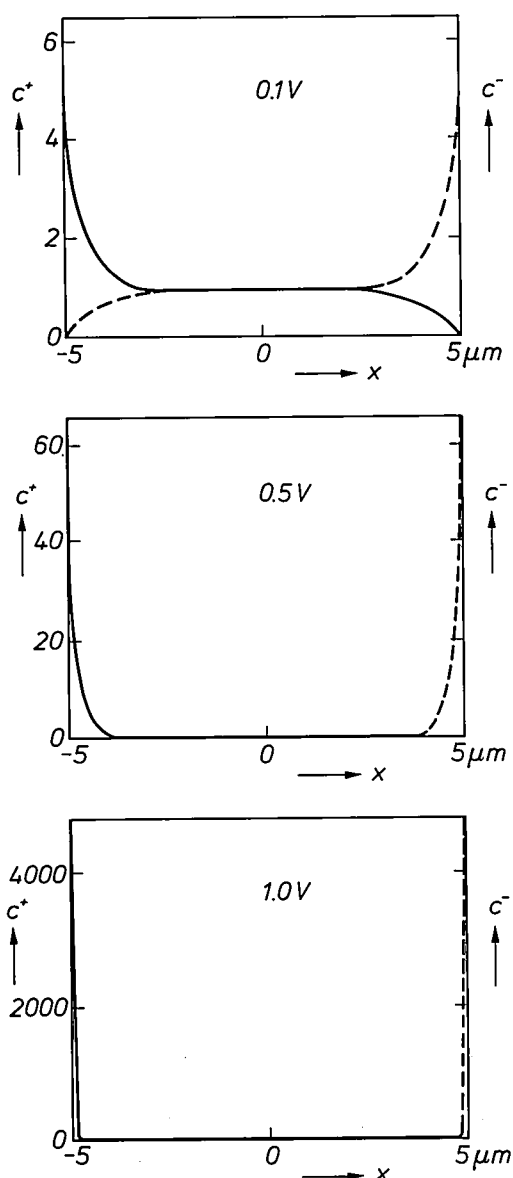


Fig. 6. Calculated distributions of cations and anions in an ECL cell with the electrodes at $x = \pm 5 \text{ }\mu\text{m}$ and with an ion concentration of $10^{-11} \text{ mol/cm}^3$, at three different voltages below the threshold voltage. The plotted concentrations c^+ (solid lines) and c^- (dashed lines) are normalized with respect to the cation and anion concentrations with no electric field. At only 1 V most of the ions are close to the electrodes.

and charge transfer of the radical ions, and the ion transport in the solution. We shall now discuss each of these subjects in more detail.

Electrostatic behaviour

In a conventional electrolyte solution the electrode potentials are almost completely screened by the ions present. The potential difference between electrodes and solution therefore remains restricted to a very small area near the electrode surfaces, while the bulk of the solution has an almost constant potential. This implies that electric fields only occur near these surfaces. As a result, charged particles in a stationary solution are transported by means of diffusion. In the new cells the situation is clearly different: the cell dimensions are very small and the ion concentration is only of the order of 10^{-9} mol/cm^3 .

When a constant voltage is applied to a solution in which the positive and negative ionic impurities have a charge z , a symmetrical ion distribution will result. The concentrations $c^+(x)$ and $c^-(x)$ of the positive and negative ions can be written as a function of the potential $\phi(x)$, from the Maxwell-Boltzmann distribution law:

$$c^+(x) = c_0 \exp\left[-\frac{\phi(x) - \phi(0)}{RT/zF}\right] \quad (2)$$

$$c^-(x) = c_0 \exp\left[\frac{\phi(x) - \phi(0)}{RT/zF}\right],$$

where R is the molar gas constant, x refers to the direction perpendicular to the electrode surfaces, and c_0 and $\phi(0)$ are the ion concentration and the potential at $x = 0$, the mid-plane of the cell. In view of the symmetry, $\phi(0)$ can be put equal to zero. The charge density $q(x)$ as a function of the position is equal to

$$q(x) = zF[c^+(x) - c^-(x)]. \quad (3)$$

The relationship between the charge density and the potential $\phi(x)$ is given by Poisson's equation:

$$\frac{d^2\phi(x)}{dx^2} = \frac{-q(x)}{\epsilon_r\epsilon_0}, \quad (4)$$

where ϵ_r is the dielectric constant of the solvent and ϵ_0 is the permittivity of free space.

In solving the above equations [9] there are two limiting cases, which depend on the ratio γ of the electrical field energy to the thermal energy of the ions at the mid-plane of the cell:

$$\gamma = \frac{\epsilon_r\epsilon_0 E_0^2}{2c_0RT}, \quad (5)$$

where E_0 is the electric field-strength at the mid-plane

[9] More details of the calculations are given in H. Schaper and E. Schnedler, *J. electroanal. Chem.*, in print.

of the cell. When the electric field energy at the mid-plane is much lower than the thermal energy ($\gamma \ll 1$), the bulk of the solution contains a relatively large proportion of the ions. The concentration c_0 decreases strongly with increasing applied voltage. The field-strength E_0 , on the other hand, increases until the electric field energy becomes much higher than the thermal energy ($\gamma \gg 1$). In this case nearly all of the ions are collected near the electrodes, and hardly any are contained in the bulk of the solution.

Fig. 6 gives the calculated distributions of the positive and negative ions for three different voltages. Even at a relatively low voltage (1 V) most of the ions are in a small region near the electrodes, whereas elsewhere the cell is practically ion-free, and comparable to an ideal dielectric. As a consequence strong electric fields can be generated in the bulk of the solution. In fig. 7 the electric field-strength E_0 is plotted against the applied voltage for different ion concentrations. It can be seen that E_0 may increase steeply with decreasing ion concentration. Above a certain threshold voltage E_0 increases linearly with voltage. This threshold voltage increases with the ion concentration. This is because at higher concentrations more ions must be drawn towards the electrodes before an electric field can be built up in the bulk.

In fig. 8 the calculated electric field-strength E_s near the electrodes is plotted against voltage. Above a cer-

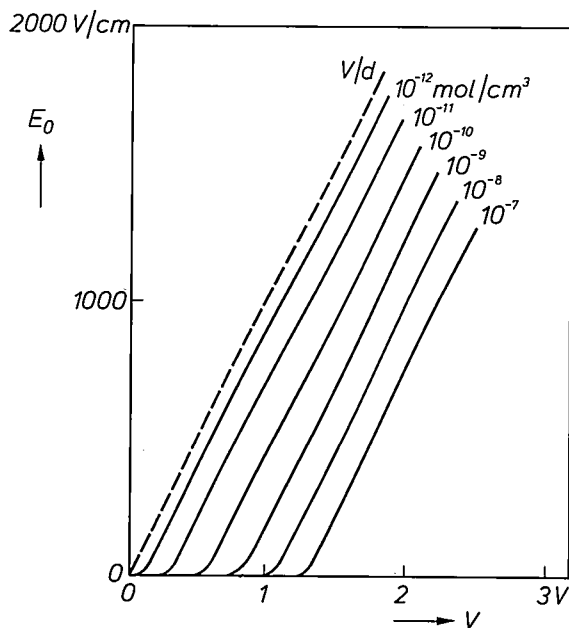


Fig. 7. Calculated mid-plane electric field-strength E_0 as a function of the applied voltage V , for an ECL cell with electrode spacing d of $10 \mu\text{m}$, at different ion concentrations. The threshold voltage above which E_0 increases linearly with V increases with the ion concentration. At a very low concentration (e.g. $10^{-12} \text{ mol/cm}^3$) the value of E_0 is no longer much different from V/d (dashed line).

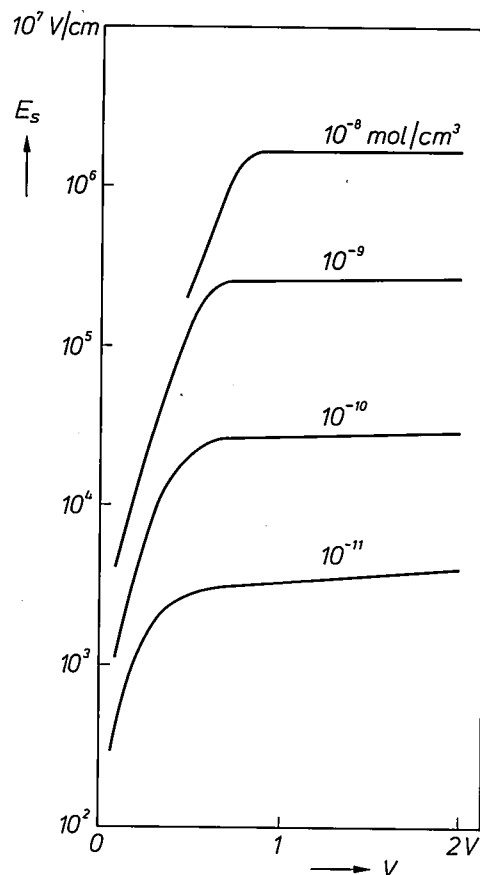


Fig. 8. Calculated electric field-strength E_s near the electrodes as a function of the applied voltage V , for an ECL cell with electrode spacing of $10 \mu\text{m}$, at different ion concentrations. While E_0 decreases when the ion concentration is increased (fig. 7), E_s increases strongly.

tain voltage E_s is practically constant, since nearly all the ions are then near the electrodes. The value of E_s is of importance in connection with the charge transfer between electrode and solution. In an electrolyte-free solution this transfer may be assumed to occur at about the same potential as in electrolyte solutions, i.e. $\leq 2 \text{ V}$. Assuming further that this potential drop must occur in a region roughly equal in size to a molecular diameter ($\approx 1 \text{ nm}$), a value for E_s of the order of 10^7 V/cm can be shown to be necessary. As can be seen in fig. 8, such a value is approached with an ion concentration as low as 10^{-8} mol/cm^3 at an applied voltage of only 2 to 3 V.

The electric field-strengths E_0 and E_s depend not only on the applied voltage V and the ion concentration c , but also on the electrode spacing d ; see fig. 9. A cell with a small electrode spacing behaves like a dielectric so that

$$E_s = E_0 = V/d. \quad (6)$$

At larger electrode spacings, E_s becomes higher than E_0 . Above a certain spacing and voltage the cell behaves like an electrolytic solution with a low ion

concentration, and then E_s is proportional to c and d :

$$E_s = \frac{zFcd}{\epsilon_r \epsilon_0} \quad (7)$$

In this case E_s is no longer dependent on the voltage. From eqs (6) and (7) the transition point d_t between the electric and the electrolytic behaviour can be derived:

$$d_t = \sqrt{\frac{\epsilon_r \epsilon_0 V}{zFc}} \quad (8)$$

The transition point thus shifts towards a lower electrode spacing if the ion concentration is higher, or the applied voltage is lower.

Ion generation and charge transfer

As shown in fig. 2, the luminescence starts at a higher voltage than the electrochemical current. Another feature is that the observed luminescence is not symmetrically distributed over the cell. In most experiments with cells containing rubrene in DME the luminescence mainly occurs near the anode. We have demonstrated this with a specially made ECL cell in which the electrodes were mounted on the same glass plate in the form of two ‘interlacing’ combs. Fig. 10 shows a photograph of the cell without luminescence and two photographs of the luminescing cell for the two polarities of the applied voltage. In both cases the luminescence appears in a very narrow region around

the anodically polarized comb, both with In_2O_3 and with gold electrodes. In another experiment the luminescence spectrum of an ordinary ECL cell with a gold electrode and a transparent In_2O_3 electrode was measured perpendicular to the electrode planes. Fig. 11 shows that the spectrum is shifted to longer wavelengths when the cathode is closer to the spectrometer than the anode. This also indicates that the luminescence is restricted to the anode region. When the anode is further away from the spectrometer than the cathode, the light emitted must pass through the rubrene solution before reaching the spectrometer. The associated self-absorption then causes a spectral shift to longer wavelengths.

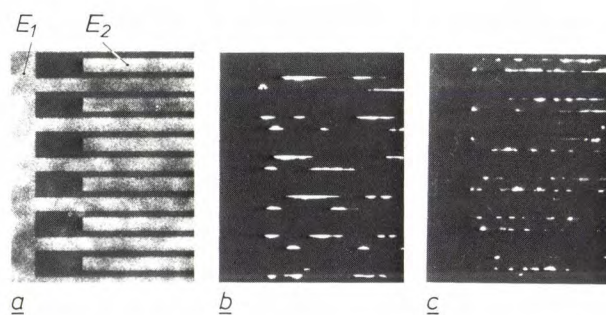


Fig. 10. Photographs of a thin-layer ECL cell with two interlacing electrode combs for localizing the luminescence. a) Electrodes E_1 and E_2 without luminescence. b) Luminescence, with the comb coming from the left (E_1) as the anode. c) Luminescence, with the comb coming from the right (E_2) as the anode. In both cases the luminescence occurs mainly near the comb acting as the anode.

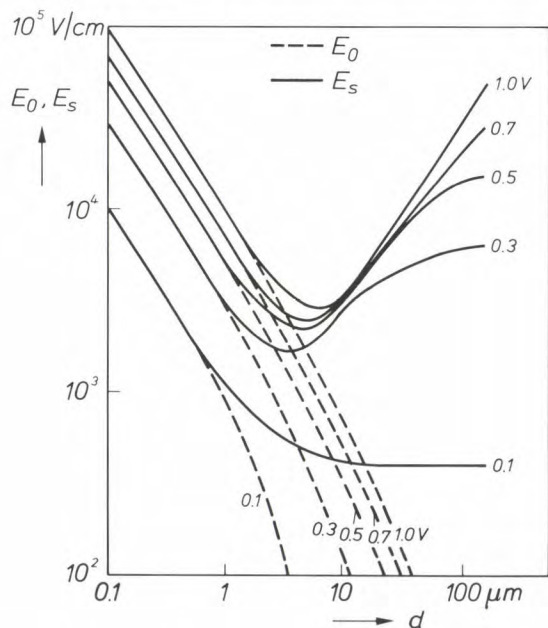


Fig. 9. Calculated electric field-strength in the mid-plane (E_0) and near the electrodes (E_s) of an ECL cell with an ion concentration of 10^{-11} mol/cm³, as a function of the electrode spacing d at different voltages. At small spacings both E_0 and E_s are proportional to $1/d$: the cell behaves like a dielectric. At large spacings the cell behaves like an electrolytic solution, and E_s increases with d if the voltage is not too low.

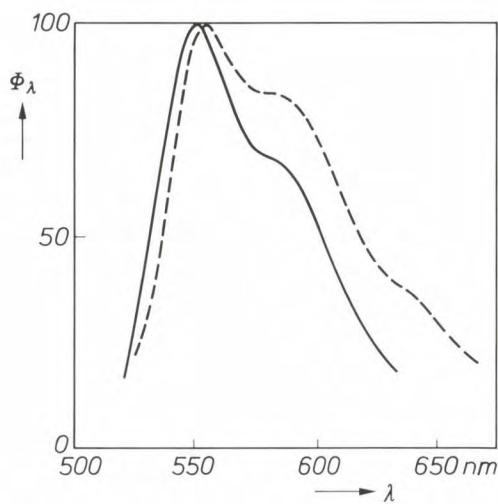


Fig. 11. Luminescence spectrum of a rubrene/DME solution in a thin-layer ECL cell with a gold electrode and a transparent In_2O_3 electrode. The spectrometer is positioned near the In_2O_3 electrode. Φ_λ is the relative emission intensity (in arbitrary units) as a function of the wavelength λ ; the maximum of both curves has been set equal to 100. When the In_2O_3 electrode is the anode, the spectrum (solid line) indicates only little self-absorption. However, when the gold electrode is the anode, the spectrum (dashed line) is shifted to longer wavelengths, owing to strong self-absorption of the solution. This also demonstrates that the luminescence is mainly excited near the anode.

These results indicate that the rubrene anions are formed in preference to the cations. In this way the symmetrical ion distribution, produced under electrostatic conditions, is distorted by the onset of the electrochemical process. The negative ions formed at the cathode are attracted to the anode by the strong electric fields (≈ 1000 V/cm) in the solution. The drift of negative ions corresponds to an electrochemical current, whose density is given by

$$J = \mu q(x) E, \quad (9)$$

where μ is the mobility of the rubrene anions, $q(x)$ is their charge density and E is the electric field-strength in the bulk of the solution. From Poisson's equation:

$$q(x) = \varepsilon_r \varepsilon_0 dE/dx, \quad (10)$$

we obtain

$$J = \varepsilon_r \varepsilon_0 \mu E dE/dx, \quad (11)$$

from which we can derive an expression for the electric field-strength $E(x)$ by integration:

$$E^2(x) = E^2(x_a) + \frac{2J(x - x_a)}{\varepsilon_r \varepsilon_0 \mu}, \quad (12)$$

where x_a refers to the position just in front of the ionic layers near the cathode. This equation shows how the electric field in the bulk of the solution is distorted by the anion drift. This distortion weakens the electric field near the cathode and strengthens it near the anode, thus stimulating the production of rubrene cations. For this reason the charge transfer between the anions and cations, and hence the luminescence as well, occurs at a voltage that is a little higher than the threshold voltage for the electrochemical current. The solution always contains an excess of anions, however. Experiments on cells with comb electrodes indicate that this is also the case at higher voltages.

Ion transport

The conventional diffusion model does not help to explain the steep linear increase in the electrochemical current and the luminescence intensity at voltages above about 3 V (fig. 2). The observations through the microscope indicate that there is liquid flow in the thin-layer cell. The luminescence patterns (figs 3 and 4) strongly resemble the hexagonal structures known from thermoconvection^[10] and also found during electrohydrodynamic investigations in apolar solvents^{[11][12]}. We therefore assume that electrohydrodynamic convection forces are also responsible for the generation of the regular hexagonal structures in electrolyte-free ECL solutions. Above a certain threshold of the electric field the uniform drift of the generated ions turns into a much faster motion of the liquid, because the surrounding liquid is dragged

along by the ions as a result of the viscosity. The flow of the liquid is structured in regular patterns^{[13][14]}, which depend on the electric field-strength. In the case of convection, the electrochemical current is no longer limited by diffusion or by field distortion (eq. (9)), because the transport of the generated ions and of the neutral molecules take place at the same velocity, as if the liquid was being stirred.

The onset of electrohydrodynamic convection can be described by a stability theory, based on Poisson's equation and on the laws of conservation of charge, mass, momentum and assuming incompressibility^{[11][13][15][16]}. The onset of convection is determined by the value of a dimensionless critical number, the electrical Rayleigh number:

$$Ra_e = \frac{\varepsilon_r \varepsilon_0 \bar{V}}{\rho \nu \nu}, \quad (13)$$

where ρ is the density of the solvent and ν the kinematic viscosity. \bar{V} , the voltage available for transport of ions in the cell, is equal to the difference between the applied voltage and the voltage necessary for the oxidation and reduction of rubrene. This equation shows that the onset of convective charge transport depends only on the applied voltage and not on the electrode spacing. This agrees with the results of fig. 5.

The order of magnitude of the average flow velocity v can be estimated by dimensional analysis:

$$v \approx \frac{\bar{V}}{d} \sqrt{\varepsilon_r \varepsilon_0 / \rho}. \quad (14)$$

The threshold voltage for convection is about 3 V and the voltage necessary for oxidation and reduction is about 2.5 V, so that \bar{V} is equal to 0.5 V. With $d = 50 \mu\text{m}$, $\rho \approx 1 \text{ g/cm}^3$ and $\varepsilon_r \approx 3.5$, the average flow velocity, just after the onset of convection, is of the order of 0.2 cm/s. Eq. (14) shows that v is proportional to \bar{V} , in agreement with the observed linear increase of the current. The interpretation of the experimental results of fig. 2 therefore leads to the conclusion that the overall degree of ionization of the rubrene solution should be constant in the linear range of the current-voltage curve. Under these conditions the degree of ionization can be estimated. With $d = 60 \mu\text{m}$, $J(3 \text{ V}) = 0.1 \text{ mA/cm}^2$, $v = 0.16 \text{ cm/s}$ and $c = 4 \times 10^{-6} \text{ mol/cm}^3$ the ionization degree β , given by

$$\beta = \frac{J}{Fcv}, \quad (15)$$

is of the order of 0.15%. This rather low value of β should not be seen as the absolute fraction of rubrene radical ions in the solution; it simply indicates the fraction of rubrene that is active in charge transport during the convection.

Applications

Detector for chromatography

Polycyclic aromatic hydrocarbons such as substituted anthracene, tetracene and pyrene are suspected of being carcinogenic. They occur as traces in smoke, plants and milk, for example, and therefore methods of detecting these compounds analytically are urgently needed. In connection with the development of a high-pressure liquid chromatograph (HPLC) for environmental analysis [17] the question arose as to whether the electrolyte-free thin-layer cells could be used as a novel type of detector for these compounds. Usually they are detected through their absorption or fluorescence when irradiated by an ultraviolet lamp, giving detection limits as low as 3.5×10^{-10} g/l [18]. Trace analysis with these methods is laborious, however, since it is difficult to separate the incident ultraviolet radiation and the measurement signal. Electrochemical detectors, on the other hand, are quite insensitive to aromatic hydrocarbons [19]. In addition, for connection to a continuously operating chromatograph the eluate first has to be mixed with a solution containing an electrolyte. Continuous operation is also very difficult with electrolyte-containing ECL cells driven by an alternating voltage [20]. Detection with electrolyte-free thin-layer ECL cells driven by direct voltage, however, seems quite promising, especially for continuous HPLC analysis [21].

Fig. 12 is a diagram of the combination of a liquid chromatograph and a thin-layer ECL cell. When a direct voltage is applied luminescence will occur in the cell if the liquid coming from the chromatograph column contains one or more electrochemiluminescent compounds. The light emitted is usually measured by means of a photomultiplier. The electrochemical current flowing in the cell during ECL is also measured. The eluents used must be free from dissolved oxygen, since otherwise the luminescence is strongly quenched. In addition hydrogen-containing

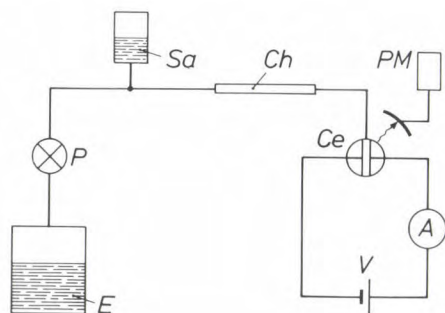


Fig. 12. Diagram of the combination of a high-pressure liquid chromatograph and a thin-layer ECL cell. *E* eluent. *P* pump. *Sa* sample. *Ch* chromatographic column. *Ce* thin-layer ECL cell. *V* direct-voltage source. *A* ammeter. *PM* photomultiplier.

impurities (e.g. water) have to be removed to avoid irreversible side reactions with the radical ions and to prevent corrosion of the cell electrodes.

The ECL cell we have developed for this application is shown diagrammatically in fig. 13. It contains two electrode plates, a spacer frame, a mounting plate, a pressure plate and a seal. The electrode plates consist of a planar glass plate coated with an electrode film. The pressure plate and one of the electrode plates

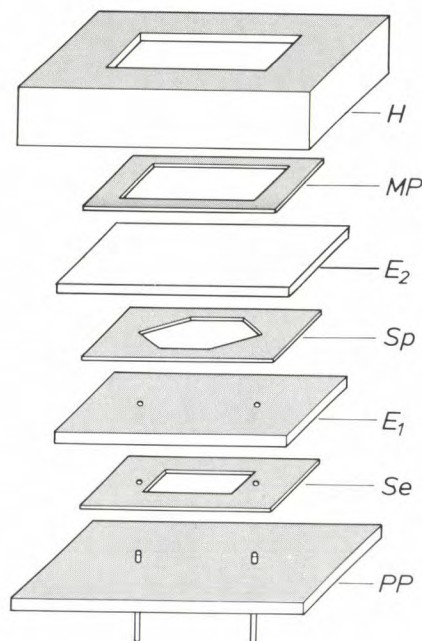


Fig. 13. Diagram showing some of the components of an ECL cell for combination with a high-pressure liquid chromatograph. *PP* pressure plate with inlet and outlet for the eluate. *Se* Teflon seal. *E₁* electrode plate consisting of a glass plate with electrode film. *Sp* spacer frame of Teflon. *E₂* electrode plate consisting of a glass plate with transparent electrode film. *MP* Teflon mounting plate. *H* housing with window for transmission of luminescence.

- [10] H. Bénard, *Rev. gén. Sci.* **11**, 1261, 1900.
 [11] J. C. Lacroix and P. Atten, *J. Electrostat.* **5**, 453, 1978.
 [12] P. Müräu and B. Singer, *J. appl. Phys.* **49**, 4820, 1978.
 [13] P. Atten and J. C. Lacroix, *J. Electrostat.* **5**, 439, 1978.
 [14] N. Felici and J. C. Lacroix, *J. Electrostat.* **5**, 135, 1978.
 [15] J. M. Schneider and P. K. Watson, *Phys. Fluids* **13**, 1948, 1970.
 [16] P. Atten and R. Moreau: *C.R. Acad. Sci. Paris A* **269**, 433 and 469, 1969, and *A* **270**, 415, 1970; *J. Mécan.* **11**, 471, 1972.
 [17] R. J. Dolphin and F. W. Willmott, *J. Chromatogr.* **149**, 161, 1978.
 [18] H. Engelhardt, *High performance liquid chromatography*, Springer, Berlin 1979.
 [19] K. Brunt, Thesis, Groningen 1980.
 [20] S. A. Cruser and A. J. Bard, *Anal. Lett.* **1**, 11, 1967.
 B. Fleet, G. F. Kirkbright and C. J. Pickford, *Talanta* **15**, 566, 1968.
 B. Fleet, P. N. Keliher, G. F. Kirkbright and C. J. Pickford, *Analyst* **94**, 847, 1969.
 B. Fleet, G. F. Kirkbright and C. J. Pickford, *Laboratory Practice* **19**, 804, 1970.
 T. M. Huret and J. T. Maloy, *J. Electrochem. Soc.* **121**, 1178, 1974.
 [21] H. Schaper, *J. electroanal. Chem.* **129**, 335, 1981.

contain holes and connections for the inlet and outlet of the eluate. The electrical contacts are made by means of terminals at the free electrode edges. The electrode facing the photomultiplier must be transparent and is normally made of pyrolytically deposited $\text{In}_2\text{O}_3:\text{Sn}$. The counter-electrode, which does not have to be transparent, usually consists of an evaporated metal film (gold, silver or aluminium). The spacer frame is made of a Teflon film with a thickness between 5 and 50 μm .

Fig. 14 shows how the electrochemical current and the luminescence intensity varied with time after 20 μl of a solution of 10^{-4} mol/l of rubrene in DME was injected into the ECL cell. For comparison, the signals were also measured without the chromatographic column. The presence of the column causes not only a considerable broadening of the measuring signals, especially at the beginning, but also a very large signal delay. The area enclosed by the curves does not change much, however.

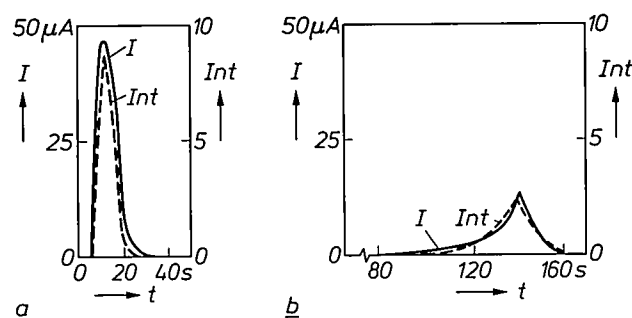


Fig. 14. Electrochemical current I and luminescence intensity Int (in arbitrary units) of a thin-layer ECL cell as a function of the time t after the injection of 20 μl of a rubrene/DME solution of 10^{-4} mol/l. The eluent is also DME. The flow rate is 0.6 ml/min. The ECL cell contains a gold cathode and an In_2O_3 anode at a spacing of 26 μm . The applied voltage is 30 V. *a*) Signals from the cell without chromatographic column. *b*) Signals from the cell in combination with a chromatographic column (fig. 12), to be used as HPLC detector. The presence of the column causes considerable broadening and delay of the signals.

The effect of the applied voltage on the electrochemical current and luminescence intensity is shown in fig. 15. We see again that voltages can be applied that far exceed the oxidation and reduction potentials of the detectable compounds without the introduction of undesirable side reactions. In fact, the ECL efficiency increases with the applied voltage as a result of the mass transport by electrohydrodynamic convection. This is a great advantage compared with conventional electrolytic cells, in which there is a considerable reduction in both the luminescence intensity and the electrochemical reversibility when the relevant electrochemical potentials are exceeded.

Other aromatic hydrocarbons such as 9,10-diphenylanthracene, pyrene and perylene gave similar results.

In view of the possible application of this ECL cell as a detector for HPLC, we also investigated the use of solvents with a polarity lower than that of DME. Less-polar solvents are preferred in chromatography because they can give a better separation. The ECL response does not vary much when rubrene, for example, is dissolved in different mixtures of DME and the less polar hexane. The luminescence intensity decreases only slightly when the percentage by volume of hexane increases to 90%. The same is true for the electrochemical current. The width of the luminescence peak and the current peak is not affected by the hexane percentage.

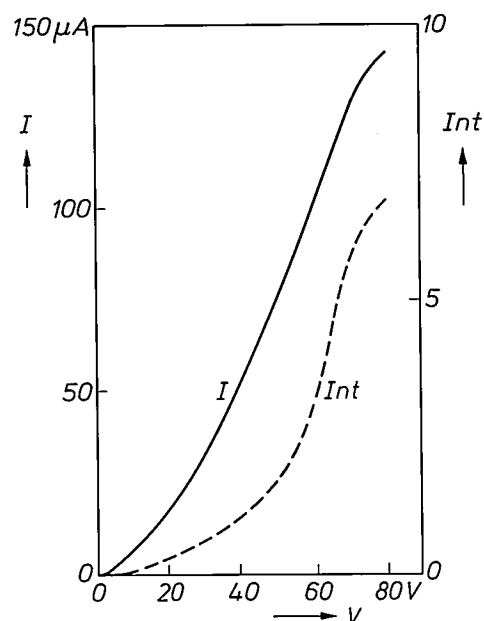


Fig. 15. Electrochemical current I and luminescence intensity Int (in arbitrary units) as a function of the applied voltage V for an HPLC/ECL combination. The eluent and the solvent consist of a mixture of DME and hexane (2:1). The sample and the measuring conditions are the same as in fig. 14. At higher voltages Int increases more rapidly than I , giving an enhanced ECL efficiency.

The luminescence intensity Int as a function of the rubrene concentration c was evaluated for DME solutions with a low rubrene content. The results are shown in fig. 16 in a log-log plot. The solid straight line corresponds to a power law of the type $Int \propto c^{1.3}$. The smallest quantity of rubrene that we could detect in these experiments was 5 μl of a solution of 10^{-8} mol/l, which corresponds to 5×10^{-14} mol or 2.7×10^{-11} g of rubrene. This compares well with the detection limits obtainable with conventional absorption and fluorescence detectors.

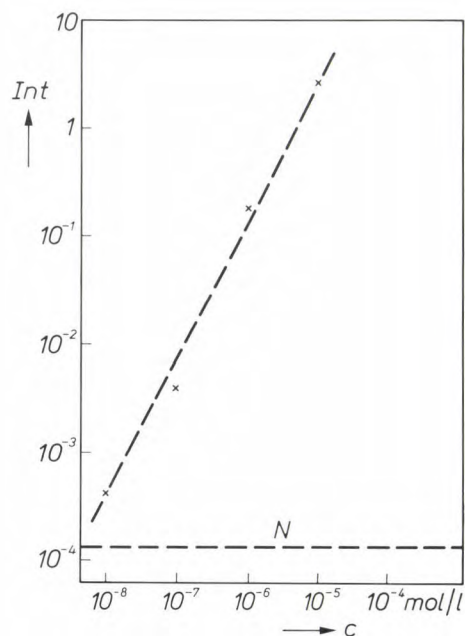


Fig. 16. Luminescence intensity Int (in arbitrary units) as a function of the rubrene concentration c in a 5 μ l DME solution, measured with an HPLC-ECL combination. The eluent is DME and the flow rate is 0.25 ml/min. The electrodes are both of In_2O_3 with a spacing of 26 μ m and a direct voltage of 10 V. The noise level N for the detection is also indicated. Concentrations as low as 10^{-8} mol/l can easily be detected.

It should be possible to improve the sensitivity even further by measuring at higher voltages or by cooling the photomultipliers to suppress the dark current. In addition, spectral analysis of the luminescence is possible, so that both identification and quantitative analysis of polycyclic aromatics should be feasible with the combination of an HPLC and an ECL cell. Since the thin-layer-cell described here can also be operated as an electrochemical detector, the detection of some non-luminescent compounds by observing the current is a further possibility.

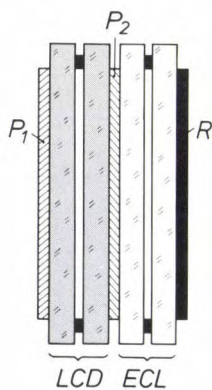


Fig. 17. Diagram of the combination of an LCD element and a thin-layer ECL cell for background illumination in dark surroundings. P_1 , P_2 polarizers. R reflecting metallic film.

Background illumination in passive displays

Liquid-crystal display (LCD) elements give the required information only when they are illuminated. For use in the dark they need additional illumination, which can be switched on briefly. The displays in digital watches and alarm clocks are examples. In these applications the LCD elements are usually illuminated by a miniature incandescent lamp. The power required for this illumination is of the order of 15 mW for watches and of the order of 100 mW for alarm clocks.



Fig. 18. Photograph of an LCD digital watch and a thin-layer ECL cell for incorporation in the watch.

When the thin-layer ECL cells described here are used for background illumination the power consumption can be lower by a factor of ten. Fig. 17 gives a diagram of the combination of an LCD element and an ECL cell. With external illumination, the light travels in turn through a polarizer P_1 and the LCD element, a second polarizer P_2 and the ECL cell, and is then reflected at the back. Depending on the arrangement of P_2 , parallel or crossed with respect to P_1 , the information appears dark against a bright background or vice versa. When there is no external illumination, the ECL cell can be switched on. The ECL light passes through P_2 , the LCD element and P_1 . Again depending on the relative arrangement of P_1 and P_2 , the information appears dark against a bright background or vice versa.

Since the distance between the electrodes in the ECL cell is less than 0.1 mm, the cell thickness is determined mainly by the thickness of the glass plates. In practice, the plates are about 0.5 mm thick, so that the cell thickness is about 1 mm. The addition of ECL

cells does not therefore significantly increase the total thickness of display units.

Fig. 18 shows an LCD digital watch and an ECL cell that can readily be incorporated behind the display unit. *Fig. 19* shows the display of two LCD digital alarm clocks, one with a miniature incandescent lamp and the other with a thin-layer ECL cell as background illumination. Even though it uses far less power, the ECL illumination gives a much clearer display.

Most of the experimental work in these investigations was carried out by Ing. grad. M. Peterek and K. H. Wilhelm.

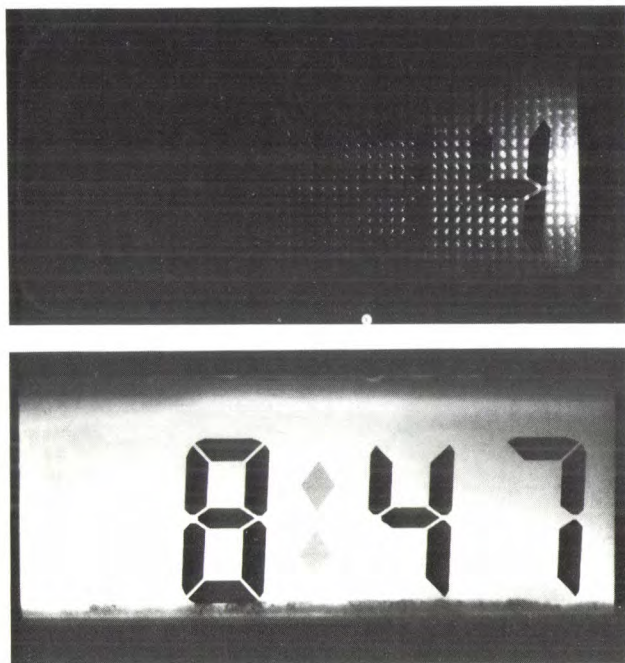


Fig. 19. Display of an LCD alarm clock with background illumination from a miniature incandescent lamp at 1.5 V and 70 mA (*above*) and from a thin-layer ECL cell at 4.5 V and 1.5 mA (*below*). Even though it requires far less power the ECL illumination provides a much better display.

Summary. Electrochemiluminescence (ECL) is the luminescence resulting from the reaction between electrochemically generated positive and negative radical ions. In general, ECL is a regenerative process; the products of the chemiluminescence reaction are also the starting compounds for the electrochemical conversion. At Philips Forschungslaboratorium Aachen (PFA) a novel type of ECL cell has been developed: it is a thin-layer cell filled with an electrolyte-free solution of a known luminescent compound (e.g. rubrene) in a weakly polar solvent (e.g. 1,2-dimethoxyethane). The cells can be driven by direct voltage. Above a threshold voltage, both the luminescence intensity and the electrochemical current increase linearly with the voltage, without leading to saturation. The luminescence occurs mainly near the anode and interesting patterns are produced in the solution as a result of mass transfer by electrohydrodynamic convection. The ECL efficiency is about 1% and the cells have a much longer life than conventional electrolytic cells. In combination with a high-pressure liquid chromatograph, ECL cells form very sensitive detectors for some polycyclic aromatic hydrocarbons. The cells can also be used as background illumination in passive display elements, as in LCD digital watches and alarm clocks.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- H. A. Algra & J. M. Robertson:** A FMR study on horizontally dipped LPE grown (La,Ga):YIG films. *J. appl. Phys.* **50**, 2173-2175, 1979 (No. 3, Part II). *E*
- J. A. Appels, M. G. Collet, P. A. H. Hart, H. M. J. Vaes & J. F. C. M. Verhoeven:** Thin layer high-voltage devices (RESURF devices). *Philips J. Res.* **35**, 1-13, 1980 (No. 1). *E*
- P. Ashburn, C. J. Bull & J. R. A. Beale:** The use of the electron-beam-induced current mode of the SEM for observing emitter/collector pipes in bipolar transistors. *J. appl. Phys.* **50**, 3472-3477, 1979 (No. 5). *R*
- G. Bartels, D. Mateika & J. M. Robertson:** Preparation of barium lead hexa-aluminate single crystal layers by the liquid phase epitaxy technique. *J. Crystal Growth* **47**, 414-420, 1979 (No. 3). *H, E*
- J. R. A. Beale:** Stored energy transmission for road vehicles. *Electronics and Power* **25**, 323-328, 1979 (May). *R*
- V. Belevitch:** Synthesis of nonconstant-impedance filter pairs. *Philips J. Res.* **35**, 104-121, 1980 (No. 2). *B*
- V. Belevitch & R. R. Wilson** (NKF Kabel, Delft): Cross-talk in twisted multiwire cables. *Philips J. Res.* **35**, 14-58, 1980 (No. 1). *B*
- F. Berz:** Comment on 'A note on the assumption of quasiequilibrium in semiconductor junction devices'. *J. appl. Phys.* **50**, 4479-4481, 1979 (No. 6). *R*
- J. W. M. Biesterbos:** Properties of amorphous rare earth-transition metal thin films relevant to thermomagnetic recording. *J. Physique* **40**, C5/274-279, 1979 (Colloque C5). *E*
- J. W. M. Biesterbos, A. G. Dirks, M. A. J. P. Farla & P. J. Grundy** (University of Salford): The microstructure and magnetic properties of amorphous Tb-Fe(O₂) thin films. *Thin Solid Films* **58**, 259-263, 1979 (No. 2). *E*
- P. Biloen, R. Bouwman, R. A. van Santen** (all with Koninklijke/Shell-Laboratorium, Amsterdam) & **H. H. Brongersma:** Surface composition of some selected group VIII/Au and group VIII/Sn alloys. *Appl. Surface Sci.* **2**, 532-542, 1979 (No. 4). *E*
- G. M. Blom:** Single Te films and Te trilayers for optical recording. *Appl. Phys. Letters* **35**, 81-83, 1979 (No. 1). *N*
- H. J. A. Bluysen** (University of Nijmegen), **L. J. van Ruyven & F. Williams** (University of Delaware, Newark): Effects of quantum confinement and compositional grading on the band structure of heterojunctions. *Solid-State Electronics* **22**, 573-579, 1979 (No. 6). *E*
- H. Bosma:** Over de invloed van micro-elektronica. *Informatie* **21**, 359-364, 1979 (No. 6). *E*
- P. W. J. M. Boumans & M. Bosveld:** A tentative listing of the sensitivities and detection limits of the most sensitive ICP lines as derived from the fitting of experimental data for an argon ICP to the intensities tabulated for the NBS copper arc. *Spectrochim. Acta* **34B**, 59-72, 1979 (No. 2/3). *E*
- P. C. P. Bouten & A. R. Miedema:** On the stable compositions in transition metal-nitrogen phase diagrams. *J. less-common Met.* **65**, 217-228, 1979 (No. 2). *E*
- G. Bouwhuis & S. Wittekoek:** Automatic alignment system for optical projection printing. *IEEE Trans.* **ED-26**, 723-728, 1979 (No. 4). *E*
- S. D. Brotherton, P. Bradley & J. Bicknell:** Electrical properties of platinum in silicon. *J. appl. Phys.* **50**, 3396-3403, 1979 (No. 5). *R*
- M. Brouha, J. E. de Jong & K. J. A. van der Weide:** Experimental verification of finite element analysis on axisymmetric deformation processes. *1979 Manuf. Engng. Trans.* — 7th North Amer. Metalworking Res. Conf. Proc. (Ann Arbor), pp. 57-64. *E*
- R. Bruno, U. Brombach & B. Steinmüller:** On calculating heating and cooling requirements. *Energy and Buildings* **2**, 197-202, 1979 (No. 3). *A*

- A. L. J. Burgmans:** Orientation and alignment relaxation of the 3^2P states of sodium measured at high vapor densities.
Phys. Rev. A **19**, 1954-1959, 1979 (No. 5). *E*
- K. H. J. Buschow & N. M. Beekmans:** Thermal stability and electronic properties of amorphous Zr-Co and Zr-Ni alloys.
Phys. Rev. B **19**, 3843-3849, 1979 (No. 8). *E*
- K. H. J. Buschow & A. M. van der Kraan** (Interuniversitair Reactor Instituut, Delft): Magnetic properties of amorphous $\text{Th}_{1-x}\text{Fe}_x$ alloys.
Phys. Stat. sol. (a) **53**, 665-669, 1979 (No. 2). *E*
- K. L. Bye & R. S. Cosier:** An X-ray double crystal topographic assessment of defects in quartz resonators.
J. Mat. Sci. **14**, 800-810, 1979 (No. 4). *R*
- W. J. Dallas:** Artifact-free region-of-interest reconstruction from coded-aperture recordings.
Optics Comm. **30**, 155-158, 1979 (No. 2). *H*
- H. Dammann:** Spectral characteristic of stepped-phase gratings.
Optik **53**, 409-417, 1979 (No. 5). *H*
- M. Davio & A. Thayse:** Implementation and transformation of algorithms based on automata, Part I: Introduction and elementary optimization problems.
Philips J. Res. **35**, 122-144, 1980 (No. 2). *B*
- E. H. L. J. Dekker, N. J. Wiegman, K. L. L. van Mierloo & R. de Werdt:** Improved propagation, stretching and annihilation of magnetic bubbles in current-access devices.
J. appl. Phys. **50**, 2277-2279, 1979 (No. 3, Part II). *E*
- P. Delsarte, Y. V. Genin & Y. G. Kamp:** Planar least squares inverse polynomials: Part I — Algebraic properties.
IEEE Trans. CAS-26, 59-66, 1979 (No. 1). *B*
- A. G. Dirks & J. R. M. Gijsbers:** Crystallization of amorphous rare earth-iron and transition metal-boron thin films.
Thin Solid Films **58**, 333-337, 1979 (No. 2). *E*
- H. Dötsch & P. E. Wigen:** Observation of the radial oscillation of a bubble lattice.
Solid State Comm. **30**, 611-616, 1979 (No. 10). *H*
- P. van Engelen & S. G. Sie:** Electron-nuclear double resonance of cubic $^{55}\text{Mn}^{2+}$ in GaP.
Solid State Comm. **30**, 515-516, 1979 (No. 8). *E*
- E. Fischmann:** A second 'original' first page of Röntgen's manuscript 'Ueber eine neue Art von Strahlen'.
Brit. J. Radiology **52**, 595-597, 1979 (No. 619). *E*
- B. Gibson:** Equalization design for a 600 MBd quantized feedback PCM repeater.
IEEE Trans. COM-27, 134-142, 1979 (No. 1). *R*
- J.-M. Goethals & J. J. Seidel** (Eindhoven University of Technology): Spherical designs.
Relations between combinatorics and other parts of mathematics, Proc. Symp. Columbus (Ohio) 1978 (Proc. Symp. Pure Math. **34**), pp. 255-272; 1979. *B*
- R. G. Gossink, H. A. M. de Grefte & H. W. Werner:** SIMS analysis of aqueous corrosion profiles in soda-lime-silica glass.
J. Amer. Ceramic Soc. **62**, 4-9, 1979 (No. 1/2). *E*
- F. Grainger & I. G. Gale:** Direct analysis of solid cadmium mercury telluride by flameless atomic absorption using interactive computer processing.
J. Mat. Sci. **14**, 1370-1374, 1979 (No. 6). *R*
- H.-J. Hagemann & H. Ihrig:** Valence change and phase stability of $3d$ -doped BaTiO_3 annealed in oxygen and hydrogen.
Phys. Rev. B **20**, 3871-3878, 1979 (No. 9). *A*
- P. A. H. Hart, T. van 't Hof & F. M. Klaassen:** Device down scaling and expected circuit performance.
IEEE J. SC-14, 343-351, 1979 (No. 2). *E*
- E. E. Havinga & P. van Pelt:** Intermolecular charge transfer, studied by electrochromism of organic molecules in polymer matrices.
Mol. Cryst. liq. Cryst. **52**, 145-155, 1979 (No. 1-4). *E*
- E. E. Havinga & P. van Pelt:** Electrochromism of organic dyes in polymer matrices.
Electro-optics and dielectrics of macromolecules and colloids, ed. B. R. Jennings, pp. 89-97; Plenum Press, New York 1979. *E*
- P. H. van der Heijden** (user of the apparatus here described), **F. Meyer, A. H. T. Sanders, H. A. J. Sanders, H. E. M. Mélotte* & H. Bouma*** (* Institute for Perception Research, Eindhoven): Writing apparatus controlled by head movements for motor handicapped people.
Appl. Ergonomics **10**, 39-45, 1979 (No. 1). *E*
- D. Hennings:** Phase equilibria and thermodynamics of the double oxide phase Cu_3TiO_4 .
J. solid State Chem. **31**, 275-279, 1980 (No. 2). *A*
- H. Hieber & K. Pape:** Umkristallisation dünner Goldschichten.
Z. Metallk. **70**, 459-466, 1979 (No. 7). *H*
- W. K. Hofker, D. P. Oosthoek, G. E. J. Eggermont, Y. Tamminga** (all with Philips Research Laboratories, Amsterdam Department) & **W. T. Stacy:** Laser irradiation of silicon containing misfit dislocations.
Appl. Phys. Letters **34**, 690-692, 1979 (No. 10). *E*
- L. Hollan & J. M. Durand:** Fast growth in GaAs VPE at low temperature and high partial pressures.
J. Crystal Growth **46**, 665-670, 1979 (No. 5). *L*
- L. Hollan, J. C. Tranchart & R. Memming:** Interpretation of selective etching of III-V compounds on the basis of semiconductor electrochemistry.
J. Electrochem. Soc. **126**, 855-859, 1979 (No. 5). *L*
- H. Ihrig & M. Klerk:** Visualization of the grain-boundary potential barriers of PTC-type BaTiO_3 ceramics by cathodoluminescence in an electron-probe microanalyser.
Appl. Phys. Letters **35**, 307-309, 1979 (No. 4). *A, E*
- L. Jacomme:** Pulse broadening in multilayered fibers.
Philips J. Res. **35**, 157-167, 1980 (No. 2). *L*

- T. Karstens, K. Kobs & R. Memming:** The triplet state in trans \rightarrow cis-photoisomerization of thioindigo. Ber. Bunsen-Ges. Phys. Chemie 83, 504-510, 1979 (No. 5). *H*
- A. J. R. de Kock, W. T. Stacy & W. M. van de Wijert:** The effect of doping on microdefect formation in as-grown dislocation-free Czochralski silicon crystals. Appl. Phys. Letters 34, 611-613, 1979 (No. 9). *E*
- K. Kosai, B. J. Fitzpatrick, H. G. Grimmeiss, R. N. Bhargava & G. F. Neumark:** Shallow acceptors and *p*-type ZnSe. Appl. Phys. Letters 35, 194-196, 1979 (No. 2). *N*
- H. K. Kuiken & P. J. Roksnoer:** Analysis of the temperature distribution in FZ silicon crystals. J. Crystal Growth 47, 29-42, 1979 (No. 1). *E*
- K. M. Lüdeke, J. Köhler & J. Kanzenbach:** A new radiation balance microwave thermograph for simultaneous and independent temperature and emissivity measurements. J. Microwave Power 14, 117-121, 1979 (No. 2). *H*
- B. J. de Maagt & G. C. J. Rouweler** (Philips Lighting Division, Eindhoven): The reaction of bromine and oxygen with a tungsten surface, Part II: Measurements on the influence of hydrogen on the reaction rate. Philips J. Res. 35, 95-103, 1980 (No. 2).
- J. Magarshack, R. Dessert, P. Harrop, V. Pauker* & R. Mauffret*** (* RTC, Suresnes): Récepteur à très faible bruit de télédiffusion par satellite à 12 GHz. C.R. Coll. Int. Espace, télécommunications spatiales et radiodiffusion par satellite — les objectifs de la prochaine décennie, Toulouse 1979, pp. 177-191. *L*
- R. Memming, F. Schröppel & U. Bringmann:** Sensitized oxidation of water by tris(2,2'-bipyridyl) ruthenium at SnO₂ electrodes. J. electroanal. Chem. 100, 307-318, 1979 (No. 1/2). *H*
- A. R. Miedema & K. A. Gingerich** (Texas A and M University, College Station): On the formation enthalpy of metallic dimers. J. Physics B 12, 2081-2095, 1979 (No. 13). *E*
- A. R. Miedema & K. A. Gingerich** (Texas A and M University, College Station): On the enthalpy of formation of diatomic intermetallic molecules. J. Physics B 12, 2255-2270, 1979 (No. 14). *E*
- A. Mitonneau & A. Mircea:** Auger de-excitation of a metastable state in GaAs. Solid State Comm. 30, 157-162, 1979 (No. 3). *L*
- M. Monnier:** Utilisation du film polyimide dans l'interconnexion de semi-conducteurs. G.P.C.P. Journées d'études des 24-25-26 avril 1978, pp. 49-54. *L*
- G. J. Naaijer:** Instalaciones solares fotovoltaicas para bombeo de agua. Mundo Electrónico No. 85, pp. 49-66, 1979. *L*
- G. F. Neumark:** Analysis of the pretransition range of the metal-insulator transition in doped semiconductors. Phys. Rev. B 20, 1519-1526, 1979 (No. 4). *N*
- A. Nicia:** Impulse response of a step-index fiber including a gap. Philips J. Res. 35, 145-156, 1980 (No. 2). *E*
- A. E. Pannenberg** (Philips Board of Management, Eindhoven): Industriële innovatie. Ned. T. Natuurk. A 45, 43-45, 1979 (No. 2).
- E. Papanikolaou** (Philips' Glass Development Centre, Eindhoven) & **A. J. H. Wachters** (Philips Information Systems and Automation Department, Eindhoven): A theory of bubble growth at chemical equilibrium with application to the hydrogen reboil in fused silica. Philips J. Res. 35, 59-93, 1980 (No. 1).
- R. E. Pearson, G. Payne & D. H. Paul:** An actively broadbanded MIC parametric amplifier. 1979 IEEE MTT-S Int. Microwave Symp. Dig., pp. 501-503. *R*
- J. F. M. Pennings & B. Bosman:** Relaxation of the surface energy of solid polymers. Colloid and Polymer Sci. 257, 720-724, 1979 (No. 7). *E*
- P. Piret:** Causal sliding block encoders with feedback. IEEE Trans. IT-25, 237-240, 1979 (No. 2). *B*
- T. E. Rozzi & G. H. in 't Veld:** Field and network analysis of interacting step discontinuities in planar dielectric waveguides. IEEE Trans. MTT-27, 303-309, 1979 (No. 4). *E*
- L. J. van Ruyven & C. A. J. Ammerlaan** (Universiteit van Amsterdam): Electron trapping in neutron transmutation doped silicon. Appl. Phys. Letters 34, 632-634, 1979 (No. 10). *E*
- S. K. Salmon:** Practical aspects of surface acoustic wave oscillators. 1979 IEEE MTT-S Int. Microwave Symp. Dig., pp. 165-167. *R*
- H. Schippers:** Analytical and numerical results for the non-stationary rotating disk flow. J. Engng Math. 13, 173-191, 1979 (No. 2). *E*
- P. J. Severin:** Calorimetric measurements of weakly absorbing materials: theory. Appl. Optics 18, 1546-1554, 1979 (No. 10). *E*
- P. J. Severin:** Calorimetriscche meting van de absorptiecoëfficiënt in optisch glas en in optische fibers. Ned. T. Natuurk. A 45, 75-78, 1979 (No. 2). *E*
- A. Shaulov, M. I. Bell** (National Bureau of Standards, Washington D.C.) & **W. A. Smith:** Direct measurement of pyroelectric figures of merit of proper and improper ferroelectrics. J. appl. Phys. 50, 4913-4919, 1979 (No. 7). *N*
- I. J. Stemp, K. H. Nicholas & H. E. Brockman:** Automatic testing and analysis of misregistrations found in semiconductor processing. IEEE Trans. ED-26, 729-732, 1979 (No. 4). *R*
- E. F. Stikvoort:** Resistive gate DMOST for low distortion mixing. Solid-State Electronics 22, 595-598, 1979 (No. 6). *E*

- S. Strijbos, A. Broese van Groenou & P. A. Vermeer** (Delft University of Technology): Recent progress in understanding die compaction of powders. *J. Amer. Ceramic Soc.* **62**, 57-59, 1979 (No. 1/2). *E*
- M. Tasto, M. Felgendreher, W. Spiesberger & P. Spiller**: Comparison of manual versus computer determination of left ventricular boundaries from x-ray cine-angiocardiograms. Roentgen-video-techniques for dynamic studies of structure and function of the heart and circulation (2nd Int. Workshop Conf., Kiel 1976), ed. P. H. Heintzen & J. H. Bürsch, pp. 168-183; Thieme, Stuttgart 1978. *H*
- J. B. Theeten & D. E. Aspnes** (Bell Laboratories, Murray Hill, N.J.): The determination of interface layers by spectroscopic ellipsometry. *Thin Solid Films* **60**, 183-192, 1979 (No. 2). *L*
- N. C. de Troye**: Morphologie einiger bekannter bipolarer digitaler Schaltungen und Prozesse. *Nachrichtentechn. Z.* **32**, 382-387, 1979 (No. 6). *E*
- M. J. Underhill & P. A. Jordan**: Split-loop method for wide-range frequency synthesiser with good dynamic performance. *Electronics Letters* **15**, 391-393, 1979 (No. 13). *R*
- M. J. Underhill & R. I. H. Scott**: Wideband frequency modulation of frequency synthesisers. *Electronics Letters* **15**, 393-394, 1979 (No. 13). *R*
- C. H. F. Velzel**: Het gedrag van halfgeleider-diode-lasers bij optische terugkoppeling. *Ned. T. Natuurk. A* **45**, 54-58, 1979 (No. 2). *E*
- J. A. Th. Verhoeven & H. van Doveren**: A vapour collect method based on Auger spectroscopy. *Mikrochim. Acta* 1979 I, 331-344 (No. 5/6). *E*
- J. Visser & J. J. Scheer**: Twenty-kelvin cryopumping in magnetron sputtering systems. *J. Vac. Sci. Technol.* **16**, 734-737, 1979 (No. 2). *E*
- J. J. Vrakking & A. Kroes**: Ion-induced Auger electron emission of Mg, Al and Si as a function of ion energy. *Surface Sci.* **84**, 153-163, 1979 (No. 1). *E*
- Q. H. F. Vrehan**: Superfluorescence experiments. *Trends in physics, 1978* (Papers 4th Gen. Conf. Eur. Phys. Soc., York), pp. 95-97; 1979. *E*
- H. J. de Wit, C. Wijenberg & C. Crevecoeur**: Impedance measurements during anodization of aluminum. *J. Electrochem. Soc.* **126**, 779-785, 1979 (No. 5). *E*

Contents of Philips Telecommunication Review **38**, No. 1, 1980:

- H. Bouwman & U. Rothgordt**: Electrostatic printing in facsimile terminals (pp. 1-6).
- L. Kool**: The Scribophone: a graphic telecommunication system (pp. 7-10).
- A. Rijbroek & J. Drupsteen**: Higher-order PCM multiplex systems (pp. 11-22).
- K. Fisher**: M294, an economic radiotelephone (pp. 23-26).
- J. de Boer & C. Hooijkamp**: The required load capacity of FDM multi-channel amplifiers if single-channel peak limiting is employed (pp. 27-36).

Contents of Philips Telecommunication Review **38**, No. 2, 1980:

- E. C. Priebee**: Air traffic control system for Singapore (pp. 41-60).
- G. C. M. de Koning & P. M. G. Westenberg**: Projecting SPC exchanges by computer (pp. 61-75).
- A. L. W. Bloemendaal**: The effect of traffic control on motor traffic efficiency (pp. 76-82).
- M. M. Jung**: Busy period distribution in an SPC processor having a clock-pulse operated gate (pp. 84-89).

Contents of Electronic Components and Applications **2**, No. 2, 1980:

- J. A. Houldsworth & W. B. Rosink**: Introduction to PWM speed control system for 3-phase AC motors (pp. 66-79).
- Bucket-brigade delay-line enhances sound reproduction (pp. 80-82).
- H. W. Evers**: Mains pollution caused by domestic appliances, Part 2 — Harmonic distortion (pp. 83-90).
- T. G. Giles**: Versatile LSI frequency synthesiser system (pp. 91-105).
- S. Y. Lau**: Integrated Schottky logic gate array (pp. 106-114).
- B. P. Bahnsen**: Digital control of radio and audio equipment, Part 6 — Voltage-controlled tuning of AM radios (pp. 115-125).

Computer-aided research on multiwire telephone cables

J. Veldhuis

Informed forecasts of the future of telecommunications include the emergence of an 'Integrated Services Digital Network' (ISDN). The idea of such an all-embracing network, in which text, sound and image would be transmitted by purely digital methods, may still seem rather visionary, at least on a world-wide scale, but it powerfully influences current thinking on existing telecommunication systems. The adaptation of these systems for the most economic yet properly engineered transition to such a 'network of the future' is a challenging topic of Philips research. The article below deals with theoretical work that is particularly relevant to this research. The main question was the extent to which the existing telephone system, with its conventional copper cable, can be adapted for the digital transmission of speech and the provision of a wide variety of new data services. The analytical method presented here, based on an advanced theoretical study of telephone cables with twisted wires, clearly illustrates that simulation by software provides a most promising research method for the study of large systems.

The copper cable and digital communication

The conventional copper cable used in existing telephone networks is once again the subject of a great deal of technological and scientific research. This tends to be overshadowed somewhat by such a revolutionary innovation as the glass-fibre cable, whose introduction has of course attracted a great deal of attention. As is often the case, however, this new development does not signify the immediate demise of the existing technology^[1], nor even that all the research on it is complete. Much practical and theoretical work remains to be done.

Research on the copper cable has two objectives: making better telephone cables and making better use of existing telephone networks. The first objective is of major importance to the cable manufacturer. In his view, a telephone cable of high quality should give no more than 1 to 2 dB attenuation per kilometre of cable (at frequencies up to 100 kHz) in each of the grouped transmission channels. 'High quality' also implies that interference between the signals travelling along the

individual transmission channels is negligible. The manufacturer therefore seeks to construct his cables so as to minimize signal leakage between neighbouring channels ('crosstalk'). In addition to low attenuation in the direction of transmission, his aim is to achieve high *isolation between* the transmission channels (more than 40 dB at frequencies up to 100 kHz in cables up to 5 kilometres long).

The second objective, making better use of existing cables, is of particular interest to those who wish to transmit signals other than analog speech signals on the existing networks. A good example of such extended use is the transmission of a wide variety of data with the aid of digital signals. Indeed, this is

^[1] The Netherlands Postal, Telegraph and Telephone service, which is responsible for the Dutch telephone system, requires for reasons of operational reliability that subscriber connections should be energized from the telephone exchanges. The glass-fibre cable does not at the moment appear to be suitable for this power-supply function. Research is however being done in this field (R. C. Miller and R. B. Lawry, Bell Syst. tech. J. 58, 1735, 1979). It is not yet entirely clear how the glass-fibre cable will ultimately be used, along with the copper cable, in local telephone networks.

already happening, though only on a small scale. Individual users who have an appropriate terminal can rent telephone lines to transmit their own data to a central computer, and receive data from the computer. In future all the ordinary subscribers may be offered digital communication facilities through the national telephone network.

A consequence of the new applications will be that the frequencies of the electrical signals in the existing local network will extend over a larger bandwidth than has so far been necessary. This will become most apparent when eventually the analog speech signals of subscribers are also converted into digital signals. Such a situation will require transmission channels with a bandwidth of the order of 100 kHz, some 20 times the bandwidth required for analog speech signals alone (4 kHz).

A disadvantage of such a larger bandwidth is that it entails greatly increased crosstalk between the transmission channels, owing to the higher frequencies in the signal. The attenuation and phase shift also deteriorate rapidly as the frequency increases. Consequently, the signals with a broad frequency spectrum will be much more seriously mutilated than the present analog speech signals, with their bandwidth of only 4 kHz.

The only real limiting factors in the plans to provide telephone subscribers with digital communication and new data-communication services are the existing network and its cables. It is not easy to obtain an overall picture of the limitations introduced into a cable network because of crosstalk and attenuation. Making a complete series of tests for every cable in the network would be an almost impossible task.

A simpler and much cheaper method of investigating the effect of cable limitations is to use a computer to simulate the telephone network as a system, with the emphasis on transmission and crosstalk. A simulation program of this type is based on a formal model of the cable taken as an elementary structural unit. Network calculations will then give a good simulation of the transmission behaviour of the cable and the interference due to crosstalk. Worst-case situations in the cable network can then be determined.

To keep the simulation flexible the software should be organized in modules, in such a way that the individual blocks corresponding to parts of the network, e.g. cables, switchgear, transmitters, etc., can easily be included.

The treatment in this article is confined to the cable block; a substantial part of the software for this block has now been completed. The calculations for this block are made at three hierarchical levels, each with its own cable model. At each level the program takes

particular input data, and employs the model to arrive at a result at the appropriate level. These results in turn form the block of input data at the adjacent level. This division into levels simplifies the cable calculations, which is important for efficiently debugging the programs in the initial phase.

In addition to these three levels, the schematic arrangement adopted has a fourth level. At this level the *complete* telephone network can be simulated by a complicated product of various transfer matrices whose characteristic data come from models of the three main parts of the network (cables, switchgear, transmitters/receivers). Improvement and expansion of the communication facilities will eventually be attained at the fourth level. For the work described in this article the fourth level only acts as the environment for the cable-simulation program. This simulation program has been completed at the *third* level.

In the sections that follow, the geometry of the multiwire telephone cables used in the Netherlands will first be described, and a number of interactions between the conductors will be discussed. Attention will then be turned to the three levels in the model calculations, in particular to finding a cable model and deriving from it expressions for the primary cable parameters at the first level. These form the basis for the network calculations, which eventually result in signal-transfer matrices. The final section deals with the software developed and the values of the primary and secondary cable parameters calculated with this software. A detailed example of one program module is included, mainly to demonstrate the functional structure of the program and to show that certain general quality criteria for software are satisfied. The article concludes with a comparison between values that we have calculated and the results of cable calculations and measurements carried out elsewhere.

Geometry and interference

Conductor configuration

The geometry of a multiwire telephone cable is mainly determined by the distribution of the wires over the cross-section of the cable. *Fig. 1* shows the cross-section of a commonly used type of telephone cable in the Netherlands. The cable contains a large number of conductor wires (there may be several hundred). The wires are distributed in groups of four, called 'quads', in concentric layers. There can be as many as six or seven layers, including the central part of the cable. From the inside outwards the number of quads increases by six per layer. In each quad the four wires are twisted together at a particular pitch,

called the group pitch. The wires do not therefore run parallel to the axis of the cable but are wound concentrically in the form of interwoven helices. In their turn the quads are wound as helices in the individual layers with a particular pitch per layer, called the layer pitch. The layer pitches differ from one another, and the group pitches also differ for adjacent quads in the same layer.

This relatively complicated structure with group and layer pitches has the advantage that if the relations between the pitches are properly chosen the crosstalk effects between pairs of wires can be greatly reduced. This is true whether the crosstalk is due to inductive or capacitive coupling. Because the wires are twisted, the structure is compact and strong — another advantage of some importance.

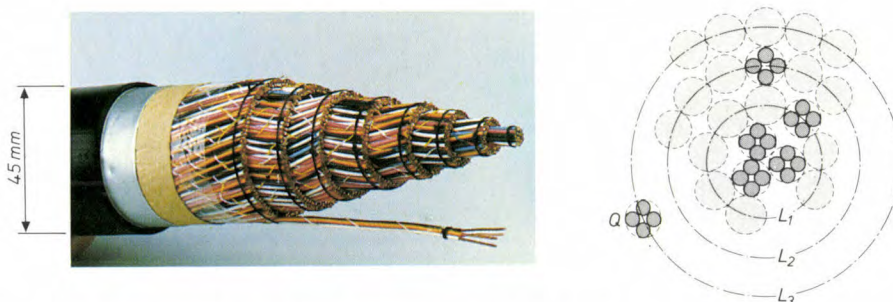


Fig. 1. A telephone cable with quad geometry (made by NKF Kabel B.V., Waddinxveen, The Netherlands). A quad (Q) consists of four conductors. The central part of the cable has three quads; the next layer (L1) has 9 quads, then there is a layer with 15 quads (L2), a layer with 21 quads (L3), and so on. The entire configuration is contained inside a plastic-encapsulated screening shield of aluminium strip. The type shown here is widely used in the Netherlands for local networks as the cable between exchanges and subscriber groups. (The cable sample has been made available by the Creative Services and Publicity Division of NKF Kabel B.V.)

Wires within a quad, and also quads themselves, can be used for making a transmission channel (in telephone cables often called a circuit). It is possible to use a wire for more than one transmission channel at the same time (a familiar example is given in *fig. 2*). In this case one quad and a conducting shield form four transmission channels, which can operate simultaneously and independently of each other.

In cross-section the wires of a quad form the corners of a square. Two circuits, the side circuits, each consist of two conductors diagonally opposite each other. Another circuit, the phantom circuit, is obtained by taking each side pair as a single conductor. The fourth circuit, the asymmetrical circuit, is produced by taking the shield as one conductor and the full quad as the other.

The case of *fig. 2* is based on the general proposition that a cable with N conductors inside a conducting shield can provide a total of N independent transmission channels^[2]. The transmission channels in a telephone cable are not chosen arbitrarily, of course.

The aim is to select the wires for the required transmission channels in such a way as to achieve optimum transmission with minimum crosstalk in the cable.

Many telephone subscribers in the Netherlands are connected to the system by a single quad. The telephone set is connected to the network by one of the two side circuits (*fig. 2*). The remaining connections are held in reserve.

In this article, calculating the transmission or crosstalk in a cable is taken to refer to the transmission along or crosstalk between the individual circuits just described. The electrical quantities such as potential, current and also the cable parameters in the rest of this article always refer to a particular circuit, unless it is explicitly stated that the quantity relates to a particular wire.

The electrical quantities in the circuits and the corresponding quantities for the wires are linearly related in a way that can be described by a matrix equation. In the case of *fig. 2*, for example, there are four circuit potentials (V_{S1} , V_{S2} , V_F , V_A), referring to the two side circuits, the phantom circuit, and the asymmetrical circuit in that order. The matrix equation is:

$$\begin{bmatrix} V_{S1} \\ V_{S2} \\ V_F \\ V_A \end{bmatrix} = \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ V_4 \end{bmatrix},$$

where V_1 , V_2 , V_3 and V_4 are the potentials of the individual wires (defined with respect to the sheath). The same matrix applies to the relation between the charges on the circuits and wires. The four wires and the shield are shown as perfect conductors in *fig. 2*; this assumption is made when the capacitances of circuits (and also between circuits) are calculated, as ratios of charges and potentials. On the other hand, finite conductivity is assumed when dealing with the currents and voltages in the longitudinal direction of a cable; this assumption enables the resistances in the longitudinal direction

[2] W. Klein, *Die Theorie des Nebensprechens auf Leitungen*, Springer, Berlin 1955.

to be derived in addition to the self-inductances and the mutual inductances. This dual approach to the cable calculations will be explained in more detail later.

Interactive effects

Although the electrical phenomena in multiwire telephone cables are described by Maxwell's classical equations, it is not so easy to obtain a clear picture of all the possible interactive effects inside and between the conductors. The principal effects are listed in

reduced) most strongly on the side of the cross-section that lies closest to an adjacent conductor. It is related to the classical skin effect, in which each conductor is affected by its own magnetic field rather than the field from adjacent conductors. Consequently, for the skin effect, the current density in the cross-section is only radially dependent, with a maximum at the outside of the conductor. If the smallest spacing between the conductors in a cable is at least five times the radius of the conductor, the proximity effect is negligible in comparison with the skin effect. In the determination of

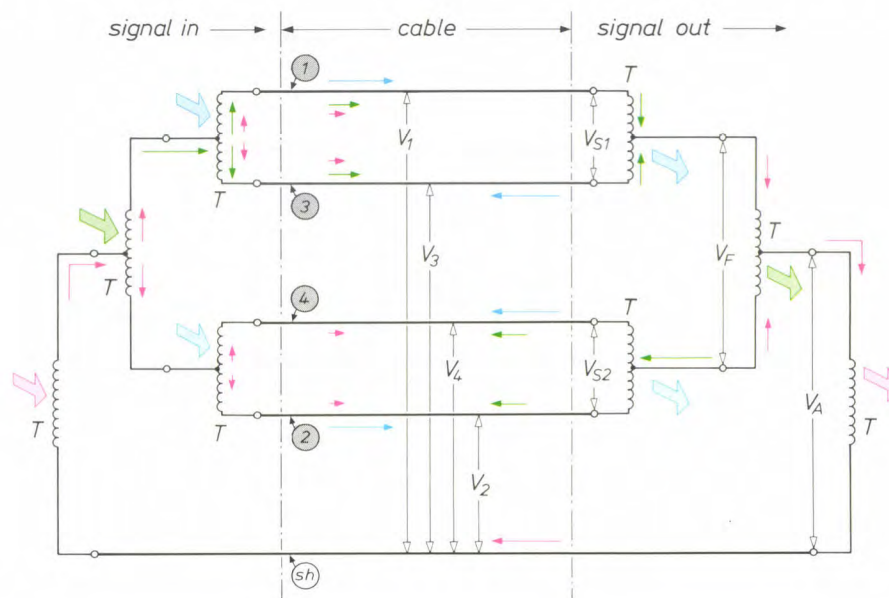


Fig. 2. Simple example of commonly used conductors in a multiwire telephone cable. The cable contains four wires (numbered), which are screened by a conducting shield (*sh*). Four independent circuits are formed: the two side circuits (*blue*), the phantom circuit (*green*), and the asymmetrical circuit (*red*). Each conductor is used for three signal currents. The input signals are supplied via the primary windings (not shown) of four input transformers *T*. The secondary windings of these transformers are connected by terminals to the conductors and the shield. In a similar way the signals at the end of the cable are delivered to the outside world via terminals and four output transformers *T* (whose secondary windings are not shown). Only one signal is effective on a particular transmission channel; any other signals cancel each other out by balancing. V_{S1} , V_{S2} circuit potential of side circuit. V_F circuit potential of phantom circuit. V_A circuit potential of asymmetrical circuit. V_1, \dots, V_4 conductor potential (defined with respect to *sh*).

Table I. Some of these effects increase strongly with frequency. An example is the proximity effect in the conductors — a current-concentration phenomenon that affects the resistances and self-inductances of conductors as well as the mutual inductances between conductors; the effect increases strongly with frequency [4].

This effect, highly undesirable at high frequencies in telephone cables, is an asymmetry in the current distribution over the cross-section of a cable conductor, caused by the magnetic field of the currents in adjacent conductors. The essential feature is that the current density in the conductor is increased (or

capacitances in telephone cables, which is essentially an electrostatic-field problem, a second proximity effect has to be taken into account. This is an asymmetry in the charge distribution (compare with the current distribution above) on a conductor as a consequence of charges in adjacent conductors. To avoid confusion this effect is referred to in Table I as the 'electrostatic proximity effect' with the symbol P_q . All these effects ultimately contribute to an increase in the attenuation and phase shift in the conductors, measured per metre of cable. The proximity effect P and the skin effect S increase with the signal frequency, but P_q is independent of frequency.

H. E. Martin's cage effect ^[3], mentioned in Table I, also increases the attenuation per metre; more precisely, it increases the effective capacitance per metre of a circuit. The effect arises because each circuit in a telephone cable consisting of twisted wires is effectively enclosed by a Faraday cage. This cage — not to be confused with the cable shield — is formed by the group of other circuits most closely surrounding the circuit under consideration. All these cages behave like conducting shields at the same constant potential. If the wires forming a cage are twisted together cor-

separation'. If the cage encloses two circuits instead of one — this could happen if non-ideal twisting of the wires effectively makes one of the adjacent circuits not part of the cage — then the induction charge will also appreciably increase the capacitive coupling between *both* of the enclosed circuits. V. Belevitch has recognized this extension of the effect, and has generalized Martin's theory and at the same time succeeded in finding an explanation for the occurrence of certain non-ideal twisting pitches ^[6], an effect known in practice but not really understood previously.

Table I. Influence of skin effect, shield, proximity effects and cage effect on the primary cable parameters in a telephone cable with the geometry given in fig. 1. *S* skin effect. *P* proximity effect (current). *P_q* electrostatic proximity effect (charge). *K* cage effect as described by Martin ^[3]. A separate column indicates which of the primary cable parameters are frequency-dependent and which are not.

Primary cable parameter of circuit <i>m</i> , or between circuits <i>m</i> and <i>n</i> (defined per metre)	Change due to					frequency dependent
	<i>S</i>	Shield	<i>P</i>	<i>P_q</i>	<i>K</i>	
Resistance <i>R_m</i>	Yes	Yes	Yes	No	No	Yes
Self-inductance <i>L_m</i>	Yes	Yes	Yes	No	No	Yes
Capacitance <i>C_m</i>	^[a] No	Yes ^[c]	No	Yes	Yes	No
Shunt conductance <i>G_m</i>	^[a] No	Yes ^[c]	No	Yes	Yes	Yes
Mutual resistance <i>R_{mn}</i>	^[b] Yes	Yes	Yes	No	No	Yes
Mutual inductance <i>L_{mn}</i>	^[b] Yes	Yes	Yes	No	No	Yes
Capacitance <i>C_{mn}</i>	^[a] No	Yes ^[c]	No	Yes	Yes	No
Conductance <i>G_{mn}</i>	^[a] No	Yes ^[c]	No	Yes	Yes	Yes

^[a] Conductance and capacitance are related by the fixed relation $G = \omega C \tan \delta$, where $\omega (= 2\pi f)$ is the angular frequency and δ is the loss angle of the material between the conductors.

^[b] $R_{mn} + j\omega L_{mn} = U_m/I_n$, the ratio of the induced voltage per metre U_m to the inducing current I_n . Owing to the screening effect of the shield, and also because of the proximity effect, this ratio also has a real part, R_{mn} , here called mutual resistance (per metre).

^[c] The change brought about by the screening action of the cable shield only occurs in so far as the circuit is not screened by the cage effect *K*.

rectly, their spatial location is such that the voltages at the connection terminals of all the circuits belonging to a particular cage do not change when the central circuit in the cage is energized. (The surrounding circuits and the central circuit are then decoupled electrostatically, making capacitive crosstalk impossible.) Local static-induction charges arise on the cage, and these oppose the field-strength of the central field, which produced them. The corresponding voltage drop between the two conductors of the central circuit in turn produces an increase in the effective capacitance, for constant charge. This increase may amount to 10 to 15%. Also, because of the twisting of the wires — with a periodic variation of the distance between central circuit and cage — the induction in the longitudinal direction of the cable will be alternately strong and weak, resulting in 'longitudinal charge

With all these effects influencing the behaviour of cables, it is not surprising that it is so difficult to calculate the crosstalk between pairs of conductors properly. The amplitude of the crosstalk signals depends closely on the method adopted for twisting the wires in the cable. For these reasons the simple cable models now in use are often no more than rough approximations to the actual cable. The effects that occur in the frequency range already important today cannot be correctly described by these models, and the difficulty is only aggravated at even higher frequencies.

^[3] H.-E. Martin, Die Berechnung der Übertragungseigenschaften symmetrischer Leitungen unter Berücksichtigung des Verdrallungseffektes, Arch. elektr. Übertr. 18, 293-308, 1964.

^[4] See for example P. Grivet, The physics of transmission lines at high and very high frequencies, Vol. 1, Academic Press, London 1970.

^[6] V. Belevitch, On the theory of cross-talk between twisted pairs, Philips Res. Repts 32, 365-372, 1977.

The models

The simulation program on which we are working attempts to take into account as far as possible all the interactive effects between the conductors. A great deal of attention therefore had to be paid to the consequences of the complicated interweaving of the conductors. This was made possible by the theoretical work published in recent years, particularly by Belevitch, working with G. C. Groenendaal and R. R. Wilson [5]-[8]. Their treatment gives particular attention to the two proximity effects and the cage effect.

It provides a better understanding of the frequency dependence of the attenuation and the phase shift, which determines the transmission behaviour.

The approach can also provide a better description of the inductive and capacitive couplings, which determine the crosstalk between the pairs of conductors.

The block structure

Fig. 3 shows the block diagram of the cable-simulation system we have designed, giving the levels at which calculations can be carried out on the cables.

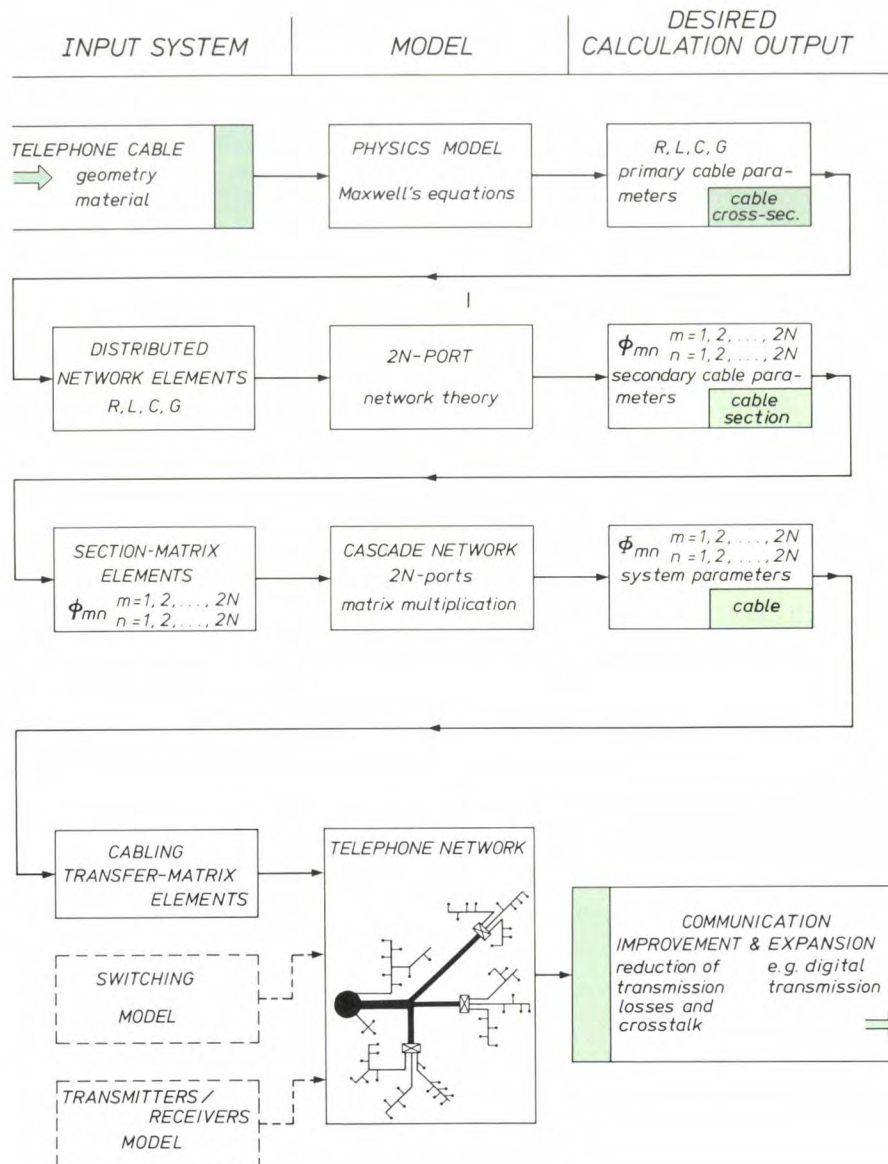


Fig. 3. Block diagram of computer-simulation software for telephone cables. The first three levels represent the calculation of the elements of the transfer matrix — the system parameters ϕ_{mn} for signal transfer and hence the final cable model. These parameters, combined with the models for the switching elements and transmitters/receivers (not dealt with here), will form a program simulating the entire telephone network, shown at the fourth and lowest level. A program of this type gives an efficient analysis of the transmission attenuation and crosstalk in existing telephone systems. *R* resistance per metre. *L* inductance per metre. *C* capacitance per metre. *G* conductance per metre. A 2*N*-port is a network consisting of *N* circuits, each with two connection terminals (a port) at the input and the output. The parameters ϕ_{mn} are defined in the text.

At the upper level the calculation starts with dimensions and materials, the 'simplest' data for establishing the geometry and the electrical properties of a telephone cable. These data are substituted in Maxwell's field equations, which serve as a model for calculating quantities such as the resistance per metre, the self-inductance per metre, etc., i.e. the primary cable parameters R , L , C and G . These quantities characterize the cable for an infinitely small element of length (i.e. at a cross-section); because of the twisting of the wires, they will not be truly constant over the length of the cable but will depend on the location of the cross-section. Owing to the complicated geometry this first step in the calculations, in which the effects in Table I enter into the picture, is especially difficult to carry out; it is the subject of a separate area of research in which a great deal of theoretical work is under way, as mentioned earlier.

At the second level the model is what is called a $2N$ -port, a network circuit built up from the quantities R , L , C and G calculated — at many cross-sections — from the first level, and referred to at this level as 'distributed network elements'. A network of this type can be used to form N independent circuits. A circuit consists of two lines, with a 'port' at the input end and at the output end of the network. In each circuit there are two potential drops: one across the input terminals and one across the output terminals, so that there are also two current levels. The network establishes a linear relation between all the currents and voltages at the input and output. This relation is mostly clearly described by a 'section matrix', a matrix of $2N \times 2N$ elements ϕ_{mn} , the secondary cable parameters.

In the simple case of two conductors without a shield the section matrix consists of four elements. If the conductors are parallel and the primary parameters thus independent of the length coordinate, the linear relation between the current (I_S) and the potential difference (V_S) at the input and the current (I_R) and the potential difference (V_R) at the output of a cable section of length Δz is not difficult to determine, and is:

$$\begin{bmatrix} V_S \\ I_S \end{bmatrix} = \begin{bmatrix} \cosh(\gamma\Delta z) & Z_c \sinh(\gamma\Delta z) \\ \frac{1}{Z_c} \sinh(\gamma\Delta z) & \cosh(\gamma\Delta z) \end{bmatrix} \begin{bmatrix} V_R \\ I_R \end{bmatrix}. \quad (1)$$

It will be evident that the pair (V_R , I_R) can be used in turn as the input signal to the next cable section. The coefficient Z_c is the characteristic impedance, and γ ($\equiv \alpha + j\beta$) is the complex propagation constant. The real part α is known as the attenuation constant and the imaginary part β as the phase constant.

In this case Z_c and γ are given by

$$Z_c = \{(R + j\omega L)/(G + j\omega C)\}^{\frac{1}{2}} \quad (2)$$

and

$$\gamma = \alpha + j\beta = \{(R + j\omega L)(G + j\omega C)\}^{\frac{1}{2}}. \quad (3)$$

The quantities Z_c , α , β are known from conventional cable measurements, and are also frequently referred to as secondary cable parameters (the name reserved in this article for the elements of the section matrix).

The objective at the second level is therefore to calculate the elements ϕ_{mn} . These can be derived from the primary cable parameters by means of ordinary network theory. The $2N \times 2N$ elements characterize a short length of cable (a 'cable section'). In the general case the derivation is much more complicated than in the example given above, in which there is only one current and one potential difference at the input and output.

At the third level the transfer matrix of a complete cable is calculated by treating the cable as a cascaded circuit of the separate cable sections and then multiplying the appropriate section matrices in the proper sequence. This results in the $2N \times 2N$ matrix elements Φ_{mn} , the 'system parameters' of the complete cable.

Calculation of transmission and crosstalk

The primary cable parameters (R , L , C and G in fig. 3) are quantities that only really become electrically significant when we consider a length of cable and not a cable cross-section, and treat that piece of cable as an electrical network. Fig. 4 shows a network that can serve as equivalent circuit for a piece of cable of length Δz . Longer pieces are equivalent to cascaded arrangements of such equivalent circuits.

A telephone cable, especially a modern high-quality cable, generally provides such good transmission channels that the intricate problem of analysing the total network, via the determination of both transmission behaviour and crosstalk, can with advantage be split into two separate problems. The transmission behaviour and the crosstalk are then determined separately. Splitting the total problem in this way makes the calculation much simpler and more efficient. It is possible to make such a split because the energy lost during transmission as crosstalk is at least an order of magnitude smaller than the energy used for the transmission.

[6] V. Belevitch, Theory of the proximity effect in multiwire cables, Philips Res. Repts 32, 16-43 and 96-117, 1977.

See also G. C. Groenendaal, R. R. Wilson and V. Belevitch, Calculation of the proximity effect in a screened pair and quad, Philips Res. Repts 32, 412-428, 1977.

[7] V. Belevitch, R. R. Wilson and G. C. Groenendaal, The capacitance of circuits in a cable with twisted quads, Philips Res. Repts 32, 297-321, 1977.

[8] V. Belevitch and R. R. Wilson, Cross-talk in twisted multiwire cables, Philips J. Res. 35, 14-58, 1980.

The transmission problem of a cable with N conductors inside a conducting shield (sh) thus reduces to the simple transmission problem — although it has to be repeated N times — of a circuit with only two wires in free space (shown in red in fig. 4). The solution to

agitation of voltage and current waves in the entire network. To find a general solution to these equations — thus giving all the secondary cable parameters ϕ_{mn} (fig. 3) — is so massive a task that numerical methods have to be used.

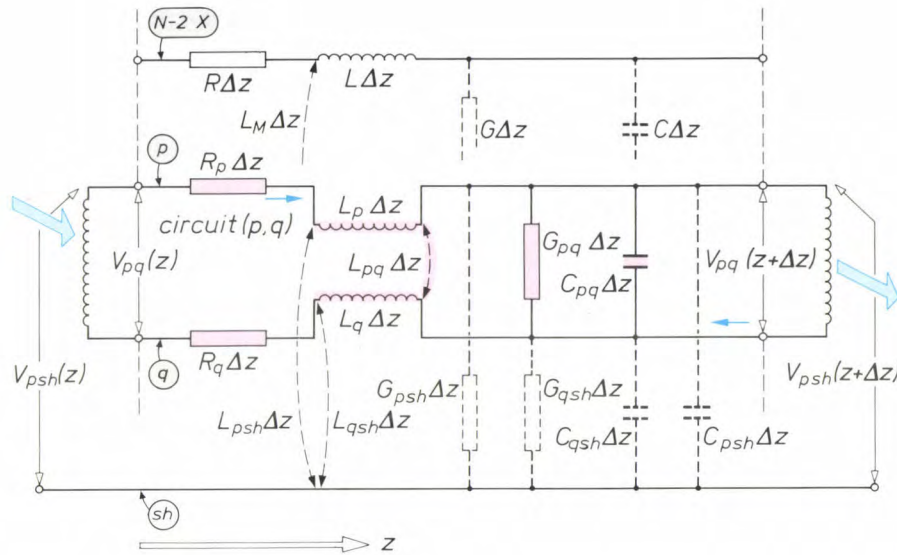


Fig. 4. Simplified equivalent circuit of a short piece of telephone cable (length Δz) in which N connection circuits are formed from N numbered wires inside a shield (sh), which at the same time represents the zero level for all the electric potentials V . The circuit (p, q) consisting of the wires p and q is shown in full. Only the network elements printed in red need be known to calculate the signal transmission in the circuit. The dashed network elements form the 'environment' of the circuit; they are necessary for calculating the crosstalk in (p, q). R, L, C, G are primary parameters, whose significance is defined in fig. 3. L_M mutual inductance per metre. All network elements refer to separate conductors.

this problem is provided by the equations (1), (2) and (3) discussed earlier. In reality, of course, the wires are not in free space, and if this simple calculation of the transmission behaviour using equations (1), (2) and (3) is to be sufficiently accurate, then the primary parameters of this single circuit (the red network elements in fig. 4) must be calculated with due allowance for all the interactive effects in Table I, as they occur in the complete cable. As mentioned earlier, the twisting of the wires makes the primary cable parameters functions of z , the coordinate of length along the cable. The values of the primary constants used in equations (1), (2) and (3) are averages over the length of the cable.

The separate crosstalk problem is a much more difficult problem in network theory, since the 'environment' does have to be taken into account by including all the network elements shown dashed in fig. 4. This would in fact become the problem of solving the generalized telegraphist's equations [9], a large array of coupled differential equations describing the prop-

The quality of modern telephone cables is so good that the coupling between the different circuits is sufficiently small to allow a considerable simplification of the network structure of the 'environment' in fig. 4. In our approach to the calculations this meant that we neglected any feedback effects from a circuit *receiving* interference to the circuit *producing* the interference, as effects of higher order.

'Indirect crosstalk', due to the induction of signals in a third circuit as an interfering intermediate stage, and often discussed previously [2], has also been omitted from our network calculations. This is because its contribution to the total crosstalk in high-quality telephone cables is negligible. The effect of the 'other' circuits is included in the calculation of the primary parameters through the cage effect. The contribution from a cage, which causes a marked increase of the capacitances between the circuits enclosed by the cage, is much more important than this crosstalk via third circuits. It would be more accurate, of course, to take *both* forms of crosstalk into account in the calculations; however, it is sufficient to include the cage effect alone if a first-order approach to crosstalk is considered acceptable.

A signal-carrying circuit therefore only experiences significant interference as a result of direct crosstalk from the other main signals on the other circuits in the cable.

We have calculated both the Near-End Crosstalk (NEXT) and the Far-End Crosstalk (FEXT), for the near and far ends of the circuit receiving interference, by inserting the primary cable parameters in N. A. Strakhov's equations [10]. The coupling parameters of importance in crosstalk (C_{mn} and L_{mn}) depend on the length coordinate (z). This dependence was taken into account by substituting our calculated coupling parameters in the appropriate equations. Since many conversations are transmitted simultaneously in a telephone cable, the crosstalk is often the sum of many contributions (with complicated statistical aspects). The annoyance caused by crosstalk in a telephone circuit depends of course on its magnitude relative to the desired speech signal. The annoyance due to the crosstalk increases with the recognition of this unwanted signal as speech.

Theoretical background

The theory used here for calculating the primary and then the secondary cable parameters is J. R. Carson's well-known quasi-stationary theory [11], which has also been widely used elsewhere. The twisting of the conductors introduces fundamental difficulties with this theory, but we were able to find a way round the difficulties.

Carson's theory shows that for perfectly parallel conductors signal transmission is possible in the form of waves of the quasi-TEM type, virtually transverse electromagnetic waves. A subsidiary condition is that the signal wavelength should be large compared with the distance between the conductors in a cross-section. Table II lists a number of field-strengths and current densities, etc., for cables carrying waves of the quasi-TEM type. These waves closely resemble pure TEM waves, which can only be excited if the material between the conductors is completely homogeneous and the wires are perfect conductors.

Transmission with quasi-TEM type waves is efficient. The mean lateral radiation of energy is virtually negligible and the attenuation in the direction of propagation is low for wires that are good conductors.

In applying the quasi-stationary theory there are two separate steps in the calculation of both the transmission and the crosstalk behaviour. The first step gives the primary cable parameters (the first level in fig. 3), by calculating them from the transverse components of the field vectors E and H . To do this, Maxwell's equations of the electromagnetic field, with their boundary conditions, are solved for a cross-section of the cable. The second step gives the secondary cable parameters (the second level in fig. 3), by deter-

mining the wave propagation along the cable. As mentioned earlier, the determination takes the form of an analysis of voltages and currents in the network model, the complete diagram in fig. 4.

In actual telephone cables there is the complication, already discussed, that the conductors are *not* parallel, which means that the quasi-stationary theory does not apply. The distances between the twisted wires vary periodically with the coordinate (z) along the length of the cable.

To find the primary cable parameters at a cross-section we have kept to the *local* values of the distances between the wires — defined between their centres. In doing so we have assumed that the wires are locally parallel to the axis of the cable and that the cross-section may be treated as an infinitely short length of an infinitely long cable with parallel wires (this is our stratagem). The quasi-stationary theory is then applicable and it can be used to solve the field equations with their boundary conditions for the cross-section. The twisting of the wires can then be taken into account by slightly altering the position of the cross-section along the cable and solving the field equations with their boundary conditions again.

We have used this method of approximation to calculate the primary cable parameters; they turn out to be periodic functions of the coordinate of length. This is in fact due to the periodic variation of the transverse distances between the wires.

In the calculations of the crosstalk problem we must take accurate account of the periodic variation of the primary cable parameters. For calculating the

Table II. Field-strengths and current densities for electromagnetic waves of the quasi-TEM type (virtually transverse), in a cable with parallel high-conductivity wires. The longitudinal direction of the cable coincides with the direction of the z -axis. Inside and outside the conductors the z -dependence of these quantities is given by $\exp(-\gamma z)$, where γ is the propagation coefficient.

Quantity	Region in the cable	
	Inside conductors	Outside conductors
E Electric field-strength	$E_z^{[a]} \gg E_x, E_y$	$E_z \ll E_x, E_y$
H Magnetic field-strength	$H_z \ll H_x, H_y$	$H_z \ll H_x, H_y$
J Current density	$J \gg 0$	≈ 0
$\partial_r D$ Displacement current density	$ \partial_r D \ll J$	$\neq 0^{[b]}$

[a] E_z inside the conductors is orders of magnitude smaller than E_x and E_y outside the conductors.

[b] The value depends on the frequency and remains small up to the microwave range.

[9] S. A. Schelkunoff, Conversion of Maxwell's equations into generalized telegraphist's equations, Bell Syst. tech. J. 34, 995-1043, 1955.

[10] N. A. Strakhov, Crosstalk on multipair cable — theoretical aspects, in: NTC 73, Conf. Rec. Nat. Telecomm. Conf., Atlanta 1973, Vol. I, pp. 8B/1-7.

[11] J. R. Carson, The rigorous and approximate theories of electrical transmission along wires, Bell Syst. tech. J. 7, 11-25, 1928.

transmission behaviour, on the other hand, this variation is not so important. Consequently in the transmission calculations the cable parameters can first be averaged over z , so that we can after all substitute values for the primary cable parameters that are independent of z in the network equations (1, 2, 3).

A diagram showing the method of solving the field equations, with their boundary conditions, for a cable cross-section is given in fig. 5. It can be seen that the calculation of the primary cable parameters R and L is kept separate from the calculation of C and G .

The first calculation aims primarily at determining the N magnetic vector potential fields ($\{A_z\}_m$, $m = 1, 2, \dots, N$ in fig. 5a),

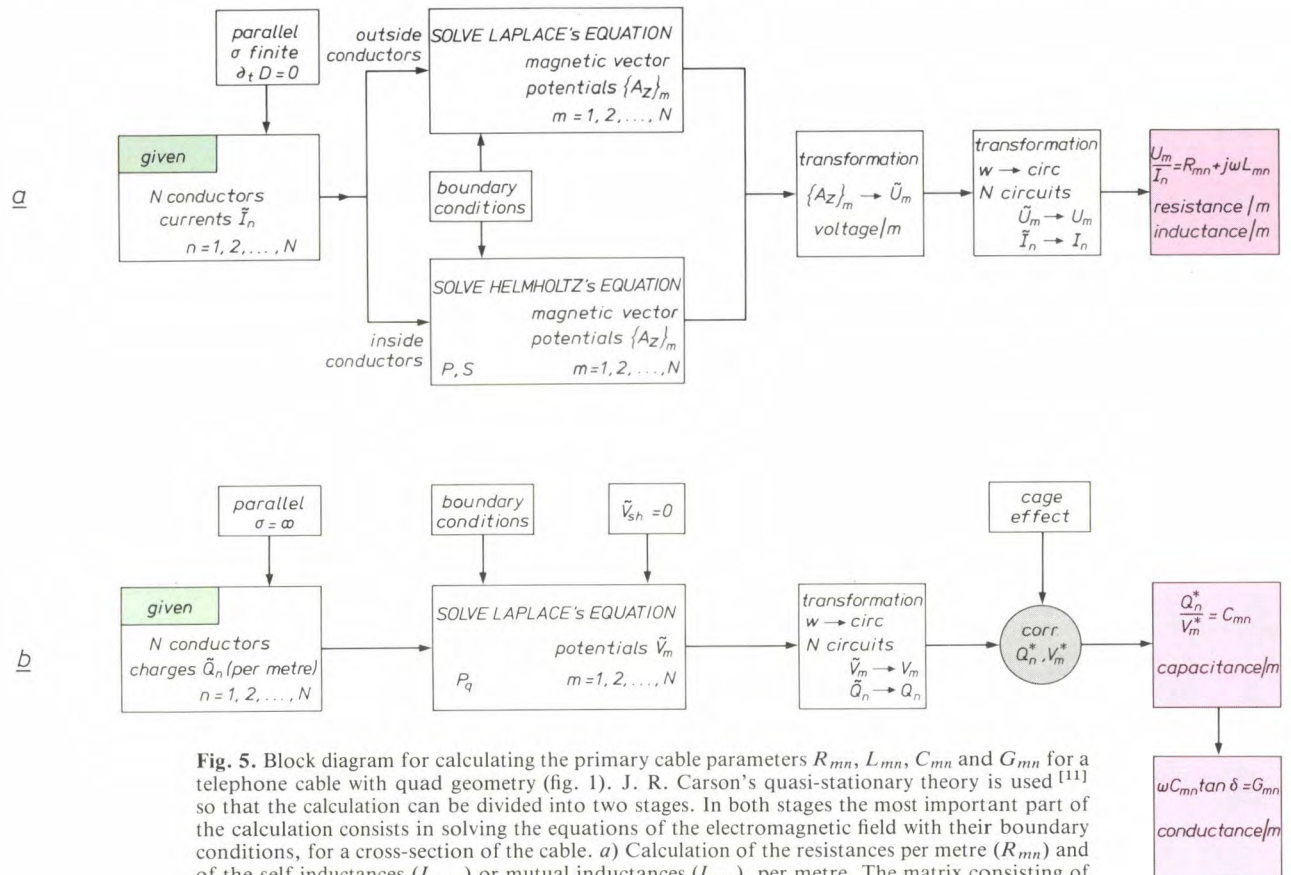


Fig. 5. Block diagram for calculating the primary cable parameters R_{mn} , L_{mn} , C_{mn} and G_{mn} for a telephone cable with quad geometry (fig. 1). J. R. Carson's quasi-stationary theory is used [11] so that the calculation can be divided into two stages. In both stages the most important part of the calculation consists in solving the equations of the electromagnetic field with their boundary conditions, for a cross-section of the cable. *a*) Calculation of the resistances per metre (R_{mn}) and of the self-inductances (L_{mm}) or mutual inductances (L_{mn}), per metre. The matrix consisting of the elements $R_{mn} + j\omega L_{mn}$ is the impedance matrix (fig. 6a). The essential feature of the calculation is the solution of Laplace's equation outside the conductors and Helmholtz's equation inside the conductors. The solution to Helmholtz's equation also provides as a result the proximity effect P and the skin effect S (Table I). σ conductivity. $\partial_t D$ displacement current density. The symbol \sim indicates that the quantity has been defined for a conductor (without the symbol it has been defined for a circuit). $w \rightarrow$ circ transformation of quantities for a conductor to corresponding quantities for a circuit. *b*) Calculation of the capacitances per metre (C_{mn}) and the conductances per metre (G_{mn}). The essential feature of the calculation is the solution of Laplace's equation outside the conductors. One of the results is the electrostatic proximity effect P_q . *corr* correction for the cage effect (K in Table I). δ loss angle of the material outside the conductors. The cable sheath acts as a shield at zero potential. (Note: R_{mm} , L_{mm} , C_{mm} , G_{mm} are called R_m , L_m , C_m , G_m in Table I.)

The resultant loss of accuracy in the transmission calculations is insignificant.

A limiting condition in the calculations is that all the pitch lengths in the cable have to be small with respect to the signal wavelength. So far this condition has always been satisfied, because in practice attenuation and crosstalk increase so rapidly as wavelength decreases that a cable becomes useless at frequencies well below those at which our approximate quasi-stationary theory loses its validity.

where N quasi-direct currents in the wire cores (\tilde{I}_n) form the input data. The voltage drop per metre (U_m) along the conductors can be calculated from the magnetic fields. Dividing by the current gives the required complex impedance, which can then be used to find the resistance per metre, the mutual inductance per metre and the self-inductance per metre.

The main purpose of the second calculation is to determine N electrostatic potential fields (\tilde{V}_m , $m = 1, 2, \dots, N$ in fig. 5b), with the charge per metre on each wire (\tilde{Q}_n) as the input data. The cage effect, which only occurs in wires that are twisted, is included in the calculation as a special boundary condition, which states that circuits forming part of a cage must everywhere have the constant

potential of the cable shield (zero). As soon as the complete potential distribution is known, all the capacitances (per metre) can now readily be calculated.

The separation of the calculations can be seen as a further simplification of the quasi-stationary theory. In this approach the transverse field associated with the TEM waves is approximated by taking the individual contributions from a static electric field and a quasi-static magnetic field. The two fields are thus effectively 'decoupled' by omitting from Maxwell's equations contributions from the magnetic field and from the electric field that are very small. The first field, the electric field, is therefore calculated for the simpler case of perfect conductors ($\sigma = \infty$); this is a good approximation, since the z-component of the electric field inside the conductors is extremely small (see Table II). The second field, the magnetic field outside the conductors, is calculated with the displacement currents set equal to zero ($\partial_t D = 0$, in fig. 5a); the field is then due entirely to the currents in the conductors. This condition is correct for all signal frequencies below those in the microwave band. It should be noted in passing that the current-density distributions of the quasi-direct currents in the calculations in fig. 5a are not uniform over the cross-sections of the wires: the current concentration resulting from the proximity effect and the skin effect are both found when solving the equations for the magnetic vector potential (Helmholtz's equation) inside the conductors.

The program

The strategy

In designing our cable-simulation program we followed the distribution of the calculations over three levels as illustrated in fig. 3. The programming method adopted is known as 'stepwise refinement' ^[12], which means that a further redistribution of the separate parts of the calculation is introduced at each level by splitting it up into increasingly refined modules. One of the main objects of such a strategy is to make the program as 'user-friendly' as possible; the extent to which this has been achieved was checked against five quality criteria for readily communicable and user-friendly software (Table III). Since, as we have seen, there are two distinct objectives in research on telephone cables, it was a challenge for us to ensure that the program would be user-friendly in regard to either objective. The manufacturer can incorporate any likely modifications to a cable design into the program without too much difficulty, while the engineer who has to consider the possibilities for data signals can easily calculate the transmission behaviour of existing cables.

The stepwise-refinement procedure offers some appreciable advantages. One is that a module can be modified, considerably if necessary, without introducing the need for significant changes in the rest of the program. We found many advantages in using modules, especially while the program was being written; the scheme adopted made the program easy to

grasp and therefore easy to communicate. Another important aspect of the modular structure is that it is easy to check that the program is correct. The check has two aspects: making sure that the relations between modules are as intended, and ensuring that the individual modules are operating to specification. We found that such checks took little time or trouble.

To meet the criteria for 'efficiency' and 'comprehensibility' (Table III) we obviously had to give careful consideration to the choice of programming language. The entire program was written in PASCAL, except for a few modules in FORTRAN (see next section). The program was run on a VAX11/780 mini-computer, which operates with words of 32 bits. The version of PASCAL used had already been provided by the computer manufacturer with facilities for compiling each of the modules separately into machine instructions. This saved a great deal of time when we came to write our program (and it still saves time whenever we have to make modifications).

Choice of language

When we were choosing the programming language there were two main conditions. The first was to have a language that could easily handle data *types*, so that we could specify the geometrical structure of a telephone cable as freely as possible. The second condition was that we should be able to solve a wide variety of numerical problems simply and rapidly, since the calculations for determining the primary cable parameters are extensive. There were many languages suitable for 'structured' programming, and because of the conditions just described we finally decided to use PASCAL. One important consideration was that PASCAL has much to offer for type classification, thus helping to ensure the comprehensibility of the

Table III. Quality criteria used in designing the cable-simulation program, mainly to ensure smooth communication and 'user friendliness'.

Criterion	Description
Correctness	Each module should meet its specification exactly.
Completeness	If properly used the program should run correctly; each error should be clearly signalled.
Flexibility	Small program modifications should be possible 'on the spot', and should not therefore require extensive rewriting of other parts of the program.
Clarity	The program should be readable and understandable.
Efficiency	For all task implementations fast standard algorithms are preferable to 'clever' (but unclear) <i>ad hoc</i> solutions; for interactive operations the time needed per task unit should be of the order of seconds.

^[12] N. Wirth, Program development by stepwise refinement, Comm. ACM 14, 221-227, 1971.

program. PASCAL does have some disadvantages, such as the absence of variables with the type of a complex number, so that calculations with complex numbers required extra programming. Another disadvantage was that the limits of declared arrays cannot be dynamically changed during execution, and this was also compensated for by additional programming work.

At the places in the program where extensive calculations are required, e.g. to solve a large set of linear equations, a separate processor, the 'array processor', was used for greater efficiency. This processor, which is exceptionally fast, only accepts modules written in FORTRAN. At these places the program calls up routines in FORTRAN, even though it is written in PASCAL itself. This is possible largely because of the modular structure of the program.

Example of a module

Fig. 6 shows a module, called ZMATRIX, whose complete program text is printed in fig. 6a. The main operation carried out in this module is the determination of the impedance matrix at a cross-section of a cable. In other words, this module establishes, for a particular cross-section of a telephone cable, the equations — in the form of a matrix — relating the voltages on the individual wires to the currents in these wires (fig. 5a). The calculations are complicated and include a number of steps, each occupying an independent submodule. The flow chart in fig. 6b shows the sequence of these sub-operations. We shall now see how a number of the quality requirements listed in Table III reveal themselves in this typical program module.

Correctness. The program module (fig. 6a) starts at the top with the specifications of the module. In this particular case merely specifying the main operation, as is done here, is sufficient. The specification should always be kept as simple as possible, because the cor-

```

(*-----*)
(* In a cable cross section the impedance matrix ZMAT is calculated *)
(* ZMAT relates the (per unit length) longitudinal voltage drop and *)
(* the current in the wires. *)
(*-----*)

%INCLUDE 'global.dat'

PROCEDURE filla (VAR mat :typ_A;
                VAR w :typ_Warr;
                VAR L :typ_Larr); EXTERN;
PROCEDURE fillb (VAR b :typ_rhs;
                VAR w :typ_Warr); EXTERN;
PROCEDURE lamn (VAR L :typ_Larr); EXTERN;
PROCEDURE screen (VAR w :typ_Warr); EXTERN;
PROCEDURE fillzmat (VAR zmat :typ_z;
                  VAR w :typ_Warr;
                  VAR rhs :typ_rhs); EXTERN;
PROCEDURE APsolv (VAR bar :typ_A;
                 VAR rhs :typ_rhs;
                 n,m,st :integer); FORTRAN;

PROCEDURE zmatrix(VAR ZMAT: typ_z);

CONST UK = 1;

VAR A : typ_A;
    b : typ_rhs;
    Lns : typ_Larr;
    Wp : typ_Warr;
    st,N,M: integer;

BEGIN
    N:=dim; (* global constants dim and stmax *)
    M:=stmax;

    lamn(Lns); (* fill the lambda-matrix Lns *)
    screen(Wp); (* fill screen factor array Wp *)

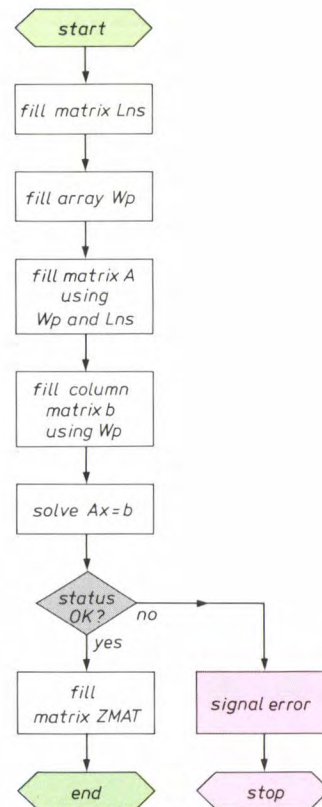
    filla(A,Wp,Lns); (* fill matrix A (NxN) using Wp,Lns *)
    fillb(b,Wp); (* fill matrix b (NxM) using Wp *)

    APsolv(A,b,N,M,st); (* solve Ax=b in Array Processor *)
    (* x is returned in b, status in st *)

    IF st=OK THEN
        fillzmat(ZMAT,Wp,b);
    ELSE
        BEGIN WRITELN('ERROR in APsolv!');
              WRITELN('Status = ',st:2);
              HALT
        END;
    END;
END;

```

a



b

Fig. 6. Example of a module in the program for simulating telephone cable (fig. 3). This module can be used for determining the impedance matrix (fig. 5a) for a cross-section of the cable. a) The program of the module. The program starts with a list of the specifications that the module has to meet; then the procedures and variables are declared; the third segment, between BEGIN and END, contains the calculations, for which (b) represents the flow chart. In addition to the specifications the program later gives explanations in natural language as well, i.e. all lines between (* and *). The program itself is written in the high-level language PASCAL. The instruction explained by (* solve $Ax = b$ in Array Processor *) brings in a second processor. This is faster than the ordinary type, which improves the efficiency. It requires FORTRAN software, however. The subroutine used solves the set of linear equations $Ax = b$ with the aid of a standard algorithm in FORTRAN; the text is not given here. b) Flow chart for the calculation in (a). The different operations are executed sequentially where possible, for simplicity and clarity. The only branching in this example is included to provide a check (red blocks) on fundamental difficulties or computing errors in solving the equations $Ax = b$.

rectness of the module can then be demonstrated by means of simple tests. Any limiting conditions appropriate to the use of a module, which therefore form part of the specification, are included at the top. The flow chart in fig. 6*b* illustrates another objective, which is to compose the modules from a very limited number of sub-operations, whose interrelationships can be represented by a diagram with as few branches as possible.

the program is of course improved by the choice of a language such as PASCAL; with such a high-level language it is easier to relate the formulation of the problems to natural language than it would be with a computer-oriented language (e.g. an assembly language). This again increases the comprehensibility. *Efficiency.* An example of the introduction of a fast, standardized algorithm to improve the efficiency of the program is the use of the array processor mentioned

Table IV. Comparison of calculated and measured values of attenuation and phase shift in a telephone cable with strong proximity effect ^[13].

Signal frequency (kHz)	Quantity	Attenuation (dB/km)			Phase shift (rad/km)		
	Measured (Lenahan ^[13])	Calculated (Lenahan)	Calculated (ours)	Measured (Lenahan ^[13])	Calculated (Lenahan)	Calculated (ours)	
5	2.55	2.56	2.56	0.42	0.42	0.42	
10	3.30	3.28	3.28	0.69	0.70	0.70	
50	7.40	7.46	7.46	2.36	2.35	2.35	
100	11.23	11.22	11.21	4.71	4.69	4.68	
500	27.96	27.53	27.52	19.45	19.51	19.48	
1000	41.95	41.83	41.80	36.85	36.98	36.94	
5000	111.73	112.72	112.87	168.70	168.74	168.60	
10000	166.25	168.02	167.93	328.85	328.83	328.61	

Completeness. The blocks 'status OK' and 'error signal' in fig. 6*b* are necessary to ensure that the text of the program is complete in this respect. In solving the set of linear equations $Ax = b$ the first step is to determine whether all the basic conditions (e.g. $\det A \neq 0$) have been satisfied and then to check whether the calculations are proceeding correctly. Any errors are signalled and this stops the program. If there are no errors, the matrix ZMAT can then be set up or 'filled', completing the operation of the module.

Flexibility. The modular structure is the best guarantee of flexibility. In the case of fig. 6 the set $Ax = b$ is solved by the numerical routine APsolv, a separate sub-module. If another method of solution seemed better, and assuming that it could also be embodied in a sub-module, called SOLV here, then we need only replace sub-module APsolv by its counterpart SOLV, and compile the text of the module ZMATRIX again. (At the level of the ZMATRIX module the available PASCAL version offers separate compilation.)

Comprehensibility. In addition to the specification the program text in fig. 6*a* also indicates the sub-operations that have to be carried out in succession. This is in fact a kind of description of fig. 6*b* in words, which improves the comprehensibility. In addition to the simplicity of the flow charts, the comprehensibility of

above for solving the large set of linear equations $Ax = b$ in fig. 6*b*. At places such as this the program calls up routines in FORTRAN, since the array processor can only work in FORTRAN. The incorporation of program text in another language is facilitated by the modular structure. In this way some parts of the cable-simulation program with a great deal of calculation can be handled rapidly and by standardized methods, which greatly improves the efficiency and also improves the comprehensibility.

Results and their verification

We used the quasi-stationary theory described earlier and the simulation program to calculate the primary and secondary cable parameters of a number of experimental cables made by NKF. We verified the results by comparing them with measurements on the cables. The calculations and the measurements were found to agree well, both for cables with parallel wires and for cables with twisted wires. (The differences amounted to no more than 5% of the calculated value.) Some of the differences are attributable to small variations in the manufacturing processes, which are difficult to avoid. The conclusion is that the cable-simulation program gives results with an error of no more than a few per cent.

As a further check on our program we made a comparative calculation of the attenuation and phase shift in a cable made by Bell Laboratories on which T. A. Lenahan had carried out calculations and measurements^[13]. Some of the results are listed in *Table IV*. In this cable, which contains only two wires inside a closely fitting shield, the proximity effect is strongly pronounced. Here again the agreement was entirely satisfactory. In the meantime a number of calculations have been made on cables with much more complicated wire geometries, and these have also given satisfactory results. The simulation of a 'real' telephone network using computer software in the man-

ner illustrated in fig. 3 is therefore something that stands a good chance of becoming reality in the not too distant future.

Summary. A study of multiwire telephone cables with quad geometry, with the twin objectives of improving the quality of the conventional copper cable with twisted wires, and of introducing digital communications via the national telephone system, is most easily carried out by computer simulation. The modular software for cable simulation, with 'stepwise refinement' and with emphasis on transmission behaviour and crosstalk (which are particularly important in local networks) is written in PASCAL for a VAX11/780 minicomputer. Primary (R, L, C, G) and secondary cable parameters and signal-transfer matrices of complete cables can be calculated with this software for given dimensions and materials. Two proximity effects (due to currents and charges), the skin effect and Martin's cage effect in the generalized form given by Belevitch are included in the calculations. A program module described as an example demonstrates the 'user-friendly' nature of the software, as well as other qualities such as efficiency and error-signalling facilities. The method of calculation is based on Carson's quasi-stationary theory of electrical transmission along wires, adapted to twisted conductors. Results of cable calculations are given along with their verification by measurements.

^[13] T. A. Lenahan, The theory of uniform cables — Part I: Calculation of propagation parameters, Part II: Calculation of charge components, and also: Experimental test of propagation-parameter calculations for shielded balanced pair cables, Bell Syst. tech. J. 56, 597-610, 611-625 and 627-636, 1977.

The TRAPATT oscillator

R. Davies, B. H. Newton and J. G. Summers

It is not uncommon in the history of a new semiconductor device for many years to elapse between the first reports of exciting research results and the eventual emergence of a mature component. The TRAPATT diode is no exception, and the promise of ten years ago is only now reaching fulfilment. This diode can provide high peak powers at microwave frequencies of several GHz, but careful circuit design is necessary to capitalize on the full capabilities of the device. Research scientists at Philips Research Laboratories, Redhill, and development engineers at Mullard Ltd have cooperated in a systematic experimental investigation of TRAPATT diodes and the oscillator circuits in which they function. As a result, the TRAPATT oscillator can now be considered as a serious contender for pulse-transmitter applications in the frequency range 2-4 GHz.

Introduction

Since the early sixties, many solid-state diodes have been studied as potentially useful sources of microwave energy. Sometimes these devices have been used as low-power oscillators requiring separate external amplification to achieve useful power levels. However, many studies have focused on their characteristics as fundamental microwave oscillators producing sufficient power directly. Both continuous and pulsed operation have been investigated.

The IMPATT (IMPact Avalanche and Transit Time) diode^[1] was the subject of considerable research over a period of about ten years, starting in 1965. This is a p-n junction diode that has a negative resistance when reverse biased; its operation depends on the transit time of carriers generated in an avalanche region passing through the device. During studies of this diode, a mode of operation was observed that differed markedly from the normal IMPATT mode. This new behaviour was characterized by its very high electrical efficiency, high power output and lower frequency of oscillation, and was referred to as the 'anomalous mode'^[2].

Out of this work was developed the TRAPATT-diode oscillator, where the acronym 'TRAPATT' stands for TRApped Plasma Avalanche Triggered Transit. This is the name by which the anomalous

mode is now known, and describes the operation of the diode when incorporated in a suitable circuit. Several computer simulations have successfully explained the basic TRAPATT mechanism^[3], but no detailed consideration has been given to the interaction between the diode and the circuit. Consequently, an empirical approach to device and oscillator design has been followed that, until recently, has led to sub-optimal designs and over-complex circuits, giving rise to unreliable operation.

Initial experimental results using this new type of diode were very encouraging, and high peak powers were obtained at frequencies up to about 10 GHz, with d.c. to r.f. conversion efficiencies greater than 30%. The present-day performance of the TRAPATT diode and three contemporary devices is illustrated in *fig. 1*, for single diodes in pulsed oscillators. Although not offering the very highest output powers, the TRAPATT diode is more efficient than its competitors, and can operate at higher duty cycles. Its operating frequency, however, is limited to the lower end of the microwave spectrum. Despite the early promise of the device, and its attractions as a source for use in

R. Davies, Ph.D., B. H. Newton, Ph.D., and J. G. Summers, B.Sc., are with Philips Research Laboratories (PRL), Redhill, Surrey, England.

^[1] D. de Nobel and M. T. Vlaardingerbroek, IMPATT diodes, Philips tech. Rev. 32, 328-344, 1971.

^[2] P. J. de Waard, Anomalous oscillations with an IMPATT diode, Philips tech. Rev. 32, 361-369, 1971.

^[3] See for example B. C. DeLoach, Jr., and D. L. Scharfetter, Device physics of TRAPATT oscillators, IEEE Trans. ED-17, 9-21, 1970.

pulsed radars, there has been no significant impact beyond the results of the initial stages of development. Very few, if any, of the early TRAPATT oscillators have been manufactured in any quantity, and there are few remaining TRAPATT studies.

This failure to make an impact on microwave systems was due to difficulties in establishing and maintaining coherent oscillations. (By 'coherent' is meant a fixed-frequency oscillation in which the fundamental and any harmonics have amplitude and phase relationships that are constant.) The lack of progress can be largely attributed to serious shortcomings in the widely accepted design for the oscillator circuit. Circuits based on this design required empirical alignment, performance was unreliable, and coherent operation into practical r.f. loads was difficult to maintain. This article describes a new design that overcomes these three disadvantages, providing a firm basis from which TRAPATT oscillators can be further developed, to the point where routine production is possible.

In parallel with the progress in circuit design, diode designs have also advanced, to a stage where a range of silicon-planar devices can be mass-produced at low cost. This progress has resulted from a coordinated research programme that has highlighted critical features in a number of associated topics. We now have a fuller understanding of the underlying physical processes involved in the TRAPATT mechanism, leading to a better appreciation of optimum 'geometry', ther-

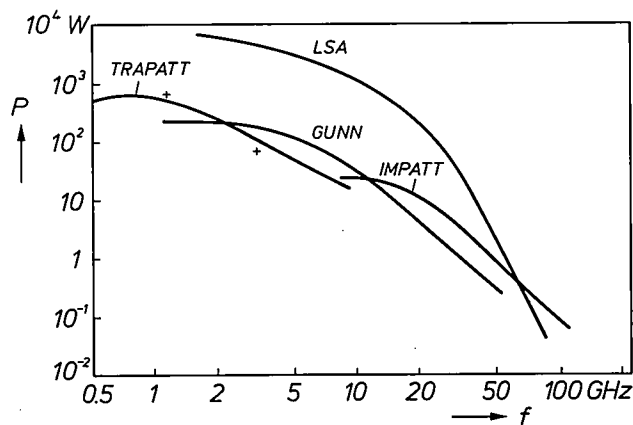


Fig. 1. Experimental results are summarized in this diagram as peak output power P plotted against frequency f for four kinds of microwave diode. These are the TRAPATT (Trapped Plasma Avalanche Triggered Transit) diode, the IMPATT (Impact Avalanche and Transit Time) diode, the LSA (Limited Space Charge Accumulation) diode, and the Gunn diode. Although the TRAPATT diode may not appear to offer the highest peak powers in the band of interest (S-band, 2-4 GHz), it is more efficient than the Gunn diode, and can operate at higher duty cycles than the LSA device. These results are for single devices. Comparable powers (+) can be obtained using bipolar transistors to amplify the output of a low-power microwave source, provided a combination of many transistors is used.

mal design and choice of silicon material. Our studies have led to a new diode structure that is more reproducible, more efficient, has lower thermal resistance, and together with the new circuit results in an oscillator that is considerably less temperature-sensitive. Present-day devices perform reliably at duty cycles of up to 3% with heat-sink temperatures in excess of 100 °C, and can be produced in a range of chip sizes.

The larger-area diodes and the new circuit have impedances that closely match each other, and considerably higher r.f. powers can be obtained. Peak power levels of 300 W at 0.1% duty cycle and 250 W at 2% duty cycle can already be achieved at S-band (2-4 GHz) with single diodes, and increased power from interconnected diodes has been demonstrated.

TRAPATT oscillator circuits

Basic principles

The TRAPATT oscillator is based on a reverse-biased silicon p-n junction that is operated at avalanche breakdown. This breakdown occurs when the electric field-strength within the depletion region of the junction is sufficiently high that ionization processes give rise to a cumulative multiplication of charge carriers (holes and electrons). Fig. 2 shows the basic TRAPATT diode, in which the n^+ substrate is bonded directly to the package, and the second connection is made via a gold wire, thermocompression-bonded to the upper contact. The active junction is at the interface of the p^+ and n regions.

The principles of circuit operation and the origin of the device acronym can be understood most easily by considering the response of a p-n junction, biased just below avalanche breakdown, to an applied voltage impulse. The impulse causes the electric field-strength within the depletion region of the junction to rise

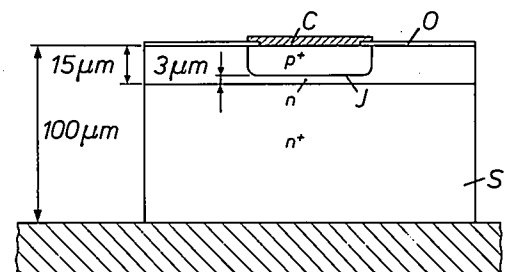


Fig. 2. The basic TRAPATT diode is manufactured using suitably doped silicon and conventional processing techniques. It consists of a p^+n junction J lying on an n^+ substrate S . The package is connected directly to the substrate, and an ohmic contact C is provided for a wire connection, through the surface oxide O . Typical vertical dimensions are shown in the figure (not drawn to scale). The diode area is chosen to suit the particular application. This article deals mainly with diodes of area 4×10^{-4} cm², and higher-power devices with areas of up to 12×10^{-4} cm².

sufficiently for avalanche multiplication to occur. As the voltage applied increases, the avalanching region spreads through the junction depletion region. When this avalanching region expands faster than the maximum velocity at which the carriers can be extracted (the saturated drift velocity), an extremely large number of carriers is generated in the avalanche region, which is filled with a dense, electrically neutral concentration of these charge carriers. This is known as the 'trapped plasma' state.

When the associated high space-charge density becomes sufficiently large, the electric field within the device collapses to almost zero. Consequently the carriers are extracted very slowly, at a rate much lower than that corresponding to the standard drift velocity. The electric field within the device gradually recovers as the carriers are extracted. To provide a suitable recovery characteristic and to minimize resistive losses,

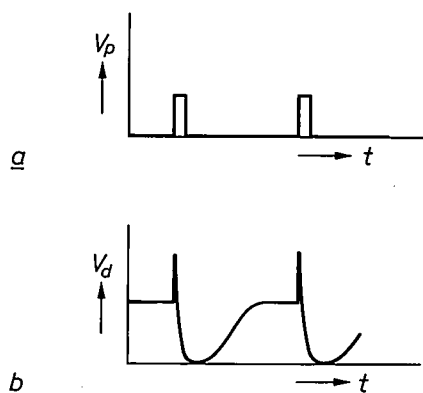


Fig. 3. Diagrams to illustrate voltage variation with time t during TRAPATT oscillation. *a*) The assumed train of positive voltage pulses V_p incident on the diode from the oscillator circuit. *b*) The device voltage V_d , showing the transition to avalanche conditions in response to the voltage pulse, and the subsequent recovery. The oscillator circuit is so designed that the negative-going pulses from the diode are inverted and reflected back towards the diode, thus forming the train of positive pulses assumed here.

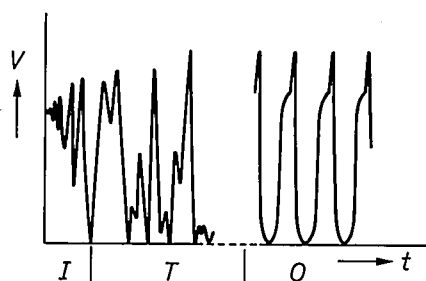


Fig. 4. Typical variation of the diode voltage V against time t , during the establishment of coherent TRAPATT oscillations. Three phases can be distinguished: the initiation phase *I* in which an IMPATT oscillation possibly occurs, leading to a transient phase *T* of non-coherent trapped-plasma oscillations, and the resulting TRAPATT oscillation *O* itself. Circuit designs should take account of this three-stage behaviour.

TRAPATT diodes have an epitaxial-layer thickness and doping such that the depletion region extends throughout the low-doped layer before avalanche breakdown is reached [4].

Fig. 3 illustrates the response of the reverse-biased diode to an assumed train of positive voltage pulses, incident on the diode from the oscillator circuit. The circuit is designed to provide the necessary sequence of voltage-triggering pulses, once the initial pulse has been generated. This initial pulse originates within the device itself, possibly as the result of IMPATT oscillations. This has not been verified conclusively, although several suggestions have been made [5]. Computer models have shown the TRAPATT oscillation developing from a high-frequency IMPATT mode as indicated in fig. 4. To introduce the circuit concepts, let us assume the incidence of a trigger pulse, and see how the circuit provides the conditions for self-sustained oscillations.

The basic circuit is shown schematically in fig. 5; it is known as a Time Delay Triggered (or TDT) circuit. It consists of the TRAPATT diode coupled to a delay line and a lowpass filter. We have already seen how the application of a voltage impulse causes the voltage across the diode to collapse. This collapsing voltage couples to the circuit delay line and a negative voltage pulse is launched along the line. This pulse travels along the line until it reaches the filter, which is designed to provide a resistive termination at the microwave output frequency and a low impedance at the harmonics of this frequency. (We shall see later that energy may be extracted at harmonics of the fundamental frequency, but for the present we consider only fundamental operation.) On reaching the filter, the negative pulse is inverted, partly reflected, and partly transmitted to the load. The matching is ar-

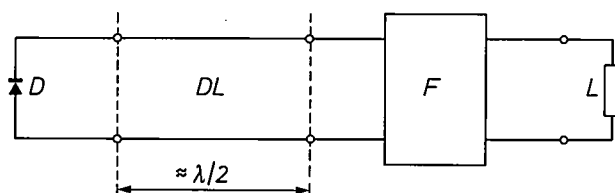


Fig. 5. The three fundamental components of the Time Delay Triggered, or TDT, circuit are the TRAPATT diode D , a delay line DL , and a lowpass filter F . The length of the delay line determines the frequency of oscillation. Since each trigger pulse travels from the diode to the filter and back, the line should be of a length equal to half the wavelength at the desired frequency of oscillation. (In practice it is slightly less than this, to allow for delays in the diode itself.) A suitable external load L is connected to the output of the filter.

[4] A detailed description can be found in G. Gibbons, *Avalanche-diode microwave oscillators*, Clarendon Press, Oxford 1973.

[5] See for example B. B. van Iperen, Efficiency limitation by transverse instability in Si IMPATT diodes, *Proc. IEEE* 62, 284-285, 1974.

ranged to maximize the power dissipated in the load at the fundamental frequency, while ensuring that the returning positive pulse is large enough to trigger the diode directly into the TRApped Plasma state, and cause the Avalanche region to undergo a Triggered Transit through the device. The subsequent collapse of the electric field within the diode generates a further negative pulse in the delay line and the process will continue until the d.c. bias is removed from the diode. This simple model shows that the TRAPATT oscillator is primarily circuit controlled and the frequency of oscillation depends mainly on the electrical length of the delay line. This length is approximately one half of a wavelength at the fundamental TRAPATT frequency.

The TRAPATT oscillations, once started, will continue in the presence of the d.c. bias. The power available is limited by the temperature rise in the diode structure, even with an efficient heat sink. To limit the diode temperature to a maximum of 200 °C at large bias currents, the oscillator is operated in a pulsed mode. In this, the bias is applied as a series of pulses; large peak powers can then be obtained. Pulse widths of about 100 ns and duty cycles of a few per cent are commonly used, corresponding to a pulse-repetition frequency of a few hundred kHz. Our investigations have been confined to this mode of operation.

Circuit design

Apart from the diode, the only major components of the TDT circuit in fig. 5 are the delay line and the filter. It might appear that all that is necessary to design a circuit is to choose the length and impedance of the line, and the match and frequency passband of the filter. A microstrip circuit designed according to this simple concept would be as shown schematically in fig. 6. The length of the delay line would be determined from the desired period of oscillation, and the filter characteristic would be optimized experimentally. In practice, TRAPATT oscillations were initially observed in such a circuit, and the theoretical explana-

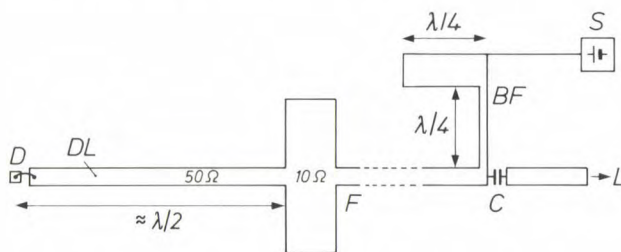


Fig. 6. This shows the basic microstrip oscillator layout, with a uniform delay line *DL* and a lowpass filter *F* of unspecified geometry. The diode *D* receives its d.c. bias via the delay line, which is coupled to a d.c. source *S* via a bias filter *BF*. A blocking capacitor *C* is included at the output to prevent the bias from reaching the load *L*.

tion was developed only later. Until recently almost all the reported practical results were measured in such circuits. Typical coaxial and microstrip versions are shown in fig. 7. This circuit 'design' is based on the following three criteria:

The circuit must

- support the small-signal oscillations that grow to trigger the first cycle,
- provide enough locally stored charge to drive the plasma generation process,
- reflect the subsequent trigger pulses that maintain the oscillation.

We find, using diodes with the correct doping profile, that the first criterion is satisfied in practical circuits that comply with the second and third criteria. We deduce this from the ease with which the device voltage collapses, which we interpret as indicating TRAPATT action.

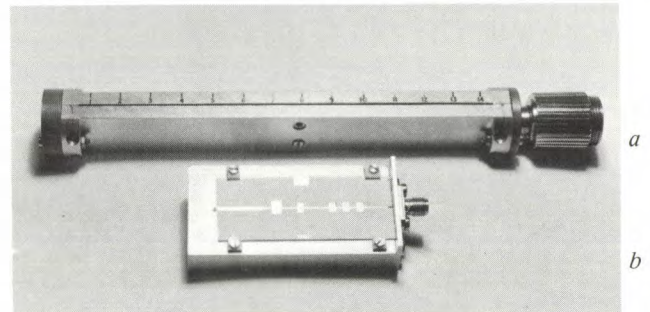


Fig. 7. Two versions of the original TDT oscillator. *a*) The coaxial version, and *b*) the microstrip version in which the delay line and filter geometries can be clearly seen. The delay line has a 30Ω section adjacent to the diode, which was preferred to direct connection to a 50Ω line. The resulting 30Ω to 50Ω step just to the left of the first large low-impedance section was found to be effective in preventing the diode from breaking down during the recovery period.

The second criterion is confirmed by the improvement in oscillator performance that results from reducing the delay-line impedance near the diode.

Satisfaction of the third criterion has been widely interpreted by TRAPATT circuit designers as requiring two circuit elements. The first is a delay line that is slightly shorter (by about 10% at S-band) than a half-wavelength at the TRAPATT frequency, to define the period of the trigger pulse. This shortening is to allow for the delay incurred by the device during the charging and avalanche-multiplication period. The second element is a filter terminating the delay line. This provides a resistive termination at the output frequency (which is usually the fundamental TRAPATT frequency but might be at a higher harmonic frequency), and a low impedance at other frequencies harmonically related to the fundamental.

There are two major shortcomings of this circuit. One is that when using a broadband-match load,

adjustment for coherent oscillation is completely empirical. The other is that subsequent operation into practical loads (e.g. isolators, bandpass filters etc.) is usually incoherent, necessitating a further adjustment that often degrades the efficiency. Before the TRAPATT oscillator could become more than a 'laboratory curiosity', a circuit design that gave predictable and reproducible interaction with the diode was essential. Such a design has resulted from time-domain studies we have made of TRAPATT circuits and associated components.

If practical loads such as filters and isolators are measured in the time domain with, for instance, a reflectometer, it is found that they have broadly similar characteristics. These loads can present a reactive mismatch to the oscillator lowpass filter, giving rise to spurious trigger pulses. In such cases, the diode may receive two or more trigger pulses per cycle, one wanted and the others unwanted. There will be no adverse effect if the wanted pulse dominates, or if the pulses are nearly coincident. However, the usual result is that coherent oscillations cannot be established.

It is the low-frequency energy content of the output signal during the transient phase that gives rise to this problem, and once this is realized, the cure becomes simple. A diplexer, connected between the oscillator and load, is used to divert the low-frequency signal to a second, matched, load. A suitable component for performing this diplexing function is the travelling-wave directional filter (TWDF)^[6]. This is a four-port component consisting of two transmission lines coupled via a resonant ring; although a resonant device, it has a constant broadband input resistance. Use of this component gives complete suppression of the spurious pulses, and has allowed operation into practical loads such as isolators and filters. A useful 'spin-off' is that the d.c. bias for the TRAPATT diode can also be fed through the TWDF, in such a way that none of the r.f. output power is dissipated in the bias circuit.

A 30Ω to 50Ω step was usually built into the delay line between diode and filter, as in practice this made it easier to establish coherent oscillations.

In the light of these results the design criteria have been modified to read as follows:

The circuit must

- support the small-signal oscillations that grow to trigger the first cycle,
- generate the optimum voltage-time response in both the transient and steady-state phases,
- reflect the subsequent trigger pulses that maintain the oscillation, without introducing spurious trigger pulses.

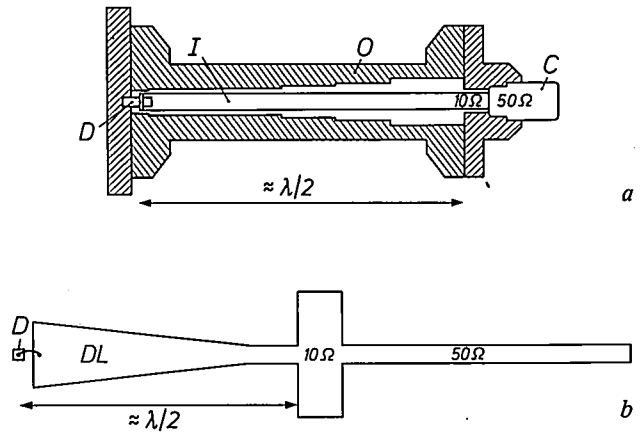


Fig. 8. Diagrams illustrating the construction of the new oscillator in its two forms: *a*) shows a cross-section of the coaxial version, in which a packaged diode *D* is connected to one end of a coaxial delay line comprising a constant-diameter inner conductor *I* enclosed in a stepped outer conductor *O*. A 10Ω section is included just before the 50Ω connector *C*. *b*) shows the microstrip version, in which a diode chip *D* is connected to one end of a linearly-tapered delay line *DL*. A 10Ω section is included before the 50Ω line, which leads to a coaxial connector (not shown). Both versions would have their output connected to the load via a travelling-wave directional filter.

The simple circuits shown schematically in *fig. 8* have resulted from the application of these new criteria. They differ from the conventional TRAPATT oscillator in two significant aspects: the delay line is stepped or tapered to prevent the diode voltage from recovering too rapidly, and there is only one element in the output filter so that spurious trigger pulses are eliminated from the circuit. This design therefore takes into account the way in which coherent TRAPATT oscillations grow during the transient phase, and produces a circuit that does not suffer from the disadvantages inherent in the earlier designs.

Oscillator construction and performance

Microstrip and coaxial versions of the new design of oscillator are shown in *fig. 9*. Whereas the microstrip version incorporates a tapered delay line, the coaxial oscillator has a delay line with five discrete impedance levels. The basic electrical characteristics of these oscillators are, however, identical, and are summarized in *Table I* for oscillators using single diodes of area $4 \times 10^{-4} \text{ cm}^2$.

The inherently lower losses in the coaxial transmission-line allow larger area devices to be used, with a consequent increase in output power, and so more attention has been paid to this version than to the microstrip one. The coaxial version is also more suitable for testing performance reproducibility since it contains packaged devices.

[6] S. B. Cohn and F. S. Coale, Directional channel-separation filters, *Proc. IRE* 44, 1018-1024, 1956.

The reproducibility of the power and frequency characteristics has been measured using a fixed-tuned version of such an oscillator. The measurements were made using diodes that were inverted so that their p+

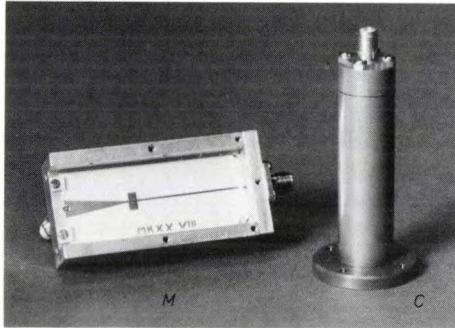


Fig. 9. Practical embodiments of the new oscillator design. The tapered delay line used in the microstrip version *M* is clearly visible. The enclosed coaxial version *C* incorporates a similar non-uniform line having five distinct steps. The basic electrical characteristics of these two designs are identical.

Table I. Basic electrical characteristics of oscillators made to the new design, using diodes of area $4 \times 10^{-4} \text{ cm}^2$. The peak power includes the losses in the TWDF, and is measured at the output of an isolator which is connected to the output of the oscillator during the measurement.

Frequency:	S-band
Peak power:	60 W minimum
Efficiency:	30% typical
Duty cycle:	2% maximum
Pulse length:	300 ns maximum
df/dT :	$-450 \text{ kHz } ^\circ\text{C}^{-1}$ maximum
Temperature range:	$-54 \text{ }^\circ\text{C}$ to $+95 \text{ }^\circ\text{C}$

contact (see fig. 2) could be bonded directly to the heat sink to minimize the thermal resistance. Fourteen such 'flip-chip' diodes of area $4 \times 10^{-4} \text{ cm}^2$, and from the same batch, were inserted sequentially. For each device, power and frequency were measured for a 4 A, 0.1% duty cycle. An acceptably small spread of $\pm 5\%$ in power was measured; the frequency spread was $\pm 23 \text{ MHz}$.

Using the coaxial-oscillator design, it has been possible to study how the oscillator characteristics relate to the circuit load at the fundamental and harmonic frequencies, and to investigate the tolerance to changes in drive current. The resistive loading at the fundamental frequency was varied by adjusting the length of the low-impedance step and the configuration of the delay line. A current drive of 4 A was used at a duty cycle of 0.1%. Fig. 10 shows how the peak output power and the efficiency increase with increasing resistive loading. It was not possible to maintain coherent operation for resistive loads in excess of 5.3Ω , because premature avalanching was brought about by the excessively large output voltage across the diode. It was also found that coherent oscillations were not critically dependent on the loading at the harmonic frequencies. The tolerance of the oscillator to variation in the drive conditions is also given in fig. 10, which shows how the power output, efficiency, and frequency vary with increasing current drive for a constant circuit configuration. The power output and efficiency increase until coherence is destroyed by premature avalanching at a peak current slightly in excess of 5 A.

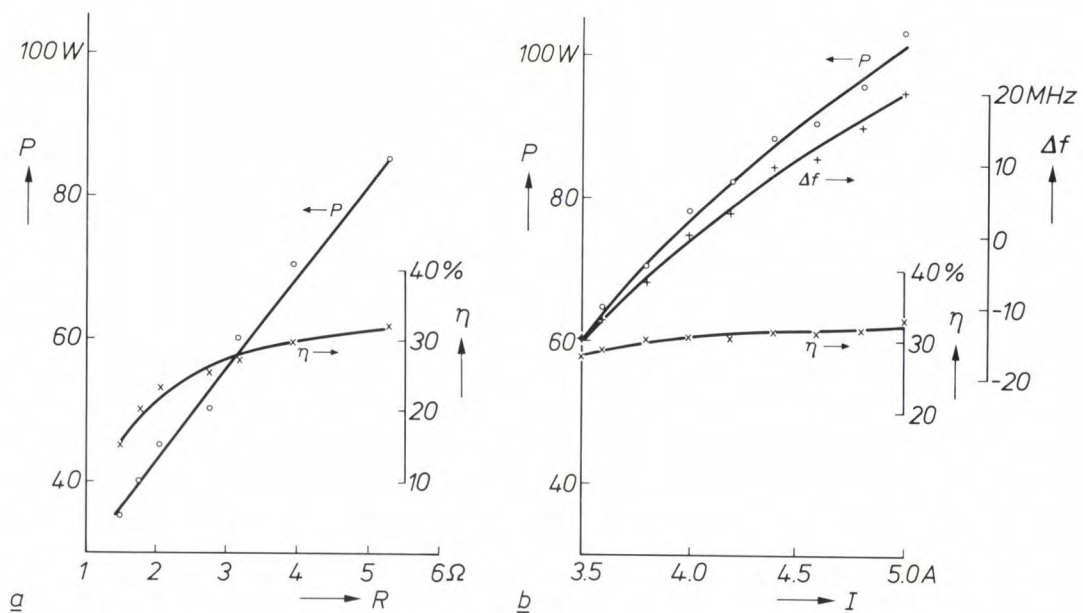


Fig. 10. The dependence of peak power *P* and efficiency η on *a*) the resistive load *R*, and on *b*) the drive current *I*, also showing the frequency variation Δf . In both cases, further increase of the independent variable caused loss of coherence.

Oscillator tuning

It is necessary to provide a simple means of tuning an oscillator if it is to be used in a practical system operating at a specified frequency. Oscillators are commonly tuned by varying the fundamental resonance of the frequency-determining components. The TRAPATT oscillator is extremely rich in harmonics, and it can be tuned by modifying either the fundamental resonance or the resonance of one or more of the harmonics, or by a combination of these methods.

The frequency of a TRAPATT oscillator is determined by the delay introduced by the circuit and the device. To tune the oscillator, the circuit delay can be adjusted by varying the electrical length of the delay line, or the device delay can be varied via the harmonic terminations.

It is worth noting that electronic tuning by means of a varactor diode (which has a voltage-dependent capacitance) is not particularly suited to a TRAPATT oscillator. This is because the varactor diode is essentially a small-signal tuning element, and the high peak power available from TRAPATT oscillators limits the method to very small tuning ranges.

Mechanical tuning

Mechanical tuning is achieved most readily by adjusting the termination of as many harmonics as possible. In the TDT circuit this is accomplished by adjusting the length of the delay line via the position of the low-impedance step. A coaxial oscillator that can be tuned in this way is shown in *fig. 11*. The figure also shows the expected linear relationship between frequency and delay-line length, with an accompanying power variation of about 10 W over the tuning range.

Magnetic tuning

There are two methods of tuning in which an external magnetic field is varied. The required field change in each method is small, and could be obtained from a small field coil. The tuning rate is reasonably high, and may reach several tens of megahertz per microsecond. Both methods of tuning have been investigated, using the same basic TDT oscillator.

The first method is tuning by variation of effective permeability, and was originally investigated for the tuning of IMPATT oscillators^[7] and later for TRAPATT oscillators^[8]. In our experiments we used a microstrip oscillator constructed on a ferrite substrate^[9]. The electrical length of the delay line is a function of the permeability of the substrate material, and it follows that the line length may be varied by the application of an external magnetic field. The construction of the prototype oscillator is shown in *fig. 12*.

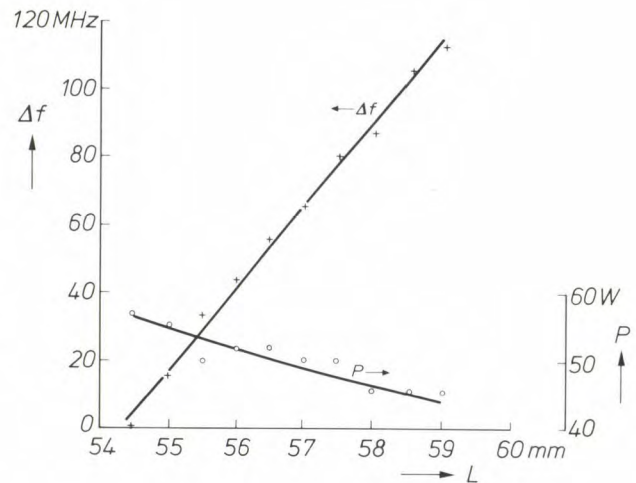
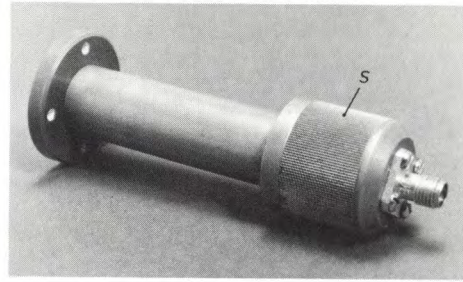


Fig. 11. The photograph shows a mechanically-tuned coaxial oscillator in which the position of the low-impedance step can be manually adjusted by means of the knurled section *S*. A linear variation of frequency Δf of over 100 MHz is obtained when the delay-line length L is changed by about 4.5 mm. There is a reduction in output power P as the frequency rises, but this is acceptably small.

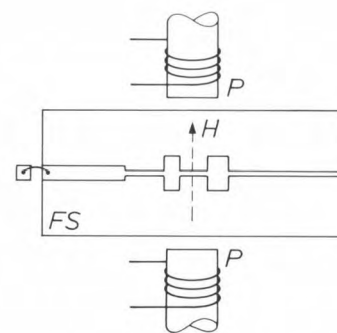


Fig. 12. The TRAPATT oscillator on its ferrite substrate *FS* is located between the poles *P* of an electromagnet. The oscillator is tuned by varying the applied magnetic field *H*.

[7] B. Glance, A magnetically tunable microstrip IMPATT oscillator, *IEEE Trans. MTT-21*, 425-426, 1973.

[8] S.-G. Liu, Magnetically tunable TRAPATT oscillator, 1974 IEEE Int. Solid-State Circuits Conf. Dig. tech. Papers, pp. 98-99.

[9] P. L. Booth, S. R. Longley and B. H. Newton, Frequency tuning of microstrip TRAPATT oscillators, *IEEE Trans. MTT-29*, 6-10, 1981.

Varying the magnetic field by 60 kA m^{-1} gave the results shown in *fig. 13*. The parameter on the graph is the angle θ between the normal to the magnetic field direction and the plane of the substrate. Different settings of θ were investigated to see if an optimum value could be found. The widest tuning range was 61 MHz for $\theta = 18^\circ$, starting at a frequency of 2.19 GHz. The output power at this setting remained substantially flat (13 W peak $\pm 0.5 \text{ dB}$) over the frequency range and the frequency spectrum was clean and symmetrical. The tuning curves show a pronounced 'dip' and are very nonlinear. We have found that a qualitative explanation of the shape of the curves is possible if the demagnetization effect of the substrate is considered.

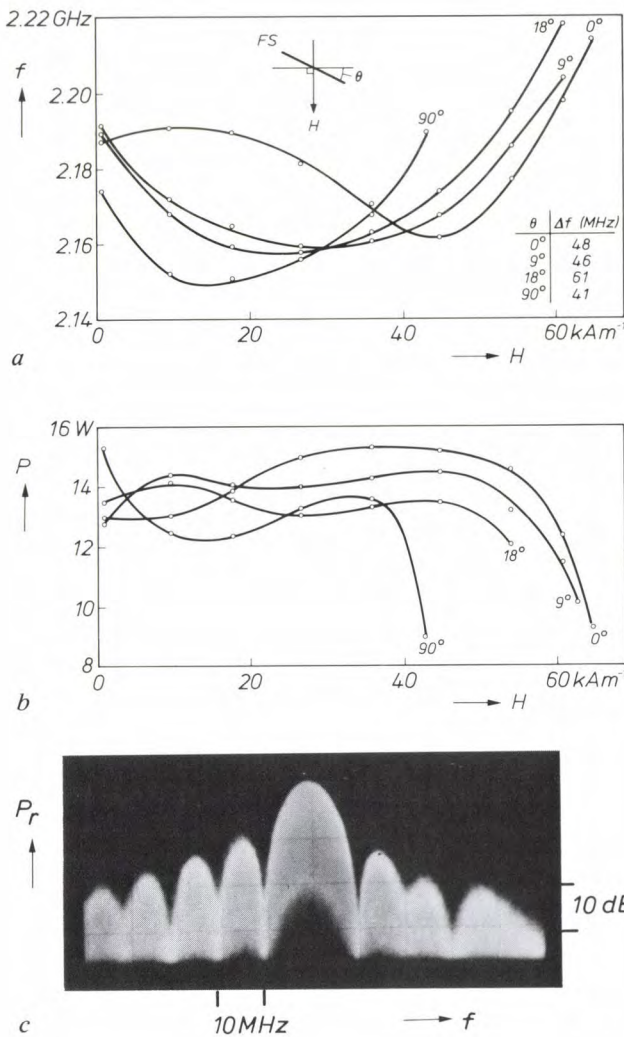


Fig. 13. Experimental results for magnetic tuning by varying the effective permeability. The angle θ between the ferrite substrate FS and the normal to the direction of the magnetic field H was set successively to 0° , 9° , 18° and 90° . *a*) The optimum setting of 18° gave the largest change of frequency f . *b*) This setting also corresponded to the smallest variation in peak output power P . The oscillator was adjusted for good spectral purity rather than for maximum output power. *c*) The power-output spectrum is clean and symmetrical, as can be seen on this photograph of a measured relative power P_r versus frequency f .

The second magnetic-tuning method makes use of the harmonic terminations. By computer modelling, it has been shown^{[10][11]} that frequency tuning is possible by varying the amplitude or phase of the fundamental, 2nd, 3rd and 4th harmonics. This simulation predicted that tuning by phase variation at the 3rd harmonic gives the largest tuning range with the least variation in output power. Altering the amplitude of the 3rd harmonic causes a smaller change in frequency.

Magnetic tuning via the harmonic terminations can be achieved by using polycrystalline yttrium iron garnet (YIG) spheres as a ferrimagnetic resonator, and this is the method we have adopted^[12]. The practical oscillator had two YIG spheres of 0.93 mm diameter. They were semi-loop coupled approximately one quarter of a wavelength from the open-circuit end of an air-spaced microstrip transmission line connected in parallel with the TRAPATT diode chip. The diode was also connected in parallel with the standard microstrip oscillator circuit. This is shown in *fig. 14*. The magnetic field required to tune the ferrimagnetic resonator, which primarily controls the oscillator frequency, was applied perpendicular to the substrate and parallel to the plane containing the semi-loops. The layout of the oscillator was adjusted to give the broadest possible tuning range.

Using this method, we obtained a fundamental-frequency tuning range of over 120 MHz at a centre frequency of 2.46 GHz, by biasing the spheres to resonate near to the third harmonic. The tuning curves for power and frequency are illustrated in *fig. 15*. It is also possible to use a simpler arrangement with a single sphere, but this has a smaller tuning range of about 90 MHz, and produces about 10 W less output power. The d.c. to r.f. conversion efficiency of these circuits is typically between 13% and 23%.

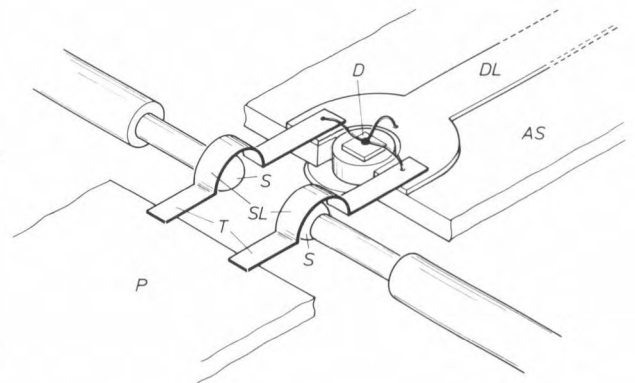


Fig. 14. This figure shows the positioning of the two YIG spheres S , with their semi-loops SL . The quartz plate P provides a mechanical support for the air-spaced transmission line T . The TRAPATT diode D lies at the end of the delay line DL , on an alumina substrate AS .

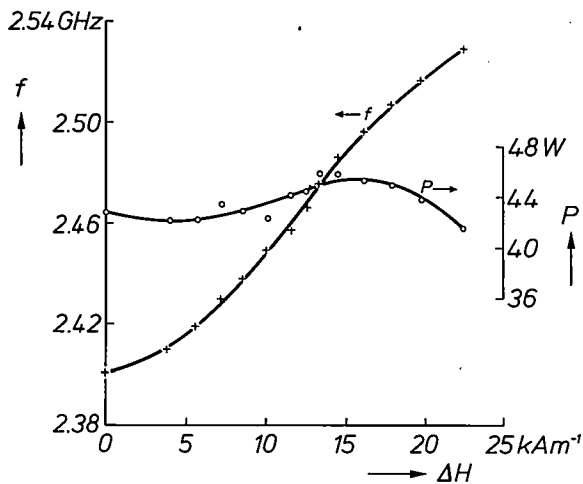


Fig. 15. Variation of frequency f and peak output power P with change in magnetic field ΔH . The nearly linear relationship between frequency and field results from the fact that the sphere resonance frequency also increases linearly with the field.

Compared with the variable-permeability approach, this novel method provides somewhat larger tuning ranges with improved linearity, and demonstrates the feasibility of individually tuning the harmonics of a TRAPATT oscillator.

Frequency-stabilized TRAPATT oscillators

The frequency of a TRAPATT oscillator is a function of circuit and device delays, each of which is temperature dependent. In practice the variation with temperature of the circuit delay is negligible. The device temperature dependence leads to an increased delay with increasing temperature, which results in a negative variation of frequency with temperature, df/dT , of up to $450 \text{ kHz } ^\circ\text{C}^{-1}$ at S-band. This frequency variation manifests itself in two distinct ways. There is a variation or 'chirp' during the r.f. pulse, due to heating of the device by the passage of current, and there is a variation with ambient temperature. Three circuit techniques have been studied for reducing the oscillator frequency dependence on temperature: injection locking, temperature compensation and temperature stabilization. Of these, only injection locking affects the characteristics of the chirp.

Injection-locked oscillators

Typically, at an operating current density of 10 kA cm^{-2} , the chirp during the $0.5 \mu\text{s}$ bias of a free-running 2.47 GHz oscillator was 6 MHz . For most radar transmitters, however, a frequency stability of 1 MHz or better is required. Injection locking is a technique whereby a free-running high-power oscillator is coupled to a stable, low-power oscillator in such

a way that a stable, high-power oscillation is produced. Fig. 16 shows a typical 'locking' characteristic for a microstrip TRAPATT oscillator having a $0.5 \mu\text{s}$ bias pulse and a pulsed locking source. An almost linear relationship exists between the logarithm of the locking range and the locking gain from 15 to 20 dB. At higher gains the relationship becomes nonlinear owing to the inherent frequency chirp. A system has been devised in which the locking source is itself phase-controlled using a discriminator at 1.3 GHz and a crystal reference [13]. Such a system can have a frequency stability of $\pm 5 \text{ ppm}$ or better.

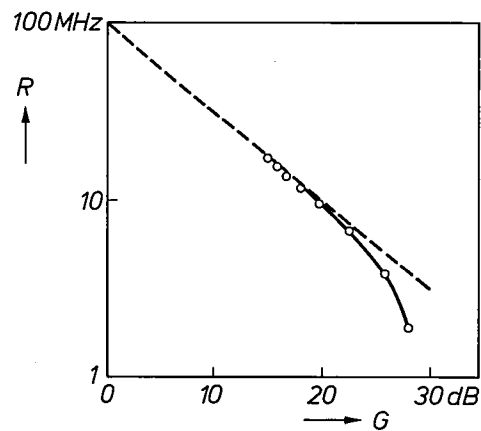


Fig. 16. The locking characteristic is the usual way of representing the behaviour of a free-running oscillator that is coupled to a low-power locking oscillator via a non-reciprocal device such as a circulator. Such circuits produce a stable, high-power oscillation. The locking range R is that range of locking-oscillator frequencies for which the free-running oscillator becomes synchronized. The locking gain G is the oscillator power ratio (free-running/locking). A simple theory predicts that a linear relationship exists between these two parameters, the locking range reducing by 1 decade for every 20 dB increase in locking gain. This is shown by the dashed line in the figure. The solid line shows that the measured performance of a TRAPATT oscillator is in good agreement with this theory. (The measurements apply to a 2.47 GHz microstrip oscillator with a peak power output of 13.3 W , and operating at a duty cycle of 0.1% .)

Temperature-compensated oscillators

The change of frequency of a TRAPATT oscillator with temperature can be compensated, by mechanically retuning the oscillator in sympathy with changes in ambient temperature. This is most easily achieved by arranging for the position of the low-impedance

- [10] R. J. Trew, G. I. Haddad and N. A. Masnari, The operation of S-band TRAPATT oscillators with tuning at multiple harmonic frequencies, *IEEE Trans. MTT-23*, 1043-1047, 1975.
 [11] R. J. Trew, Properties of S-band TRAPATT diode oscillators, thesis, University of Michigan 1975.
 [12] S. R. Longley and P. L. Booth, Frequency tuning of TRAPATT oscillators using ferrimagnetic resonators, 8th Eur. Microwave Conf., Paris 1978, pp. 790-794.
 [13] B. H. Newton and G. Payne, A rugged phase-locked C-band source, Mullard Research Labs Annual Review 1973, pp. 88-94.

step to be temperature-dependent. The thermal coefficient of expansion of different materials can be exploited for this; plastic and metal are used in the compensated oscillator shown in *fig. 17*. This oscillator changed by only 4 MHz in frequency as the ambient temperature was varied from 20 °C to 100 °C, compared with 30 MHz for the uncompensated oscillator [*].

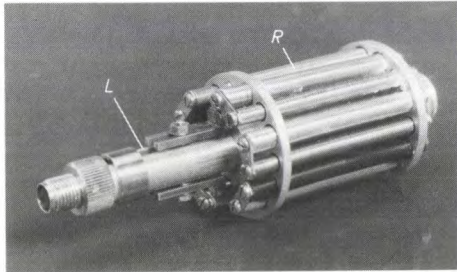


Fig. 17. The thermal expansion of different materials can be used to compensate for the change in frequency of this oscillator as the temperature varies. The rods *R* that surround the coaxial line *L* are divided into two groups of seven each, the groups acting in parallel. In each group, the rods are connected serially (end-to-end), in a zig-zag fashion to save space. One extremity of the group is fixed to the body of the oscillator, and the other to a moveable low-impedance step. By using alternate rods of Invar and a chosen material, an additive expansion is obtained that can be made to match the oscillator characteristics very closely. In this example, the materials used are Delrin and aluminium.

Temperature-stabilized oscillators

The temperature of the diode itself can be stabilized at an elevated value, by using the heating effect of a small reverse current that is allowed to flow through the diode. An external control circuit is employed that decreases the current as the ambient temperature rises. It is sufficient to maintain a simple linear relationship between this current and the temperature. Because the current is applied between oscillator pulses, and is at least an order of magnitude smaller than the threshold current for TRAPATT oscillation, no spurious r.f. power is generated. This technique gives a stabilized frequency over a wide range of ambient temperatures.

Device manufacture

The initial theoretical and experimental work demonstrated the potential of the TRAPATT diode for the efficient generation of high peak power at microwave frequencies. However, it became clear that to produce acceptable devices, careful design was needed to accommodate the large electric field and high current density, both inherent to the TRAPATT mechanism. So in addition to investigating the electrical design, it has been necessary to consider reliability, performance reproducibility and high-yield processing.

The result is that devices have been available for some time for circuit investigations, and are now being developed for specific systems applications.

Diode structures

Unless special precautions are taken, the electric field-strength at the edge of an avalanche diode can be higher than elsewhere, leading to localized premature breakdown ('edge breakdown'), which destroys the device. It is therefore of prime importance to design the diode so that the field is low at the edges, and uniform over the junction region. Two constructions have been investigated that can be made to satisfy this design requirement. These are the mesa and deep-diffused planar constructions.

The SEM (Scanning Electron Microscope) picture in *fig. 18* shows a typical mesa diode. The edge of the

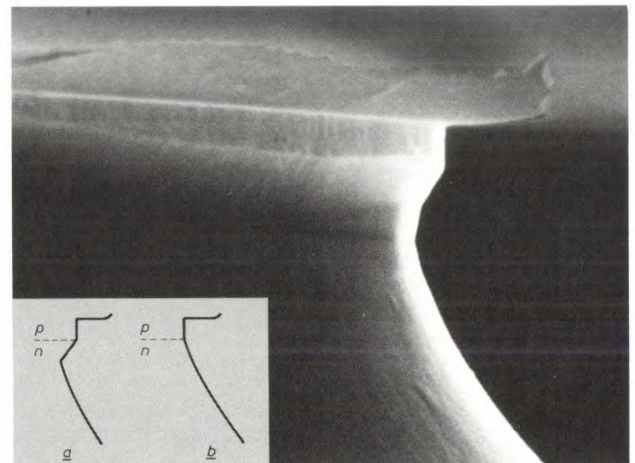


Fig. 18. An SEM picture showing a mesa diode profiled to prevent edge breakdown of the junction. The insert shows a positive bevel (*a*) that ensures that, in the active part of the device, the p-n junction has the largest area. This is not the case for a negative bevel, (*b*). The photograph corresponds to stage *C* in *fig. 19*, without passivation.

diode is bevelled positively so that the p-n junction is always the largest area in the diode. For maximum electrical efficiency, the bevel angle is critical and the surface of the mesa must be smooth. If the mesa is wet-etched these conditions can be difficult to control because of electrochemical action in the vicinity of the p-n junction. In extreme cases this can produce a localized negative bevel. An investigation of etching effects using a dry gas plasma has shown that the bevel angle can be controlled accurately depending on the composition of the etching gas. This has made possible a new approach to mesa technology, in which the diode can be etched from the junction face as opposed to conventional etching from the substrate. Using the processes shown in *fig. 19*, mesas can be accurately defined from the junction face by whole-slice processing

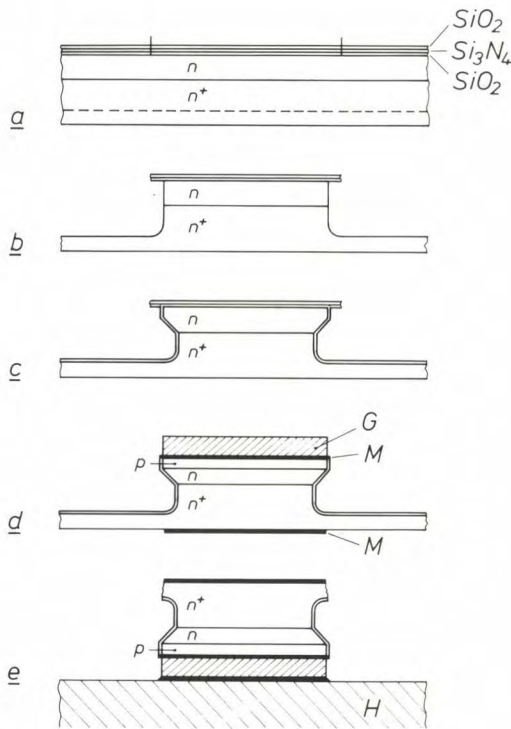


Fig. 19. Processing steps used in producing a plasma-etched mesa diode: *a*) n^+ substrate with n -type epitaxial layer, and masking layers SiO_2 , Si_3N_4 , SiO_2 . Photolithography is used to define the mesa mask. *b*) Mesa definition by plasma etching. *c*) Mesa shaping by selective plasma etching to give the mesa the required outline. SiO_2 passivation is then applied. *d*) p -region formed either by diffusion or ion implantation. Metallization M applied for ohmic contacts, plus a plated gold layer G . *e*) Device inverted on to main heat sink H .

either before or after heat-sink attachment. The resulting devices are completely free of any surface contamination produced during the etching process. An advantage of plasma etching is that junction design and edge-breakdown criteria can be optimized independently during processing.

In 1968, a method was described [14] for making planar IMPATT diodes with improved breakdown characteristics. This was achieved by diffusing the junction through an oxide window to a depth sufficient to prevent edge breakdown. The same method can be used to make planar TRAPATT diodes. In practice, the junction depth is determined by the doping of the n -type silicon and so this requirement and the junction design cannot be independently optimized. The end regions of the diode are therefore characterized by graded profiles, and the depth of the active junction region makes it difficult to provide the diode with an effective heat sink. However, the diode can be made by conventional silicon-planar processing and has been shown to be ideal for high-yield batch processing. Provided long-life metallization schemes are incorporated, this results in reproducible and reliable devices.

A simple, but effective, flip-chip device has been evolved using planar processes for all the fabrication including the heat sinks. This device is well-suited for applications with duty cycles of less than 2%, and typical performance figures are shown in Table I. The limiting factor in the structure is the thick diffused p -region that lies between the n -region, where the heat is generated, and the gold heat sink. This gives a minimum chip thermal resistance of 10°C W^{-1} for a diode area of $4 \times 10^{-4} \text{ cm}^2$. Additionally, the gold heat sink degrades the performance of the diode by reducing the d.c. to r.f. conversion efficiency by about 5%. This occurs because, in the region where the heat sink overlays the oxide, there is a significant oxide capacitance in parallel with the junction, and this capacitance increases the charging time of the diode during the charge-generation period. Thus the design always involves a compromise between thermal spreading resistance and electrical efficiency.

To overcome the electrical and thermal limitations of the flip-chip diode, a more ideal 'thinned' structure has been produced. This is necessary for operation at higher duty cycles. The new device, which is shown in fig. 20, has a lower thermal resistance (5°C W^{-1} for a diode of area $4 \times 10^{-4} \text{ cm}^2$) but the electrical efficiency

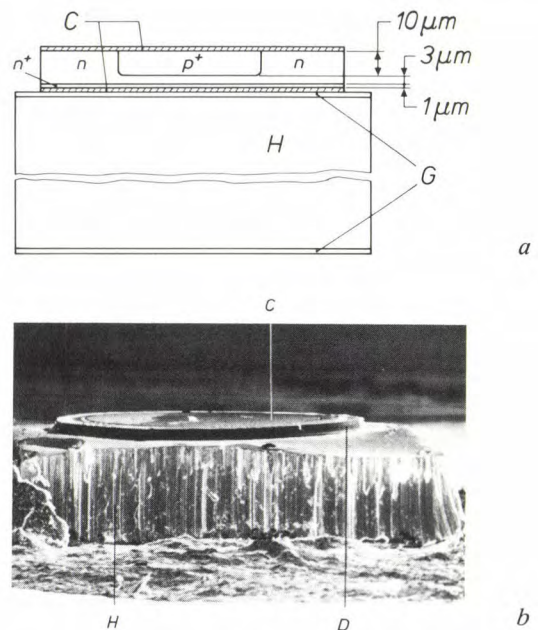


Fig. 20. The 'thinned' TRAPATT diode. *a*) A schematic cross-section showing the separation of only $1 \mu\text{m}$ between the junction region and the silver heat sink H . Ohmic contacts C are provided on each side of the diode, and the heat sink is gold plated (G) above and below. *b*) An SEM picture of such a diode D , showing the silver heat sink H and the upper ohmic contact C .

[*] The temperature-compensated oscillator was designed and measured by G. Tubridy.

[14] H. G. Kock, D. de Nobel, M. T. Vlaardingerbroek and P. J. de Waard, Continuous-wave planar avalanche diode with restricted depletion layer, Proc. IEEE 56, 105, 1968.

is not degraded by the heat-sink arrangement. The device is identical to the planar diode, except that all but 1 μm of the n^+ substrate is removed, and the remaining layer is connected directly to the heat sink. The distance between the active junction and the heat sink is thus reduced from 10 μm to 1 μm , and the associated decrease in thermal resistance has been measured to be greater than 50%. As the main heat sink is connected directly to the n^+ substrate, the stray heat-sink capacitance is completely eliminated, giving an improvement of at least 5% in the electrical efficiency of the diode. A further advantage is that removal of the bulk of the substrate reduces the series resistance to less than 1 Ω , and so the circuit loading is improved. It is then possible to achieve impedance matching to larger area diodes, and so higher output power can be obtained than for other designs. The thinned diode can be designed so that its oscillating frequency is less temperature-dependent than that for a conventional device. A df/dT of only 100 $\text{kHz } ^\circ\text{C}^{-1}$ is possible, compared with 450 $\text{kHz } ^\circ\text{C}^{-1}$ for flip-chip devices. Some thinned devices have shown almost zero df/dT , and work is in progress to improve the diode design until this valuable characteristic becomes standard.

To produce thinned devices, it has been necessary to develop a technology in which whole silicon slices are reduced to a uniform thickness of $15 \mu\text{m} \pm 1 \mu\text{m}$ before the heat-sink plating is formed. A high-yield process has been established, and devices of this type are in development at Mullard Hazel Grove.

The performance and likely applications of flip-chip and thinned diodes are given in *Table II*.

Table II. A comparison between the performance and likely applications of flip-chip and thinned diodes.

	Flip-chip	Thinned device
Output power:	> 60 W (2-3 GHz)	> 60 W (2-3 GHz)
Duty cycle:	2% maximum	> 2%
Efficiency:	35% maximum	> 35%
Applications:	Typically low cost radars, such as altimeters, portable radars, etc.	Suitable for applications requiring longer duty cycles, such as airborne radars. Also capable of further development for phased arrays and marine radars

High-power oscillators

There are basically two ways of increasing the output power of a TRAPATT oscillator, either by changes to the diode structure, or by using a combination of diodes. Of course, both these approaches can also be used at the same time. Changes to the diode

structure include increasing the area of the device and improving the d.c. to r.f. conversion efficiency. Diodes can also be combined in a number of ways, and this is possible either at chip level (an interconnected arrangement of chips within one encapsulation), or at circuit level (an interconnected arrangement of packaged devices).

An important quantity underlying this subject is the impedance of the device itself. For efficient transfer of power, a match must be maintained between the diode or diode combination, and the rest of the oscillator circuit. There is a finite minimum circuit impedance that imposes a limitation on some of these approaches, despite advances in oscillator circuit design. A diode with an intrinsically higher impedance therefore has a clear advantage.

Improved devices

Peak output power is roughly proportional to device area, the practical limitation to power being determined by the lowest impedance that can be matched by the microwave circuit.

All the results that have been discussed so far are based on diodes with an area of $4 \times 10^{-4} \text{ cm}^2$. We have investigated larger devices; in particular, devices having no heat sink, and with junction areas of up to $12 \times 10^{-4} \text{ cm}^2$, have been measured in an S-band coaxial circuit. The oscillator was optimized by adjusting the length of the low-impedance step to maximize the output power. The power, efficiency and associated drive current are shown in *Table III*^[**]. These results all correspond to the same operating frequency, and since the efficiency is maintained in the large-area devices it would appear that appreciably larger devices could be used before circuit losses become significant.

Optimization of the charge generation period of the diode would increase efficiency. With the present diodes there is a compromise in the design between charge generation and charge removal. Refinement of the carrier avalanche processes gives an improvement in the dynamics of plasma formation, and increases the density of the trapped plasma. These effects can be demonstrated by optical illumination of the diode, and one approach may be to integrate a solid-state

Table III. Experimental results for a range of diodes with different junction areas.

Diode area (cm^2)	Peak power (W)	Efficiency (%)	Current (A)
4×10^{-4}	150	22	10
8×10^{-4}	240	25	14
12×10^{-4}	320	28	16

laser for independent control and modulation of the charge generation.

Combinations of devices

Various circuit and device options exist for parallel/series combinations that give more power than a single device. We have investigated several interconnections of two or four small-area devices, at frequencies in the range 2.1 GHz to 2.4 GHz. The circuits used have been in both coaxial and microstrip configurations, and to make power optimization and alignment easier the circuits had mechanical tuning sections as previously described. All our measurements have been made at room temperature, using broadband load resistances, with duty cycles in the range 0.05% to 1%.

Arrangements of two $2 \times 10^{-4} \text{ cm}^2$ diodes connected in parallel give peak output powers between 50 W and 70 W, with an efficiency of about 23%. The construc-

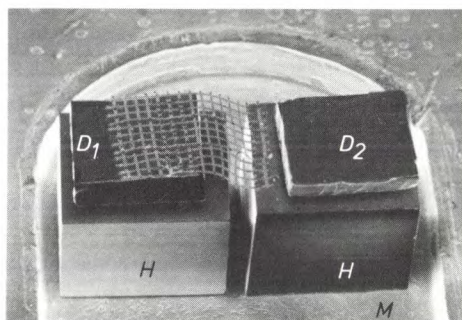


Fig. 21. Two TRAPATT diode chips D_1 and D_2 connected in series. The gold-plated type IIA (high thermal conductivity) diamond heat sinks H are thermocompression bonded to the mount M and the chips are flipped and bonded to the upper surfaces. (Note that the left-hand diamond is metallized on all six surfaces, to provide electrical connection to ground.) Connection to the circuit is made from the back contact of the right-hand chip, but is omitted here for clarity.

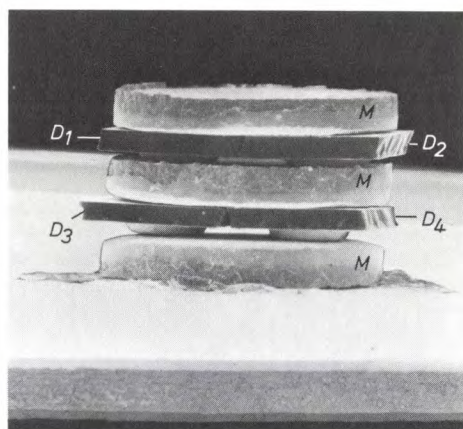


Fig. 22. An SEM picture of a series-parallel assembly of 4 diodes (in pairs D_1 - D_2 , D_3 - D_4) interleaved with gold-plated metal discs M . This combination has the same impedance as that of a single diode, but must operate at low duty cycles to avoid excessive heating of the upper diode pair.

tion of a series-chip connection in the microstrip circuit is shown in the SEM picture of *fig. 21*. The diodes are each of area $4 \times 10^{-4} \text{ cm}^2$, and this arrangement produced a peak output power of 100 W at a drive current of 4 A, with an efficiency of 20%.

Another approach, in which four diodes present an impedance of just one device, is shown in *fig. 22*. Clearly this series/parallel stacked arrangement does not give the ultimate in thermal resistance but the use of diamond offers the possibility of combining several devices and yet retaining a practical impedance level. These diodes combine to give four times the power from a single diode, i.e. a total of between 200 W and 240 W. Because the two uppermost chips have no heat sink, operation was possible only at low duty cycles between 0.05% and 0.1%. The efficiency was about 23%.

For comparison, single diodes from the same batch were operated in standard coaxial and microstrip circuits, giving typically the following power outputs at 4A drive current:

	coaxial	microstrip
area: $2 \times 10^{-4} \text{ cm}^2$	40 W	30 W
area: $4 \times 10^{-4} \text{ cm}^2$	70 W	60 W

The results show that the output power can be readily increased by using combining techniques at chip and circuit level. Series, parallel and series-parallel chip combinations are possible. It can be concluded that several hundred watts of peak power with mean power levels of the order of several watts at S-band are attainable, using these techniques with currently available TRAPATT diodes. The mean power capability of the diamond heat-sink configuration is currently under investigation.

Much of the work described in this article has been carried out with the support of Procurement Executive, Ministry of Defence, sponsored by DCVD.

Summary. The TRAPATT oscillator is a source of high peak powers at microwave frequencies in S-band (2-4 GHz). It is based on the TRAPATT diode, and is particularly suitable for pulsed applications in radar systems. Although operation is at lower microwave frequencies than some contemporary oscillators, efficiencies are high and comparatively large duty cycles are possible. The conventional TRAPATT oscillator circuit has shortcomings that make its use in practical systems intractable; in particular, loss of frequency coherence arising from spurious trigger pulses is a common problem. A new circuit has been developed that obviates these difficulties, and which has a much improved performance. Coaxial and microstrip versions of the oscillator have been made to this new design, and various aspects of their performance have been measured. The oscillator can be tuned either mechanically or magnetically, and can be made insensitive to ambient temperature variation in three different ways. Considerable attention has been paid to the design of diodes with optimum electrical and thermal characteristics; this has resulted in a new 'thinned' diode structure. Very high peak powers can be obtained using both large-area devices and novel circuit configurations.

[**] P. L. Booth and G. Tubridy supplied the results for the high-power devices.

Microwave measurement of moisture content in process materials

W. Meyer and W. Schilz

Ever since man began to engage in manufacture, the determination of the water content of raw materials and commodities has been of great economic importance. Nearly 35 years ago J. Boeke wrote in this Journal: 'The moisture content of grain, of tobacco, of butter, to name only a few examples, is also paid for by the purchaser, and large profits or losses may be involved when a kilogram of the product contains a few grams more or less of water than were estimated' [1]. This is just as true today as it was then. Nowadays, however, there are simpler methods of measurement than Boeke's rather cumbersome extraction method. The authors have succeeded in developing a microwave method that gives a measure of the relative moisture content and does not depend on incidental properties of the sample, such as density and structure.

Introduction

In the control of industrial processes it is necessary to have methods for continuously measuring the moisture content of materials such as domestic salt, fish meal, instant coffee and tobacco. This requires an electrical output signal that can be used to regulate a preceding manufacturing process by means of feedback. There are various methods that satisfy these conditions, employing effects such as the absorption of infrared radiation, the scattering of γ -rays and the conduction of alternating or direct current. Another method uses microwaves, electromagnetic waves with frequencies from 1 to 100 GHz. The difficulty with the microwave methods used until now is that only the absolute quantity of water in the sample is measured. To obtain a 'variable of state', however, that does not depend on the fortuitous properties of the particular sample, we are more interested in the relative moisture content ψ , the ratio of the mass of the water to that of the moist sample [2]. To obtain a value of ψ with the existing methods, the mass of the sample usually has to be determined separately.

The special feature of our method of measuring the moisture content by a microwave method is that the

relative moisture content ψ is directly measured and that the measurement signal is independent of the density, the dimensions and the structure of the sample. In addition, the mean moisture content of the entire sample is measured, whereas with infrared light, for example, only the surface moisture content is determined.

We have developed instruments of two types. With one type the material can be measured continuously during the manufacturing process. With this instrument there is no physical contact with the material being measured, as there is when the conductivity for alternating or direct current is measured. With the other type of measurement very rapid determinations can be made for separate samples such as a bundle of fibres or a cigarette. With an instrument of the first type a large number of measurements have been performed in industrial conditions to test the reliability of the results. The small amount of microwave radiation escaping during the measurement remains well below the statutory safety limit.

We shall first consider the basic theory, and then we shall discuss some examples of the instruments we have developed. Finally we shall present some results of measurements.

Dr W. Meyer and Dr W. Schilz are with Philips GmbH Forschungslaboratorium Hamburg, Hamburg, West Germany.

Theoretical considerations

The purpose of the measurement is to determine the relative moisture content

$$\psi = \frac{m_w}{m_d + m_w}, \tag{1}$$

where m_w is the mass of the water and m_d the mass of the dry matter in the sample. In the determination of m_d by weighing alone, the result depends on the method adopted for drying the sample. When it was necessary to determine m_d , to verify the results of measurements, we dried the sample for an hour at a temperature of 110 °C.

The extent to which the propagation of electromagnetic waves is modified by the material, compared with propagation in free space, depends on the complex relative dielectric constant

$$\epsilon_r = \epsilon' - j\epsilon'', \tag{2}$$

where ϵ' is the real part of ϵ_r and ϵ'' is the imaginary part. ϵ' determines the velocity of the propagating wave in the material; ϵ'' , also called the loss factor, determines the attenuation of the amplitude of the propagating wave as it penetrates into the material.

The values of ϵ' and ϵ'' depend on the density ρ_d of the material, the temperature and the frequency. In the microwave frequency range the complex dielectric constant of water, e.g. at 10 GHz and 25 °C

$$\epsilon_{r,w} = 65 - j35, \tag{3}$$

is much greater than that of the materials commonly measured, which might typically have $\epsilon' \approx 1.5$ and $\epsilon'' \approx 0.01$ in the dry state. The water content therefore largely determines the dielectric constant of the moist material. Disregarding for a moment the effects of temperature, we can say that in the microwave range ϵ_r depends mainly on the density and relative moisture content of the material:

$$\epsilon_r(\rho_d, \psi) = \epsilon'(\rho_d, \psi) - j\epsilon''(\rho_d, \psi). \tag{4}$$

The measurement of ϵ' or ϵ'' separately does not therefore give an unambiguous value from which ψ can be determined, since both are also dependent on the more or less fortuitous density of the sample.

Looking for a quantity that would indeed constitute a unique function of the relative moisture content ψ , we first measured the dependence of ϵ' and ϵ'' on ρ_d for tobacco (see fig. 1) at a constant water content of 18%. It was found that in this material — and, as we later found, in other materials too — $(\epsilon' - 1)$ and ϵ'' are both an approximately linear function of the density ρ_d . The line for ϵ' starts at the point (0,1) and the line for ϵ'' starts at (0,0), since ϵ_r takes the value 1 (with $\epsilon' = 1$ and $\epsilon'' = 0$) when $\rho_d = 0$. A closer look

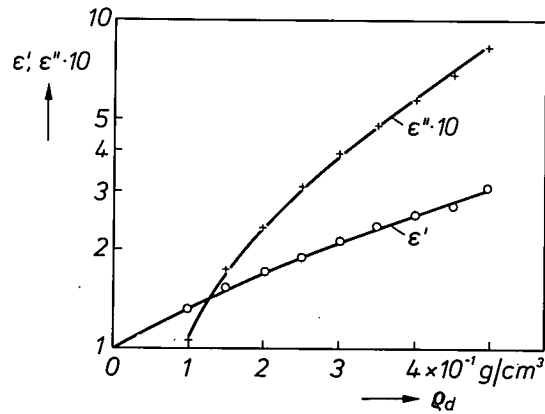


Fig. 1. The real part ϵ' and the imaginary part ϵ'' of the complex relative dielectric constant $\epsilon_r = \epsilon' - j\epsilon''$, as a function of the density ρ_d of tobacco. ϵ' and ϵ'' were measured by filling a waveguide with different tobacco samples with a relative humidity of 18% and determining the velocity and attenuation of a propagating wave. The measurements were carried out at a frequency of 12.5 GHz. The values of $(\epsilon' - 1)$ and ϵ'' both turn out to be approximately linear functions of ρ_d . Since $\epsilon_r = 1$ if $\rho_d = 0$, the line for ϵ' starts at (0,1) and that for ϵ'' at (0,0). The lines in the figure are not straight because ϵ' and ϵ'' are plotted on a logarithmic scale.

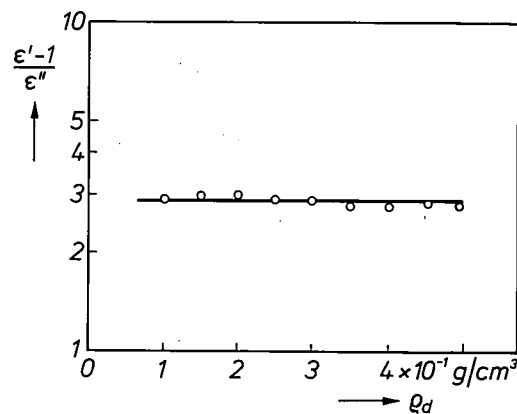


Fig. 2. Variation of the quantity $(\epsilon' - 1)/\epsilon''$ as a function of ρ_d for the results of the measurements in fig. 1. The relation corresponds approximately to a straight line parallel to the abscissa. $(\epsilon' - 1)/\epsilon''$ is thus virtually independent of the sample density.

at the behaviour of $(\epsilon' - 1)/\epsilon''$ clearly indicated that this quantity was indeed found to be practically independent of the density ρ_d ; see fig. 2. Measurements on other materials such as wheat grain, feathers, instant coffee, tea, domestic salt, fish meal and wool also showed that $(\epsilon' - 1)/\epsilon''$ was independent of the density within certain ranges. It was found^[3] that the quantity

$$A(\psi) = \frac{\epsilon' - 1}{\epsilon''} \tag{5}$$

is a measurable electrical magnitude that can yield a

[1] J. Boeke, A new electrical method for determining moisture content, Philips tech. Rev. 9, 13-15, 1947.

[2] In exceptional cases the relative moisture content may also be related to the sample without water. We shall not deal with this aspect here.

[3] W. Meyer and W. Schilz, J. Physics D 13, 1823, 1980.

unique value for the relative moisture content. Within practical limits, $A(\psi)$ is independent of ρ_d , and appears to be only slightly dependent on temperature. Fig. 3 gives the various calibration curves of $A(\psi)$ for some of the materials that have been measured in our laboratory. They start at the ordinate with different $A(\psi)$ -values, depending on the properties of the dry material, and are all a monotonically decreasing function of ψ . At moisture contents of approximately 30% the curves approach the value $A \approx 1.8$ for water alone, which follows from eq. (3). At even higher moisture contents the dielectric behaviour of water predominates, and this method is then no longer suitable for moisture measurements. The method is not very suitable for moisture contents below 3%, since most of the curves are nearly parallel to the ψ -axis in that range, corresponding to the properties of the dry material.

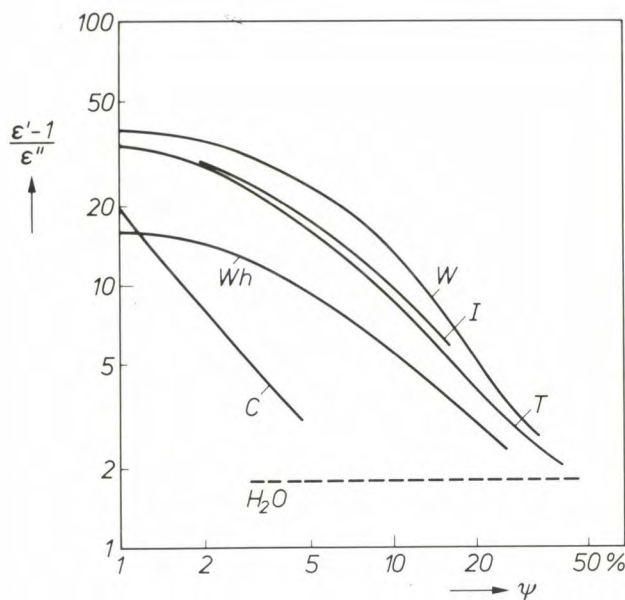


Fig. 3. Some calibration curves for different materials currently analysed at 12.5 GHz. T the curve for tobacco, C for domestic salt, I for instant coffee, W for wool and Wh for wheat grain. The quantity $A(\psi) = (\epsilon' - 1)/\epsilon''$ is plotted as a function of the relative moisture content ψ . When ψ is greater than 30%, $A(\psi)$ is no longer useful as a unique measurement value for the moisture content, since all the curves then approximate to the value $A(\psi) = 1.8$ for water. The method is not very suitable for moisture contents below 3%, since most of the curves are nearly parallel to the ψ -axis in that range, corresponding to the properties of the dry material.

At frequencies of about 10 GHz $A(\psi)$ is only affected by the free water in the sample. At lower frequencies — outside the microwave range — trouble is experienced from ion conduction and the result of a measurement may depend on (say) the salt content of the sample.

Industrial moisture measurement

A great deal has been published on the subject of microwave moisture measurement [4]. Various instruments are also on the market that measure either the attenuation or the velocity of the propagated waves in the moist material. These instruments differ not only in the microwave frequencies they use, but also in the manner in which the microwave signal is introduced into the sample. The test unit must be adapted to suit the dimensions of the sample, e.g. sheets of paper or cigarettes. Depending on the type of test unit, instruments for moisture measurement can be divided into two types:

- transmission instruments, mainly used for measuring large quantities of material of average or high moisture content, e.g. on conveyor belts; and
- cavity instruments, suitable for measurements on small samples of low moisture content, such as cotton fibres.

We shall now discuss two examples of moisture meters that we have developed, which are of the first and the second types. In both designs we applied the principle of obtaining an output signal that does not depend on the density of the sample.

A transmission-type moisture meter

Fig. 4 shows the moisture-sensor head (sometimes called an 'applicator') we have designed for moisture measurements by microwave transmission. The microwave signals are coupled to and from the sample by horn transitions. This unit can be used for measuring fluffy materials such as tea and tobacco. Fig. 5 shows how the instrument can also be used for measurements of materials on a conveyor belt. The ratio of the phase shift ϕ and the attenuation L of the amplitude of the wave propagated through the sample depends on ϵ' and ϵ'' [5]:

$$\frac{\phi}{L} = \frac{\epsilon' - 1}{\epsilon''} \frac{2\sqrt{\epsilon'}}{1 + \sqrt{\epsilon'}} \quad (6)$$

This means that when $\epsilon' \approx 1$ the ratio ϕ/L approxi-

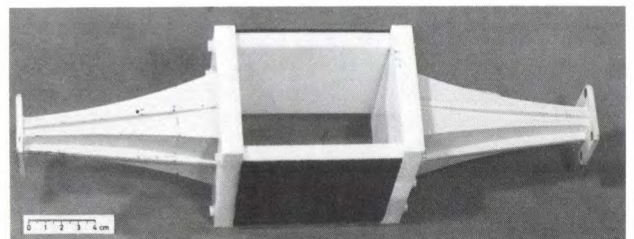


Fig. 4. The moisture-sensor head (the 'applicator') for measuring fluffy materials by means of microwave transmission. The microwave signals are coupled to and from the sample by the horn transitions.

mates to the value of $A(\psi)$ as indicated by eq. (5). Since ϵ' is not greater than 1.5 for the dry state of most of the materials commonly measured ϕ/L is a measure of the relative moisture content. A microwave bridge arrangement is used for determining ϕ and L . Fig. 6 shows the relation between ϕ/L and the relative moisture content ψ for tobacco. The curve resembles the curve for tobacco in fig. 3 after compensating for the logarithmic ψ -scale. Fig. 7 shows for moisture contents of 12 and 18% that ϕ/L is also virtually independent of the density ρ_d . Here again

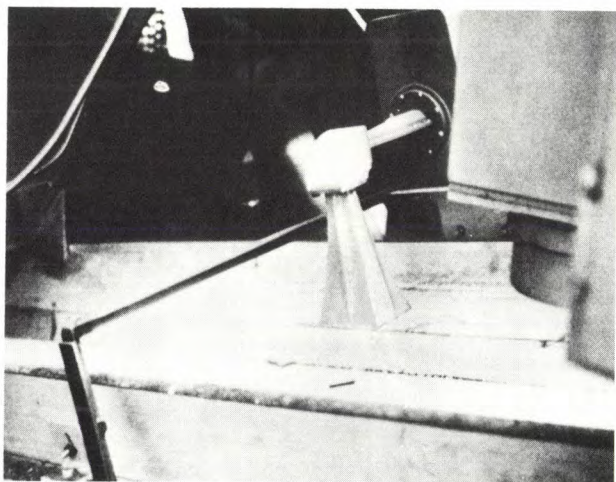


Fig. 5. Measuring the moisture content of tobacco on a conveyor belt by microwave transmission. One of the horn transitions of fig. 4 can be seen at the centre. The equipment for drying the tobacco is at the upper left. An infrared moisture meter is shown at the upper right. The arrangement was designed for comparing the microwave and infrared methods.

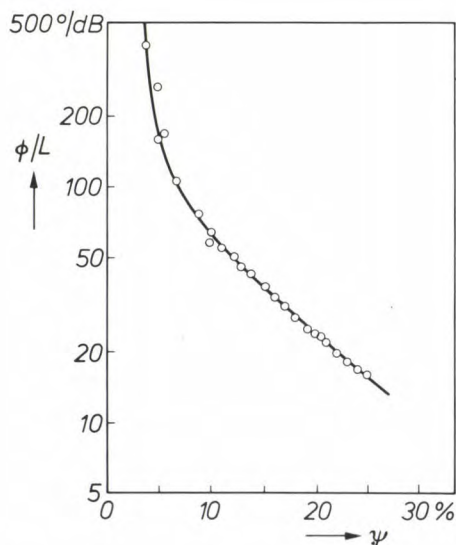


Fig. 6. The ratio of the phase shift ϕ (in degrees) and the attenuation L (in dB) as a function of the relative moisture content ψ for tobacco. The measurements were made with a transmission-type instrument, as in fig. 4; the values for ϕ and L were measured with a bridge arrangement. The shape of the curve agrees with the curve for tobacco in fig. 3, after correction for the logarithmic ψ -scale of fig. 3.

the method is seen to be of practical use for moisture contents of 3 to 30%.

For use in industry we have developed an instrument in which the electronics and the microwave circuits (except the sensor head) are combined in a single module; see fig. 8. The test unit, consisting of two

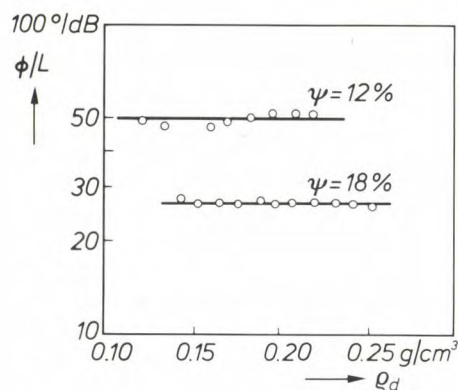


Fig. 7. The ratio ϕ/L (see caption to fig. 6) as a function of the density ρ_d of tobacco, for relative moisture contents of 12 and 18%. It can be seen that ϕ/L , like $(\epsilon' - 1)/\epsilon''$ (see fig. 2), is virtually independent of the sample density.

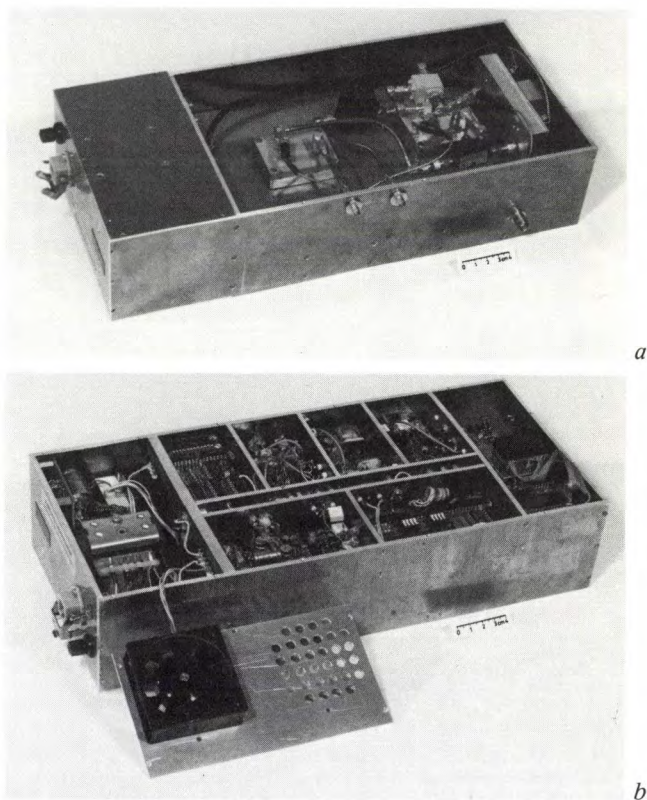


Fig. 8. Photographs of the main module of the transmission-type moisture meter. a) The upper side with the microwave unit. The two Gunn oscillators (see fig. 9) are on the right. The connections for the two coaxial cables to the sensor head (see fig. 4) can be seen just below the centre. b) The underside with the electronics.

[4] A. Kraszewski, Microwave aquametry — a bibliography, 1955-1979, J. Microwave Power 15, 298-310, 1980.
 [5] W. Meyer and W. Schilz, IEEE Trans. MTT-29, 732, 1981.

horn transitions of the type shown in fig. 4, can be connected to the unit of fig. 8 with coaxial cables. Fig. 9 gives the block diagram of the electronic circuits. From this it can be seen that there are three sections operating at different frequencies — 9 GHz, 10 MHz and 10 kHz. The microwave part of the diagram contains the sensor head, which can be placed around the sample S , and two Gunn oscillators (G_1 and VCO). One of the oscillators can be varied in frequency (VCO : voltage-controlled oscillator). The oscillators have nominal frequencies of 9 and 9.01 GHz. Their frequency difference is kept accurately constant in a phase-locked loop, using a reference oscillator G_2 at

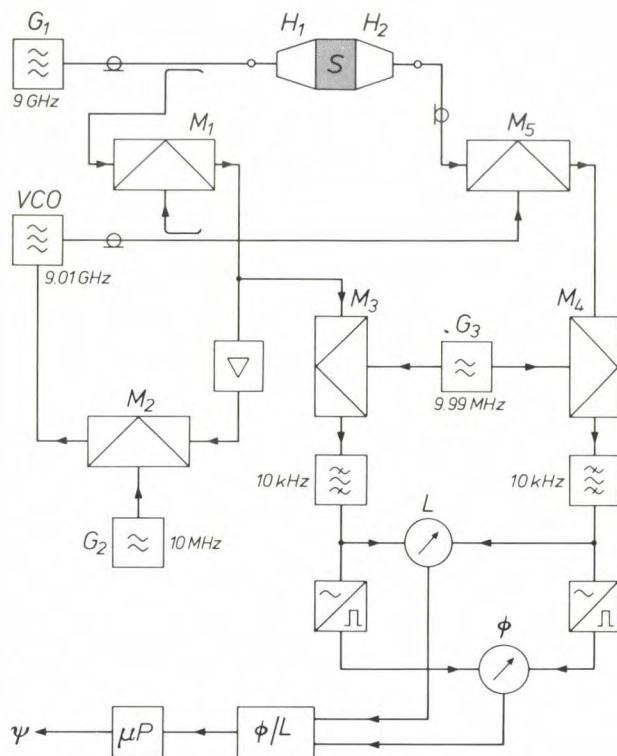


Fig. 9. Block diagram of the moisture meter of fig. 8. S sample. G_1 and VCO Gunn oscillators with frequencies of 9 and 9.01 GHz respectively. The microwave signals at 9 GHz pass through the sample via the horn transitions H_1 and H_2 of the sensor. The 9 and 9.01 GHz signals are coupled out of the coaxial cables, and a difference signal at 10 MHz is obtained in a mixer M_1 . The difference signal is compared in a phase-locked loop via an amplifier and a mixer M_2 with a 10 MHz signal from a reference oscillator G_2 with a quartz crystal. In this way the difference of the frequencies from G_1 and VCO is maintained at 10 MHz exactly. The mixer M_5 provides a 10 MHz signal that contains information about the attenuation L and phase shift ϕ of the 9 GHz signal in the sample. A second reference oscillator G_3 delivers a signal at 9.99 MHz. The mixers M_3 and M_4 provide signals at 10 kHz that contain information about the 9 GHz signal before and after the sensor head respectively. After passing through two filters, which transmit only the sine waves at 10 kHz, both signals are compared. This yields a digital signal proportional to the attenuation L . The sinusoidal 10 kHz signals are converted into pulsed signals and compared, yielding a digital signal proportional to the phase shift ϕ . The signals pass through a divider and the signal for ϕ/L is then processed in a microprocessor, taking into account the relation between ϕ/L and ψ shown in fig. 6. In this way a signal proportional to the relative moisture content ψ is obtained, which can be used for controlling the preceding stages of the manufacturing process.

10 MHz with a quartz crystal. Two signals at 10 kHz are derived from this with the aid of a second reference oscillator G_3 at 9.99 MHz and two mixer circuits (M_3 and M_4). One of the 10 kHz signals contains the information from the 9 GHz signal before the sensor head, the other contains the information from the signal after the sensor head. Filters and conventional electronic circuits are used for obtaining signals for the phase shift ϕ and the transmission loss L . After the ratio ϕ/L has been derived in a divider circuit a microprocessor finally delivers the signal that is proportional to the relative moisture content ψ by making use of a calibration curve similar to that of fig. 6. This signal may if required be used for controlling the manufacturing process. The accuracy of the instrument is about $\pm 1^\circ$ for ϕ and about ± 0.8 dB for L .

A moisture meter of the cavity type

Measurement of $A(\psi)$ in eq. (5) and using resonant cavities is relatively simple if the samples are small. In this case perturbation theory^[6] indicates that the relation between ϵ' and ϵ'' and the resonant frequencies and Q-factors is given by

$$A(\psi) = \frac{\epsilon' - 1}{\epsilon''} = 2 \frac{f_2 - f_1}{\frac{f_2}{Q_1} - \frac{f_1}{Q_2}}, \quad (7)$$

where f_1 and f_2 are the resonant frequencies, and Q_1 and Q_2 are the Q-factors before and after insertion of the sample.

The moisture meter we have built permits an even simpler measurement than by determination of the frequencies and Q-factors as mentioned above. Fig. 10 shows the sensor head. It consists of a cylindrical cavity with two Gunn oscillators coupled to it. Fig. 11 shows a number of lines of force for the electric and magnetic fields of the H_{011} and E_{012} modes of the cavity. Each Gunn oscillator is coupled to the cavity in such a way that it excites only one of the two modes. The two modes differ little in frequency and lie a little above and below 11.5 GHz. The location of the sample coincides with the centre-line of the cavity. In the E_{012} mode this location corresponds to that for the maximum amplitude of the electric field. In the H_{011} mode the amplitude of the electric field is zero at this position. Since the sample only affects the electric field, insertion of the sample does affect the E_{012} resonance, shifting it in frequency, but the H_{011} resonance remains unchanged. The signal from the H_{011} mode of the cavity is therefore applied to the local-oscillator port of a balanced mixer via a coaxial cable. The signal from the E_{012} resonance is applied to the other

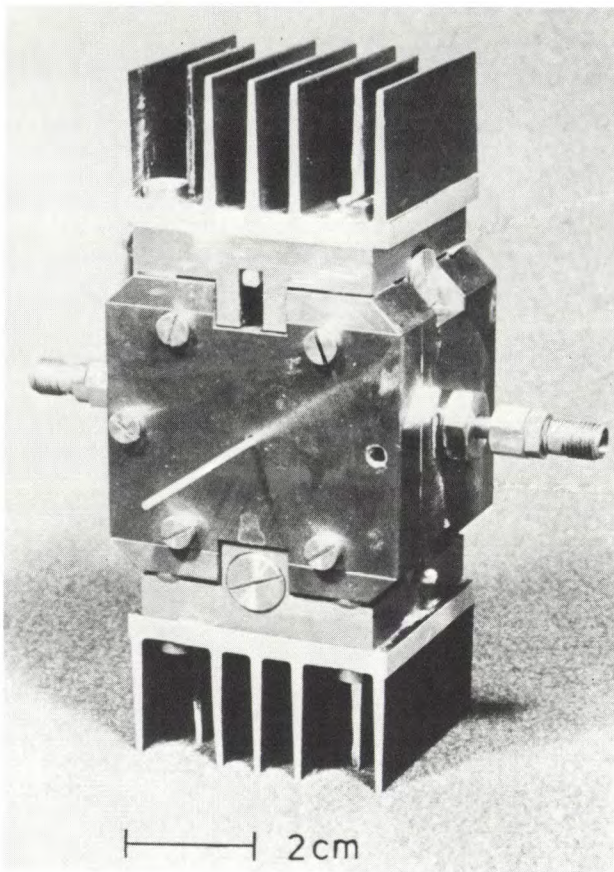


Fig. 10. The sensor head of the cavity-type moisture meter. The central part is the actual cylindrical resonant cavity. The sample, here consisting of a bunch of cotton threads, is placed in the centre of the cavity. Above and below there are two Gunn oscillators, which are coupled to the cavity by means of a 'cut-off' technique. One oscillator excites the E_{012} mode, the other excites the H_{011} mode of the cavity (see fig. 11). At the left and right are the couplings used for extracting the power from each mode, using coaxial cables.

port of the mixer via a similar cable. At the output of the mixer the power P of the E_{012} oscillation and the difference f between the frequencies of the two oscillations is measured. Since the resonances in both modes are affected in much the same way by temperature changes, f is practically independent of temperature.

A special 'cut-off' coupling method maintains the output power from the Gunn oscillators constant to within 0.1 dB over a frequency range of at least 50 MHz. The power injected into the cavity by the Gunn oscillator for the E_{012} resonance thus changes hardly at all if the frequency of this resonance is changed by the insertion of the sample. This is of fundamental importance, since it means that the power extracted from this mode of the cavity, is a monotonic function of ϵ'' — the loss factor — of the sample and hence of its water content. Since ϵ' determines the velocity of a propagating wave, the frequency of the standing wave in a cavity also depends on ϵ' . The

change Δf in the frequency difference f measured at the output of the mixer is therefore a monotonic function of $\epsilon' - 1$ of the inserted sample. We thus see that, by measuring the ratio $\Delta f/\Delta P$ after insertion of the sample we have again obtained a measure of the water content, which, like $A(\psi) = (\epsilon' - 1)/\epsilon''$ does not depend on the density. In fig. 12 both $(\epsilon' - 1)/\epsilon''$

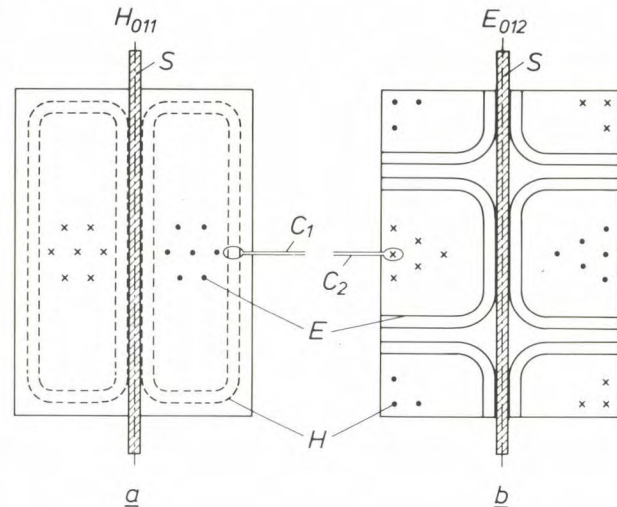


Fig. 11. a) Lines of force for the H_{011} mode of resonance, b) those for the E_{012} mode of resonance in the cavity of fig. 10. S sample, placed at the centre of the cylinder. The crosses and dots indicate lines of force perpendicular to the plane of the drawing, away from the reader and towards him, respectively, E lines of electric force; H lines of magnetic force. The electric field is always zero at the location of the sample in the H_{011} mode. Insertion of the sample thus only affects the oscillations in the E_{012} mode. C_1 schematic indication of the coupling to the coaxial cable for extracting the signal from the H_{011} mode; the loop should be considered as perpendicular to the plane of the drawing; C_2 the same for the E_{012} mode; the loop is shown in the correct position.

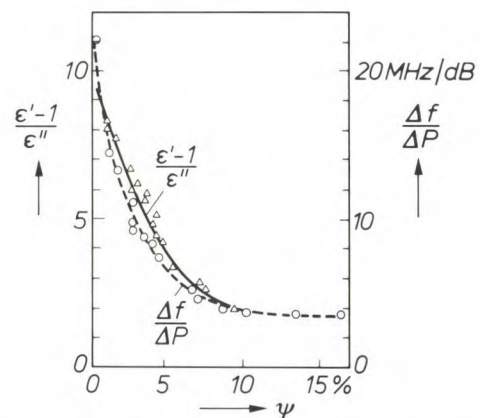


Fig. 12. Plot of $A(\psi) = (\epsilon' - 1)/\epsilon''$ and $\Delta f/\Delta P$ as a function of the relative moisture content ψ . The changes Δf and ΔP are respectively the detuning and the power loss of the oscillations in the E_{012} mode of fig. 11 after the sample has been inserted. $A(\psi)$ can be determined by means of the cavity by measuring the Q and the frequency of the oscillations in one of the possible modes before and after insertion of the sample and then applying equation (7). It is clear that the ratio $\Delta f/\Delta P$ — which is easier to determine — is just as useful a measure of the relative moisture content as $A(\psi) = (\epsilon' - 1)/\epsilon''$.

[6] W. Meyer, IEEE Trans. **MTT-25**, 1092, 1977.

— obtained by measuring the change in frequency and Q of the cavity and applying equation (7) — and $\Delta f/\Delta P$ are plotted as a function of the relative moisture content ψ . It can be seen that the functions behave in almost the same way and that the method is usable for moisture contents up to about 10%. From fig. 13 it follows that ΔP and Δf are each extremely

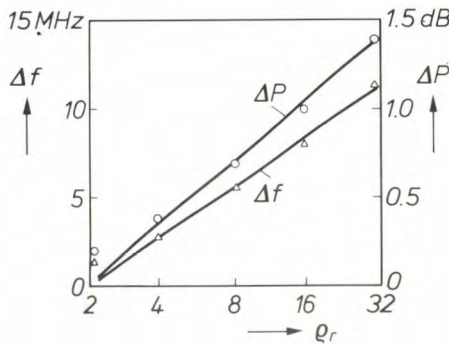


Fig. 13. The power loss ΔP and the detuning Δf (see caption to fig. 12) as a function of the relative density ρ_r (varied to $32\times$ with respect to an arbitrary density) of a bunch of cotton threads. Both quantities are strongly dependent on ρ_r and therefore not suitable separately as a measure of the moisture content of the sample. The relative moisture content of the cotton threads was 4.5% during the measurement.

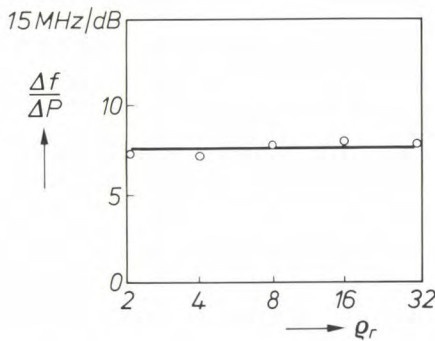


Fig. 14. The ratio $\Delta f/\Delta P$ as a function of ρ_r for the same measured points as those in fig. 13. It is clear that $\Delta f/\Delta P$ is almost independent of the density, and can readily be used as a measure of the moisture content that does not depend on the fortuitous properties of the sample. $\Delta f/\Delta P$ remains within 7.5 ± 0.2 MHz/dB, corresponding to a variation of $4.5 \pm 0.1\%$, see fig. 12, in the relative moisture content ψ , while the density varies by a factor of up to 32.

dependent on the density of the sample, of course, which in this case is cotton thread with a water content of 4.5%. The density here was varied by a factor of 32. Fig. 14 shows a plot of the ratio $\Delta f/\Delta P$ for the same measurement data; the ratio is seen to lie within 7.5 ± 0.2 MHz/dB for the entire density region. This corresponds to a variation of $4.5 \pm 0.1\%$ in the measured relative moisture content; see fig. 12.

Industrial trials

Moisture-measuring instruments of the transmission type (fig. 4) and of the cavity type (approximately corresponding to the one shown in fig. 10) were subjected to extensive trials under industrial conditions, and the results of the measurements were evaluated statistically.

The most promising results were obtained with the transmission-type instrument, mounted above and below a conveyor belt for tobacco, as illustrated in fig. 5. Also visible in that photograph is an instrument that works with the aid of infrared radiation. The output signals from both instruments were recorded, and at regular intervals samples were removed from the conveyor belt. The moisture content of these samples was determined by weighing, using the drying procedure mentioned earlier. It was found that the curves



Fig. 15. Different structures in tobacco samples. From top to bottom can be seen successively leaves, stalks and the cut finished product. Moisture measurement in the tobacco industry should be affected as little as possible by these different structures.

in both recordings agreed well with the moisture content determined by weighing. The structure of tobacco can vary considerably owing to the presence of leaf, stalks and the cut end-product; see *fig. 15*. It turned out that the results of the measurements using microwaves were less affected by the varying structures of the material than when infrared radiation was used.

The instrument shown in *fig. 4* was also tested in the fish-meal industry. Fish meal also has a widely varying structure, owing to the presence or absence of bones. Here again, a moisture-content determination is required that does not depend on the fortuitous structural and chemical properties of the sample. Ten measurements, each on 15 samples of different kinds of fish meal were carried out at 9 GHz with the moisture meter of *fig. 4*. The moisture content of each sample was then determined by weighing. The results

of the measurements are shown in *fig. 16*, the vertical lines indicating the spread of the measurements and the circles the average for each sample. The quantity plotted along the vertical axis is L/ϕ , the reciprocal of the quantity given by equation (6). A linear regression analysis gave as a result the equation

$$\psi = 5.29 \frac{L}{\phi} + 4.02 \quad (8)$$

for the calibration 'curve', with a correlation coefficient of 0.98 and a standard deviation of 0.42% for the relative moisture content. However, if the phase shift ϕ alone is plotted as a function of ψ , a 'cloud' of points is obtained showing little or no correlation. The measurement of the ratio L/ϕ is thus much more useful for determining the moisture content than measurements of ϕ or L separately.

The advantages of measuring moisture content by our microwave method during the manufacturing process can best be seen by comparing our method again with the infrared method. Both methods are non-destructive, make no mechanical contact with the sample and are continuous. Microwave signals, however, have longer wavelengths and are therefore less affected by the different structures of the sample than infrared rays. Both methods compensate for the effects of varying density in the sample, but this requires the use of a second radiation source as reference in the infrared method. The main advantage of using microwaves, however, is that the signals pass right through the sample and thus provide an average value for the relative moisture content of the entire sample. Infrared rays, on the other hand, only scan the surface of the material.

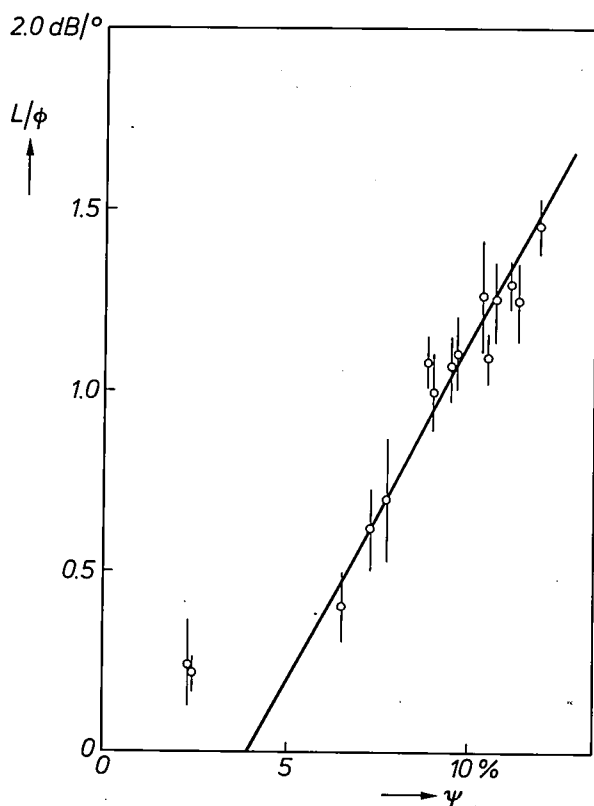


Fig. 16. The ratio L/ϕ (the reciprocal of the quantity given by equation (6)) of the attenuation and the phase shift as a function of the relative moisture content ψ . The results are given of 150 measurements, made with a microwave moisture meter as in *fig. 4*, on 15 samples of fish meal of different kinds. The ψ of each sample was then determined by weighing. The average value of the results of the measurements for each sample is indicated by a circle, while the vertical lines indicate the spread in the results. The best-fit straight line for the 'calibration curve' is given by the equation $\psi = 5.29 L/\phi + 4.02$; it was calculated by a linear-regression analysis. The correlation coefficient is 0.98 and the standard deviation is 0.42% in relative moisture content.

Summary. By using microwave methods to determine the quantity $(\epsilon' - 1)/\epsilon''$ (ϵ' and ϵ'' are the real and imaginary components of the complex dielectric constant), a measure of the relative moisture content of the material is obtained. This measure is found to be independent of the density of the material. An instrument based on microwave transmission is described, which is especially suited for measuring the moisture content of raw materials during the manufacturing process. Another type of instrument measures separate samples in a microwave cavity. In this case the detuning and attenuation of one of the resonances provides a measure of the moisture content. The article concludes with the results of a number of trials of the instruments in an industrial environment, and a comparison is made with the results of moisture meters using infrared radiation.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- P. M. van den Avoort:** Socio-technical research. Conf. Rec. Int. Conf. on Communications, Boston 1979 (ICC '79), pp. 47.3.1-47.3.5. *E*
- E. G. Berns & J. de Ridder** (Philips Lighting Division, Eindhoven): Calculation of the tungsten and heat transport in spherical gas-filled incandescent lamps. Philips J. Res. 35, 173-189, 1980 (No. 3).
- R. N. Bhargava, R. J. Seymour, B. J. Fitzpatrick & S. P. Herko:** Donor-acceptor pair bands in ZnSe. Phys. Rev. B 20, 2407-2419, 1979 (No. 6). *N*
- K. H. J. Buschow:** Crystallization and magnetic properties of amorphous Gd-Fe and Er-Fe alloys. J. less-common Met. 66, 89-97, 1979 (No. 1). *E*
- T. A. C. M. Claasen & W. F. G. Mecklenbräuer:** The Wigner distribution — a tool for time-frequency signal analysis, Part I: Continuous-time signals. Philips J. Res. 35, 217-250, 1980 (No. 3). *E*
- H. Duifhuis, L. F. Willems** (both with Institute for Perception Research, Eindhoven) & **R. J. Sluyter:** Pitch in speech: a hearing theory approach. Speech Comm. Papers 97th Meeting Acoust. Soc. Amer., ed. J. J. Wolf & D. H. Klatt, pp. 245-248, 1979. *E*
- J. M. L. Engels:** The drift velocity of excess electrons in compressed methane. J. Electrostatics 7, 233-237, 1979. *E*
- H. Hieber:** Dünne Metallschichten sind nicht stabil. Phys. Blätter 35, 455-459, 1979 (No. 10). *H*
- H. de Lang & N. H. Dekkers:** A posteriori focussing of a STEM using shadow image processing. Optik 53, 353-365, 1979 (No. 5). *E*
- J. Magarshack:** All FET 12 GHz satellite television receivers. Conf. Rec. Int. Conf. on Communications, Boston 1979 (ICC '79), pp. 26.3.1-26.3.5. *L*
- A. Mircea, D. Pons & S. Makram-Ebeid:** Depleted layer spectroscopy. New developments in semiconductor physics, Proc. Int. Summer School, Szeged 1979 (Lecture Notes in Physics 122), pp. 69-96; 1980. *L*
- R. F. Mitchell & C. K. Davis:** Traffic handling capability of trunked land mobile radio systems. Conf. Rec. Int. Conf. on Communications, Boston 1979 (ICC '79), pp. 57.2.1-57.2.5. *R*
- E. M. A. M. van der Ouderaa, H. Vrieling & A. Willemse:** COMA, a single-chip communication interface for distributed microcomputer systems. Microprocessors and their applications, ed. J. Tiberghien, G. Carlstedt & J. Lewi, pp. 305-314; North-Holland, Amsterdam 1979. *E*
- Ph. Piret:** Wire-tapping of a binary symmetric channel. Philips J. Res. 35, 251-258, 1980 (No. 3). *B*
- R. J. van de Plassche & D. Goedhart:** A monolithic 14-bit D/A converter. IEEE J. SC-14, 552-556, 1979 (No. 3). *E*
- U. Rothgordt:** Halbtondarstellung mittels nichtmechanischer Aufzeichnungsverfahren. ntz-Archiv 1, 201-204, 1979 (No. 9). *H*
- W. Schnitker & H. Rau:** The stabilizing layer on iron particles for tapes. IEEE Trans. MAG-16, 14-16, 1980 (No. 1). *A*
- M. F. H. Schuurmans & D. Polder:** Superfluorescence and amplified spontaneous emission: a unified theory. Physics Letters 72A, 306-308, 1979 (No. 4/5). *E*
- J. M. Shannon:** A majority-carrier camel diode. Appl. Phys. Letters 35, 63-65, 1979 (No. 1). *R*
- A. Thayse:** Implementation and transformation of algorithms based on automata, Part II: Synthesis of evaluation programs. Philips J. Res. 35, 190-216, 1980 (No. 3). *B*

Gilles Holst, pioneer of industrial research in the Netherlands

H. B. G. Casimir

On the occasion of the twentieth anniversary of the Eindhoven University of Technology in 1976 it was decided to institute an annual 'Holst Lecture'. Its organization was made possible by a donation from the Philips Company. The central theme of each lecture was to be the development of the technological sciences in interaction with the natural sciences and mathematics on the one hand, and with their industrial applications and the consequences of such applications on the other. The choice of name pays tribute to the memory of Prof. Dr G. Holst who, as a distinguished physicist in the Netherlands and founder of Philips Research Laboratories, did much to promote that interaction.

Since then there have been five Holst Lectures. The speakers at the first four lectures were Prof. Alexander King, Prof. C. Freeman, Prof. Dr C. F. Frh. von Weizsäcker and Prof. Kevin Lynch. The main subject of the fifth Holst Lecture was Holst himself, and the organizers were fortunate in finding the best imaginable speaker on this subject — Prof. Dr H. B. G. Casimir, from 1942 a member of Holst's staff and later one of his successors.

The opportunity to publish Professor Casimir's paper here (with a very few minor changes, mainly his own), gives the editors particular pleasure in as much as Holst was in fact the founder of Philips Technical Review.

Who was Gilles Holst? Those who knew him, who felt his influence, who had the privilege to work under his leadership, will never forget him. But their number is beginning to diminish. Anyone undertaking an in-depth study of the history of scientific research in the Netherlands, and in particular of the history of industrial scientific research, will undoubtedly come across his name, but there are few engaged in such studies. It is good to use the occasion of this annual lecture to perpetuate his name, which is in danger of becoming to many no more than a name. The name of Van der Waals immediately conjures up the equation of state, the condensation of gases, the Van der Waals forces between molecules; when Kamerlingh Onnes is mentioned we think of low temperatures, of liquid

helium, of superconductivity; Lorentz is associated in our minds with electrons, with the Lorentz force, the Lorentz transformation, the Lorentz contraction; Zeeman with the Zeeman effect. In textbooks you will probably look in vain for the name of Holst: although the work he published did play a role in the general progress of science, it cannot be counted among the permanent pillars of modern physics. And yet, he was one of the great men of his generation. I am therefore happy to have this opportunity today to try and recreate for you the figure of Holst.

Gilles Holst was born at Haarlem on 20th March 1886. His father, Casper Hendrik Holst, was the manager of the Conrad Shipyards. Holst went through higher secondary education in Haarlem, took his leaving certificate in 1903 and then worked for six months in his father's shipyard and six months with

Prof. Dr H. B. G. Casimir was formerly a member of the Board of Management of N.V. Philips' Gloeilampenfabrieken and headed scientific research in the Philips group of companies worldwide.

the firm of Willem Smit. During this initiation into the business world he must already have realized that in the long run mechanical engineering was not going to satisfy him. 'I didn't like cast iron' was later one of his *obiter dicta*, and this at first sight perhaps rather mysterious remark comprises a great deal of Holst's way of thinking and working. What irked him was that, although you can make splendid calculations for mechanical structures, in the final analysis you had to allow for a substantial safety margin because not enough was known about the properties of the construction material, particularly of cast iron. It was this very problem of understanding and technologically controlling the properties of materials that was later to become a guiding principle of Holst's work.

In 1904 he went to Zürich to study at the Technological University there, initially electrical engineering. Was there still too much cast iron in electrical engineering? After two and a half years and an intermediate examination he switched to mathematics and physics. In 1908 he received the degree of *geprüfter Fachlehrer*, qualifying him as a teacher in higher education. At all events, fundamental physics interested him — and mathematics, too, for that matter. Marcel Grossmann (1878-1936) was from 1907 a professor of geometry at Zürich. He was a good friend of Einstein's and his scientific adviser in matters of invariant multidimensional geometry. I seem to recall Holst telling me that he had attended Grossmann's lectures on '*Ausdehnungslehre*' as propounded by Grassmann (1809-1877). After graduating, Holst spent about a year as an assistant of Professor H. F. Weber [*]. He then returned to the Netherlands, was for a short time an assistant of Siertsema at Delft and at the beginning of 1910 joined Kamerlingh Onnes as an assistant at Leiden.

It is perhaps amusing to mention in passing that Einstein in 1900 received the same degree as Holst did later, but that Weber did not want him as an assistant. One might wonder whether Einstein would have been the most suitable man to lead a practical physics course for up-and-coming engineers, but Einstein himself appears to have been somewhat affronted by this rejection and even went so far as to complain that Weber was intriguing against him. Even more remarkable is that Einstein at the time applied for a job with ... Kamerlingh Onnes, on a carelessly handwritten postcard, reply postage paid. Kamerlingh Onnes did not use it; I don't suppose he took the application seriously.

The Leiden cryogenic laboratory occupied a unique position in 1910. In 1908 Kamerlingh Onnes had succeeded in liquefying helium, thereby opening up a

completely new field of research. It led as early as 1911 to the discovery of what still constitutes one of the most remarkable phenomena of the solid state: superconductivity. It think it is worth while looking a little more closely at the role that Holst played in this sensational discovery.

Electrical-resistance measurements at liquid-helium temperatures (1-4 K) formed part of the research programme Kamerlingh Onnes had embarked upon. When preliminary measurements on platinum revealed that, contrary to all expectations at the time, a further decrease of resistance took place upon cooling to liquid-helium temperatures, Kamerlingh Onnes published a theoretical argument providing a plausible explanation for a very rapid decrease. The measurements were prepared and carried out by Dr Dorsman and Holst, who were given the job of measuring the resistance of mercury. At least, it is reasonable to assume that that is what happened. At the end of the first publication [1], which did establish a rapid drop in resistance upon cooling but no complete disappearance of resistance, that is to say no superconductivity, Kamerlingh Onnes wrote:

'I gratefully record my indebtedness to Dr. C. Dorsman for his intelligent assistance during the whole of this investigation, and to Mr. G. Holst, who conducted the measurements with the Wheatstone-bridge with much care.'

These or similar acknowledgements were repeated in subsequent publications, in particular those that describe superconductivity. Two years later Kamerlingh Onnes concluded a publication [2] with the words:

'Having completed the series H of my experiments with liquid helium I wish to express my thank to Mr. G. Holst, assistant at the Physical Laboratory, for the devotion with which he has helped me, and to Mr. G. J. Flim, chief of the technical department of the cryogenic laboratory, and Mr. O. Kesselring, glassblower to the laboratory, for their important help in the arrangement of the experiments and manufacturing of the apparatus.'

Much later, in a confidential document in which he recommends Holst for appointment to membership of the Netherlands Royal Academy of Sciences (Holst in fact became a member in 1926) Kamerlingh Onnes wrote:

'Of his work at Leiden it should be mentioned that he cooperated in the discovery of the superconductivity of metals and in the further research on superconductivity.'

What are we to conclude from all this? It seems clear to me that Holst carried out the resistance measurements and that he was thus the first to observe

superconductivity. May we then regard him as a co-discoverer of superconductivity? 'Of course not,' Kamerlingh Onnes would probably have replied. 'I provided the liquid helium,' he might have argued, 'I set the task of measuring the resistance of mercury — a measurement which, it is true, was not so easy, given the special circumstances, but which could nevertheless be performed with known methods. A competent experimenter, performing the task I had set him, was bound to discover superconductivity. Holst was undoubtedly competent, and devoted. For that, I, Kamerlingh Onnes, thanked him in the correct manner.' So, it seems to me, he would have reasoned. He was an authoritarian man, he believed in ranks and stations of life, though within that scheme of things he behaved correctly and with benevolence. Nowadays people think differently about these things; I believe, incidentally, that even at that time there were laboratories that pursued a different line of conduct.

How did Holst react to all this? He never spoke to me about it, but it can scarcely have been anything but a disappointment to him. When you have had an important share in one of the most sensational discoveries of that time — which was also promptly rewarded with a Nobel prize — when you have had to convince your superior that your observations were correct — obviously, everyone thought at first it was a case of short-circuiting — then it is not satisfactory to have your name mentioned only in a polite but non-committal word of thanks.

One thing is certain: in later years Holst himself never followed the example of Kamerlingh Onnes in this respect. If his name appears at the author or co-author of a publication, you can be sure that he took an active part in the research work involved. Of course, the times had changed; and in any case, Holst's collaborators were graduate independent researchers whereas those of Kamerlingh Onnes were still in a sense undergraduates, because in those days someone with a master's degree was not regarded as a full graduate. However this may be, Holst might very often have claimed, quite rightly, that he as the principal instigator of a research project was entitled to publish as a co-author. The fact that he did not do so was a noble and generous reaction to a painful experience. It must not be thought that Holst was completely indifferent as to whether or not his name was mentioned. I at least know of one case in which a member of his staff himself insisted that Holst should publish as co-author, and that was certainly appreciated.

A curious remark about Kamerlingh Onnes occurs in a letter dated 30th November 1913 which Holst wrote to his friend A. D. Fokker. He writes:

'I believe that rather underestimates the merits of the great KOH — evidently a codename for Heike Kamerlingh Onnes — because he is certainly a man who plays a big role, even if it is often only the role of bath superintendent.'

This is a characterization that Kamerlingh Onnes would not himself have appreciated. In the same letter Holst writes that he is going to Eindhoven on 1st January. With a dash of bravura he writes:

'So you see, Philips has realized where his own interest lies. That's what you would expect of a good businessman. It sounds to me very nice. I shall have a brand new laboratory to fit up and will have to do all sorts of measurements aimed at revealing to us *the* formula of the incandescent lamp. With that in view I have concluded my work at Leiden. This means that the thesis will not be what I had wished, but that can't be helped.'

That thesis, which gained him his doctorate at Zürich in July 1914, dealt with the equation of state and the thermodynamic properties of ammonia and methyl chloride — a typical Leiden topic, and a thorough piece of work to which he was probably not passionately committed. In a letter to Fokker of 10th May 1914 he even refers to it as his 'accursed thesis'. But on 24th September he writes — again to Fokker, who had meanwhile been called up — that the doctoral ceremony at Zürich, where he was presented by Pierre Weiss (the well-known authority on magnetism, who worked at Strasbourg after World War I) and which was attended by Kamerlingh Onnes, had been a jolly party. Shortly after the doctoral ceremony, while Holst was still looking around in Switzerland, World War I broke out and Holst, with some difficulty, had to take a roundabout way home. So much for his apprentice years.

On the eve of World War I Philips' Gloeilampenfabrieken was a firm of some substance and enjoyed an established reputation. Production was going well — the work of Gerard Philips — and the selling organization was thriving — the work of Gerard's younger brother, Anton. From a technological point of view, however, the firm was in a vulnerable position. It had, it is true, been able to adopt processes from abroad, to master them, make them suitable for mass production and to establish accurate manufacturing formulae, but the really new ideas, such as the metal-filament lamp and the gas-filled lamp, came

[*] The unit of magnetic flux, Wb, is connected with W. E. Weber, 1804-1891. (Ed.)

[1] H. Kamerlingh Onnes, Comm. Phys. Lab. Univ. Leiden No. 119, p. 26, 1911.

[2] H. Kamerlingh Onnes, Comm. Phys. Lab. Univ. Leiden No. 133d, p. 68, 1913.

from elsewhere. In the Netherlands up to 1910 there was no patent law in force, and that made it easy to take over inventions made abroad. After 1910 it became desirable for Philips to build up their own position in patents. However this may be — and here I should remark that there are no documents extant that give any real clue to Gerard Philips's thinking on this matter — Gerard Philips decided that he, like General Electric, had to have a physics research laboratory. Finding the right man seems to have been quite a problem. W. J. de Haas, who later succeeded Kamerlingh Onnes, liked to relate that he was offered the job, but that after taking a look around in Eindhoven, he had declined the honour. I rather doubt the truth of this story. It is also likely that there were negotiations with the Dr Dorsman mentioned earlier, which led to nothing. At all events, on Thursday 23rd October 1913 the following advertisement appeared in *De Nieuwe Rotterdammer*:

Wanted, a skilled, young
Doctor of Physics,
must be a good experimenter.

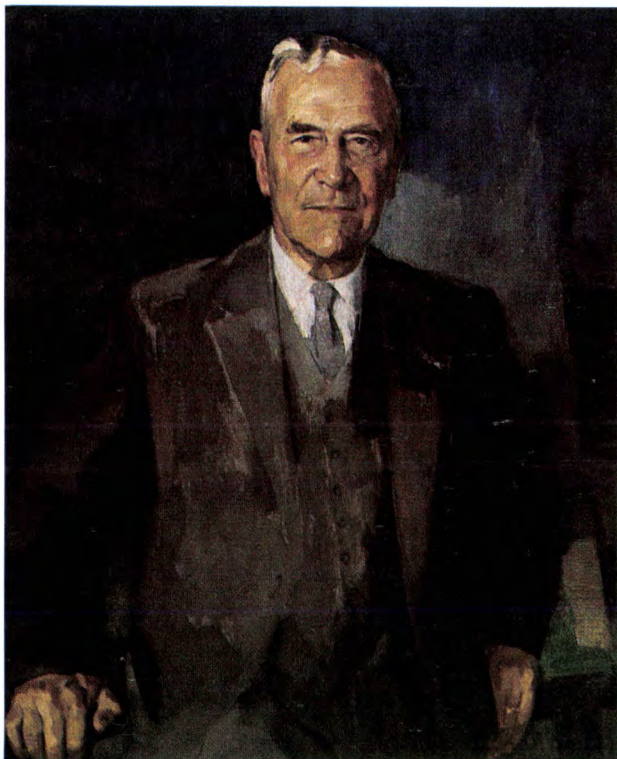
On 25th October Holst sent in his application for the job, or rather he applied for an interview. Gerard Philips and he soon reached agreement, and on 2nd January he was put on the payroll. Notice that it was 2nd January. New Year's Day was a public holiday and Gerard Philips was not the man to start a business relationship by paying someone a day for nothing. This rather tight-fisted way of doing business was continued into later years. My colleague Rinia used to complain that he thought he had been taken on as from 1st November, but at the end of the first month he was paid for only 29 days: the 1st of November was a public holiday. On his 25th anniversary with the Company he was handed three Dutch half-crowns by way of back payment.

After only a few months the laboratory staff was doubled: on 16th April Dr Ekko Oosterhuis joined the Company. In the letter of 10th May, already cited, Holst comments: 'It is much fitter that there are now two of us. You can discuss things much more easily and also — which is the view of the boss — better use is made of the instruments.' In his letter of 24th September Holst is already speaking of a triumvirate: Sophus Weber, a Danish physicist, had joined the team.

If it were my intention to lay stress on Holst's own direct contributions to physics, I would have to speak in particular about the early years, let us say the first ten years, of our Research Laboratories. It was in those years that Holst personally did a great deal of

research work: on photometry, on the properties of tungsten, on the sputtering of metals, on thin films. At first this work had a direct connection with lamp manufacture, but gradually it acquired a broader scope. I shall not deal here with Holst's own work. I assume that it was useful and of good substance, but as I remarked at the outset, it was not of permanent significance. And, as the laboratory expanded, it receded more and more into the background. The great significance of Holst lay not in his own research work but in the way in which he organized and gave leadership to research at Philips. I have already spoken of Holst's unselfish attitude with regard to publications. More in general one can say that Holst set aside his personal ambitions: the work of the Laboratories was more important to him than personal fame. This kind of renunciation is never easy, and definitely not for a man like Holst who, as I have said, had been closely involved in exceptionally important fundamental research and who undoubtedly possessed the capacities for going on doing pioneering work himself. 'I used to be quite a good experimenter,' he once said to me, and his words carried a touch of wistfulness. Van der Pol, one of his best-known colleagues, never wanted to make such a sacrifice; he remained first and foremost a specialist in the mathematical treatment of 'radio'.

I shall now speak first about areas of research that Holst himself selected and on which he was able to set people to work. First of all, gas discharges. The study of the conduction of electricity by rarefied gases and of the phenomena associated with it had a respectable history. Towards the end of the last century it had already led to the discovery of X-rays and of the electron. Holst believed in the possibilities of entirely new light sources, and his conviction was not just some airy notion without foundation. Holst had in fact realized that it was Bohr's atomic theory that had made it possible to understand the phenomena, in particular the emission of light, and *therefore* there were bound to be technological applications. And indeed, this research bore much fruit, including contributions to fundamental physics, as for example the work of Dorgelo and of Penning (and here the work of Holst and Oosterhuis themselves should not be forgotten), and contributions in the field of applications, such as the sodium lamp, the super-high-pressure mercury vapour lamp and the Penning gauge, the principle of which also underlies modern high-vacuum pumps. As regards the sodium lamp, I was told the following story. When it was demonstrated for the first time, the observers looked at the strange yellow light with rather puzzled faces: it was certainly an efficient light source, but what could you do with it?



'Gilles Holst, ... in the painting by Van de Molengraft, ...'

Until someone produced a small railway timetable, one of those booklets that was notorious for its almost illegible small print. Contrast and legibility were surprisingly good. 'Right,' said Holst, 'road lighting'. And that was what it became.

A second area of research was that of thermionic valves (or tubes) and radio. This was not so much concerned with completely new concepts, such as quantum theory. Classical physics was adequate for describing the movement of electrons in electromagnetic fields. The starting point here was rather technological than fundamental: thanks to the manufacture of incandescent lamps Philips had mastered vacuum technology and the technology of glass-to-metal seals. The first valves were in fact commonly called radio lamps. Holst pursued the work in this field, even though Gerard Philips was apparently not very keen on it. But Anton was, and when Gerard retired in 1922 — he was then over 63 — and Anton took over the reins, it became easier to persuade the Philips top management that radio could become a very important field of activities for Philips.

The electrons, which do the work in a thermionic valve, are emitted by a cathode, originally a bare, incandescent tungsten filament, later an oxide cathode. The investigation of cathodes and, more generally, the

study of electron emission, became and has remained an important area of research in the Philips laboratories.

Whereas work on valves is concerned with thermionic emission, that is to say the emission of electrons from a hot cathode, work on photocells and image converters has to do with photo-emission, and this effect was also thoroughly investigated. Holst himself had a very active share in the design of the first image converters, tubes with which an infrared image is converted into a visible image.

At the beginning of the thirties Holst realized that the new quantum mechanics led to an understanding of the behaviour of electrons in solids, at that time not yet understood, and he considered that the time had come to embark upon research into the solid state. Outstanding successes were achieved in the field of permanent magnets, and even more important was the work on ferrites. You could say that there was a happy coincidence here between a tangible problem of technology and Holst's realization that the time was ripe for tackling it. Philips Research Laboratories had done a lot of work in connection with a carrier-wave system, which, as you know, is a system in which a number of telephone conversations can be sent simultaneously along one cable. Important elements in such a system are the electrical filters (or filter networks) whose job is to pass alternating current in a specific frequency band and stop everything outside that band. Filters of this type are built up from capacitors, resistors and coils (inductors). Those early coils had to be rather bulky in order to meet the requirements imposed on them. I have the impression that the annoyance this caused was mainly what gave Holst the idea that we ought to start looking for insulating materials possessing good magnetic properties. Such materials would have no eddy current losses and it would be possible to use them for high frequencies. Easier said than done. Lengthy research was needed before suitable materials were finally found among the ferrites. This turned out to be one of the greatest successes from Philips Research Laboratories.

With these few briefly outlined examples I have tried to show you how Holst selected the areas of research for the Laboratories. As far as the details are concerned, much depended on the initiative of the staff.

Holst was primarily a physicist. Though his objective was technological, physical effects and their explanation were his point of departure. His shots aiming at direct commercial possibilities were sometimes off target. For instance, he saw a great future for the electric bicycle. The success of the moped shows that

there was indeed a market for a motorized bicycle, and perhaps we shall some day have a battery that will make the electric bicycle an attractive proposition. With the batteries of those days it was far from attractive.

Holst also saw the importance of microdocumentation and considered that, even for the individual reader, a cigar box with microcards was going to oust the bookcase. Here again, he was ahead of his time, but he believed wrongly that the essential problem lay in finding a grain-free emulsion. Although interesting results were achieved in that field, they were rather irrelevant to the development of microdocumentation. The problems were to be found principally in publishing and also in reading. No-one likes reading microfilms or microcards. That was the problem to which much more attention ought to have been paid.

Sometimes the good was the enemy of the better — to modify a Dutch proverb. For instance, the fact that Philips were highly skilled in coil winding delayed the introduction of the superheterodyne principle in radio receivers. And the mechanical-optical system of sound recording developed by Philips — known as the Philips-Miller system — was so good that research into magnetic recording was started rather late.

Minutes have survived from the thirties in which Holst and some other wise men made the following predictions concerning sound and image. 'Gramophone records do not have much of a future: people will prefer to listen to the radio. Television is ruled out for general use: much too complicated and much too expensive. On the other hand there is a great future for the home film-show. Cheap copies of sound-cine-films will therefore have to be made.' Precisely the opposite has happened. Will Holst's home film-show finally materialize with the video cassette and the video long-play record?

I want to leave it at that as far as the substance is concerned. Let me now take a closer look at Holst's ideas on research management. I don't think that this is a term he would like to have used, nor did he like long disquisitions on such matters. Nevertheless, from his deeds and from a few talks he gave one can deduce some general viewpoints. In doing so we must be careful to bear in mind that at the time Holst entered Philips employment, joining an industrial company was a very unusual step for a Dutch physicist to take. As far as that goes, there were in general very few university graduates engaged in industry then. Holst was not only faced with the task of setting up a laboratory, he also had to convince his younger colleagues that working with Philips was a good thing.

I tried on a former occasion to summarize his views in the form of ten commandments. I hope that they

do not do too much violence to his ideas, and that perhaps these very 'commandments' may provide a suitable starting point for further discussion, in which it may be asked to what extent they can fully be upheld today. Here, then, are the ten commandments, with some commentary.

1) *Take on bright researchers, preferably young but with experience in academic research.*

This point was by no means self-evident. In Britain until recently the view was held in many industries that the last thing you should do was to hire a physics Ph.D. with experience in fundamental research. Such a character was ruined for practical work! Holst wanted on the contrary to achieve the highest possible scientific level.

2) *Do not pay too much attention to the details of the work they have done.*

In other words, don't hire specialists on account of their specialities. To take the case of Dr Haantjes, for example, who had taken his doctorate with a thesis on the separation of neon isotopes by distillation at low temperature, Holst set him to work on television, and with resounding success. An English colleague once asked me: 'Would you really take on a Ph.D. in nuclear physics if you didn't intend to work on nuclear physics?' 'Of course,' I answered, and I was able straight away to give him examples. And yet, one may wonder whether one can still keep to that rule in this day and age: in many subjects the mass of drudgery you have to wade through before you can do anything new is so dense that it is tempting to want to profit from a level of attainment already reached.

3) *Give your staff a lot of freedom and put up with their peculiarities.*

An important point that applies with equal force today.

4) *Let your staff publish and take part in national and international scientific activities.*

This is an essential point. Holst knew that his best people would leave — or would never have come — if they did not enjoy great freedom in this respect. But that is not the only consideration. Through publications the laboratory acquires a reputation that makes it possible to maintain good contacts with the academic world. As a result, of course, many of his staff were offered professorships. That could mean a loss in the short term, but Holst never complained about it.

5) *Avoid too rigid an organization. Let authority rest on real ability.*

Not everyone will agree with me on this. But scientists tend to be opinionated people who do not readily submit to rules, especially when the rules are imposed by someone they do not respect as an intellectual equal.

6) *Do not divide a laboratory into distinct departments — mathematics, physics, chemistry, etc. — but form multidisciplinary teams.*

In theory the university ought to be the ideal place for bringing about cooperation between different disciplines. In reality the universities have been sadly lacking in this respect, and still are. The Philips laboratories set an example of how it should be done.

7) *Allow great freedom in the choice of work, but see that the leaders in particular are aware of their responsibility to the Company.*

For Holst that responsibility was self-evident. I don't know what is thought about it today.

8) *Do not budget for an industrial laboratory 'project by project', and do not let manufacturing departments have any say in the budgeting for research programmes.*

This independence of the industrial laboratory in matters touching research programmes and budgets was for Holst a *sine qua non* for really advanced research. I am so completely in agreement that I do not think it necessary to comment on or defend the proposition.

9) *Encourage the transfer of competent older researchers from the laboratory to development and production in the factories.*

The remarkable thing is that the converse — transferring from the factory a man who wants to devote himself in his maturer years to fundamental research — only happens sporadically.

10) *Let the choice of subject be co-determined by the state of academic science.*

I have already dealt with this point at some length.

After this rather dry synopsis, let us return to Holst himself. He went through two very difficult periods, the years of the great slump (the world economic crisis) and the German occupation. During the great slump the growth of the Laboratories came to a halt. There was in fact some contraction, but Holst nevertheless succeeded in keeping the nucleus of the Laboratories intact during the years of crisis. This did a lot of good for morale in the Natuurkundig Laboratorium. During the occupation the Laboratories were a veritable oasis. The occupying authorities recognized that it would be useless to issue military assignments for the Laboratories, and, in so far as they believed in a German victory, they probably thought that the research work being done would yield fruit for them later. However this may be, the research programme carried on almost normally, and no reports of any significance were sent to Germany. Holst did not talk much to us, the younger members of his staff, about the political situation or about resistance, but he went on in his own way, not without risk to himself. For a time he

was even held prisoner. In 1944 the occupying authorities began to demand tangible results. From this as well as other points of view, the liberation in September 1944 came in the nick of time.

When Holst retired in 1946, the Laboratories had largely overcome the difficulties of the war years. There began a period of rapid expansion.

During the first ten years of his retirement Holst served on the Philips Supervisory Board, and he also served for ten years as a 'Curator' of the Delft University of Technology, seven of them as President. His relations with this university went back to an earlier date, however. He had a great influence in bringing about the introduction of the course in engineering physics leading to the degree of *Natuurkundig Ingenieur* and in 1933 an honorary doctorate in the engineering sciences was conferred upon him. As chairman of government committees he played an important role in the founding of the second University of Technology in the Netherlands, the Technische Hogeschool at Eindhoven. From 1930 to 1938 he was a Visiting Professor at Leiden.

I have come to the end of my survey of Holst's career. I hope I have shed enough light on its salient facets. Let me conclude with a few more personal remarks.

Holst was for me and for many others a great preceptor; he was not a great lecturer. As a visiting professor he often gave no more than one or two introductory lectures, and then left the rest to his staff.

He was a leader, but not a manager in the current meaning of the word; he was averse to administration and did not concern himself much with organization charts.

His ideas on industrial research and on higher education were not laid down in systematic treatises; only a few of his speeches, including his inaugural address, have been preserved.

He was not a fluent speaker, but in a discussion he could surprise you time and again with his shrewd, sometimes rather paradoxical remarks couched in aphorisms. 'I didn't like cast iron' is one I have already mentioned. 'Difficult exams make a people stupid' and 'Clever students can be allowed a long period of study; dim students should finish their studies quickly' sum up his very original views on higher education. Another of his remarks, which I found very useful, was: 'Don't believe that rogues exist with built-in rectifiers,' in other words, you can't expect someone who is dishonest to others to be always honest within the Company.

He was excellent at choosing his staff, without however going about it very systematically, being

convinced — remember the ‘third commandment’ I quoted — that people of very different natures can be equally valuable to a laboratory. He did, though, once confess to me with a sigh: ‘I don’t know why it is, but I always find it difficult to judge people with brown eyes.’

Perhaps his greatest gift was his ability to inspire his staff with enthusiasm for the things in which he himself believed. He seldom set exactly defined tasks, but precisely because he recommended something no more warmly than with ‘I should try that some time,’ or advised against it no more sharply than with ‘I wouldn’t do that if I were you,’ his words carried an uncommon force of conviction.

In human relations Holst was frank and straightforward, but also — a strange combination — rather

shy, shrinking away from intervention in personal matters.

Holst once said he was more interested in things than in people. I don’t know whether that was true. What I do know is that he made it possible for many people to do fruitful work under favourable conditions in his own Laboratories and also, thanks to the traditions he created, in laboratories elsewhere. He must have been aware that giving people meaningful work to do can also help them to overcome personal difficulties.

Gilles Holst, his gaze fixed on the future, a trifle impatient at times, about to utter one of his shrewd remarks . . . that is how we see him depicted in the painting by Van de Molengraft, and that is how he will continue to live in our memory.

Making the tracks on video tape visible with a magnetic fluid

A. M. A. Rijckaert

It has been known for more than fifty years that magnetic effects at the surface of materials can be made visible by means of a liquid containing a magnetic substance in suspension. A suspension of this kind is known as a 'magnetic colloid' or a 'magneto-fluid'; sometimes the term 'Bitter water' has been used, after F. Bitter, who published the method in 1931 [1].

Magnetic colloids have been in use for some time at Philips Research Laboratories, mainly for making the magnetic tracks visible on video tape. This is done by using an aqueous suspension of aggregates (clusters) of iron-oxide particles about $0.007 \mu\text{m}$ in diameter. The clusters must not be larger than about $0.1 \mu\text{m}$ [2] because of the very fine distribution of the magnetic patterns on the tape. The special feature of our method is the good use it makes of diffraction effects when the patterns are observed under an optical microscope. This provides a simple way of evaluating the operation of the servosystems that control the movements in the mechanical part of a video cassette recorder — in our case the Philips VR 2020 — during the writing process. The advantage here is that the effect of the servosystems for writing can be observed separately from that of the servosystems for reading.

The video signal written on the magnetic tape is approximately sinusoidal. After treatment of the tape with colloid a pattern appears on its surface — a 'Bitter pattern'. The pattern consists of two accumulations of clusters of iron-oxide particles in each sinusoidal period; this is because the forces acting on the particles are proportional to the square of the magnetic field-strength [3]. The period of the lattice structures that form the Bitter pattern is thus half that of the corresponding sinusoidal signal. Our experimental arrangement is shown in *fig. 1*. The magnetic tape treated with colloid is illuminated at a particular angle (through a bundle of glass fibres) and the patterns on the tape are observed through a microscope *M*. The light diffracted back towards the microscope is increased in intensity if the path difference s is equal to the wavelength of the incident light or to a multiple of it; see *fig. 2*. The dominant wavelength in the diffracted light thus depends on s and hence on the period of the observed lattice structure. Since the incident light is more or less white, the regions on the

tape that correspond to a different period appear differently coloured under the microscope.

The diagram of *fig. 3* shows the frequencies, dependent on the magnitude of the video signal, that occur in the frequency-modulated signal written on the tape. The information in the video signal relating to

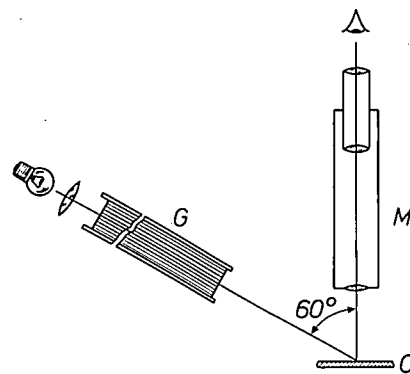


Fig. 1. Arrangement for observing Bitter patterns. *M* microscope. *O* object, a video tape treated with magnetic colloid (Bitter water). *G* bundle of glass fibres.

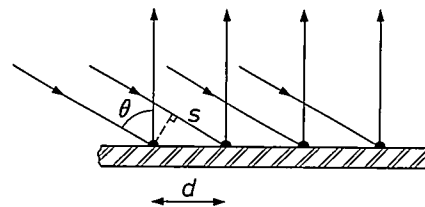


Fig. 2. Origin of interference due to the periodic structure of the Bitter patterns; s path difference after diffraction of the light rays; θ angle between the incident and the diffracted light, d 'lattice constant'.

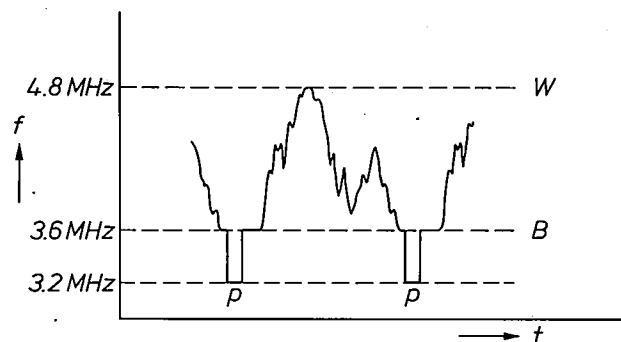


Fig. 3. The frequencies f of the frequency-modulated signal written on the tape and containing the video-signal information, as a function of time t . *W* and *B* correspond to the 'white' and 'black' levels of the video signal. The parts of the signal corresponding to the line-synchronization pulses are denoted by *p*.

[1] F. Bitter, On inhomogeneities in the magnetization of ferromagnetic materials, *Phys. Rev.* 38, 1903-1905, 1931.

[2] Magnetic liquids, *Philips tech. Rev.* 33, 293, 1973.

[3] N. H. Yeh, Ferrofluid Bitter patterns on tape, *IEEE Trans. MAG-16*, 979-981, 1980.

the successive lines of the television picture is divided up by line-synchronization pulses. (These pulses are used for synchronizing the flyback of the electron beam in the picture tube after a line has been written.) The lowest level in the signal corresponds to 3.2 MHz, the highest to 4.8 MHz. Since the magnetic heads in the VR 2020 recorder move at a velocity of 5.08 m/s, these frequencies correspond to signal wavelengths on the tape of 1.59 μm and 1.06 μm respectively. The 'lattice constant' d of the Bitter pattern thus varies from 0.80 μm to 0.53 μm .

Assuming that the light is incident on the Bitter pattern at an angle $\theta = 60^\circ$ and that the interference pattern of the first order is observed in the microscope, then the wavelength of the diffracted light is equal to

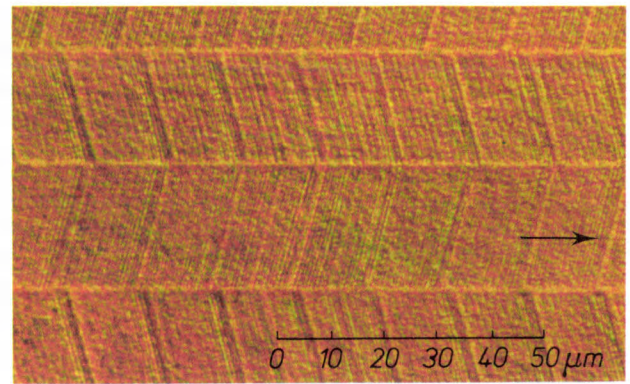


Fig. 4. The Bitter pattern, without interference. Two video tracks (with no space between them) are shown over their full width. The arrow indicates the direction of movement of the magnetic heads. The track width is approximately 22.5 μm .

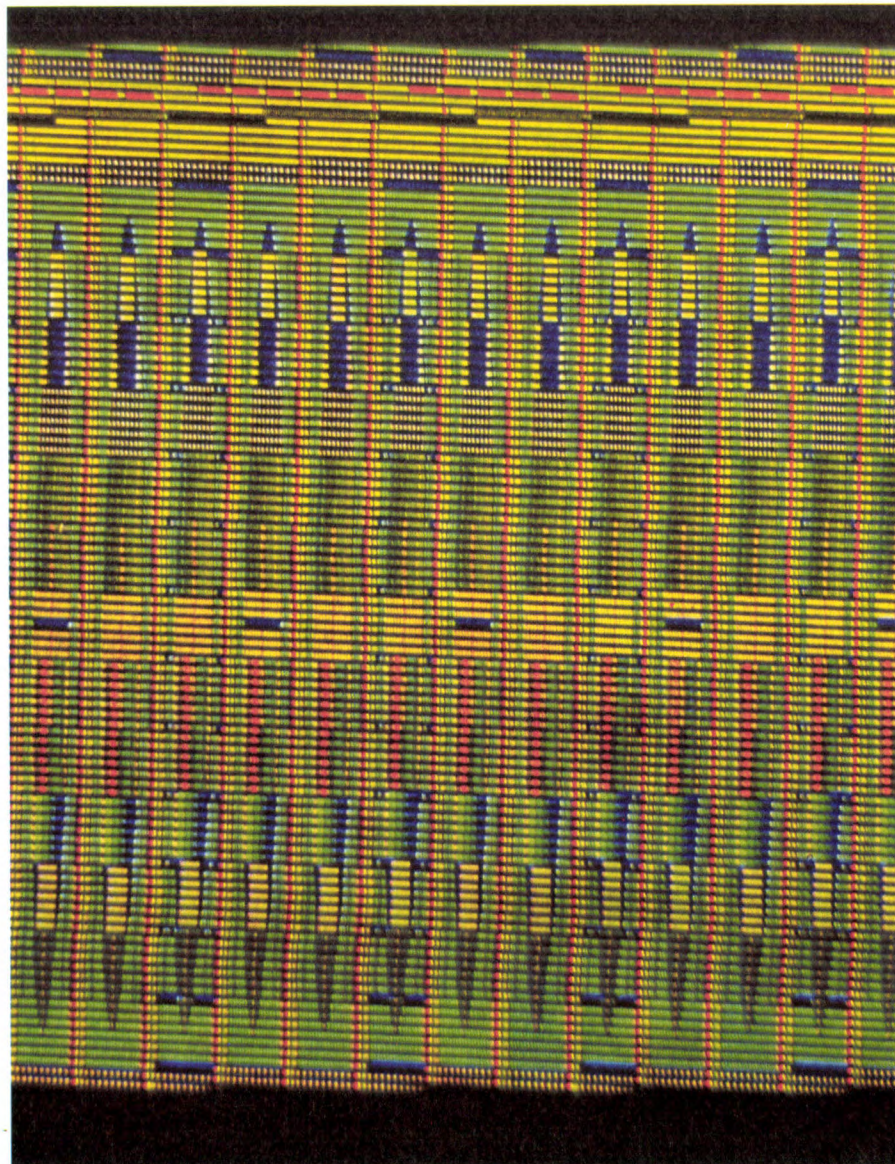


Fig. 5. The Bitter pattern with interference. The line length on the tape is 325 μm .

$\lambda = d \sin \theta = 0.87d$, as illustrated in fig. 2. The wavelengths of the light that makes the Bitter pattern visible in the microscope therefore vary from $0.69 \mu\text{m}$, for the lowest level of the line-synchronization pulse, to $0.46 \mu\text{m}$, for the highest level W . In the microscope image the synchronization pulses are thus observed as red areas, and the 'white' parts of the signal are observed as violet areas, and all the levels in between appear as other colours. The colours observed in the Bitter pattern are therefore completely unrelated to the colours in the corresponding television picture, but they are related to the value of the luminance.

The video signal in colour television contains the information about the colour value or chrominance and the brightness or luminance E_y . This quantity E_y is a 'weighted mean' of the local intensities E_R , E_G and E_B of the red, green and blue components of the television picture [4], as expressed in the well-known relation

$$E_y = 0.30 E_R + 0.59 E_G + 0.11 E_B.$$

The information about the individual quantities E_R , E_G and E_B is contained in the chrominance signal, whose carrier signal at 4.43 MHz is converted to 625 kHz in the VR2020 video recorder, so that the chrominance is not visible in the interference pattern.

Fig. 4 shows the Bitter pattern observed in the normal way, without interference, at the — rather high — magnification of $700\times$. Fig. 5 shows the Bitter pattern observed by means of interference and at a much lower magnification, $30\times$; here both the start and the end of the tracks written on the tape are visible. Fig. 6 shows a detail of this pattern, at a magnification of $50\times$. The video signal of these photographs relates to the test pattern transmitted by the Dutch television stations; see fig. 7. To clarify figs 4, 5 and 6, the method of writing the video signal on the magnetic tape in the VR2020 recorder will now be briefly described. The method follows the VIDEO2000 system.

The video signal is written on the tape in the form of oblique tracks by two magnetic heads rotating at high speed [5]. This is done by making the tape travel around a drum in a helical path through an angle of

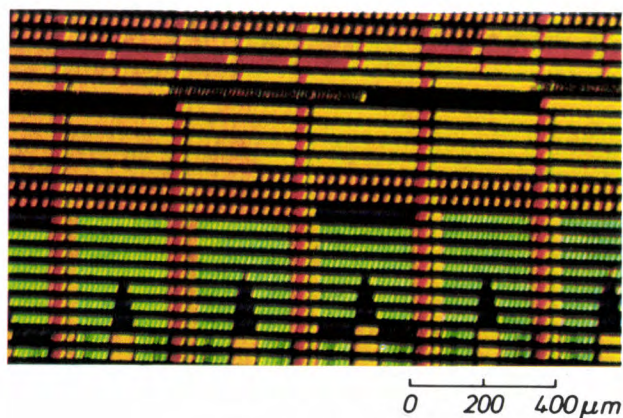


Fig. 6. Detail of fig. 5.

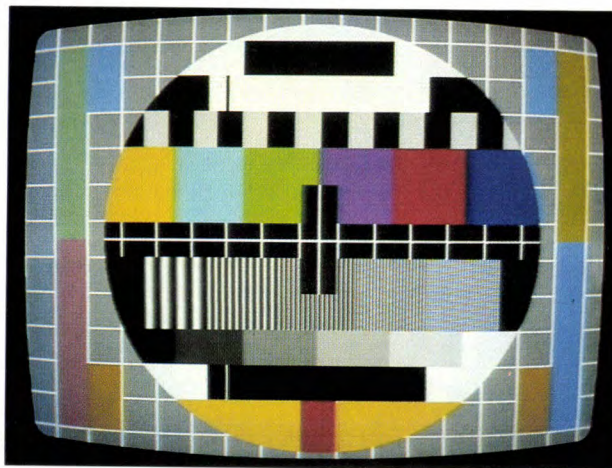


Fig. 7. The test pattern transmitted by the Dutch television stations. The photographs in figs 4, 5 and 6 give the video tracks corresponding to this pattern.

186° at 2.44 cm/s. The part of the drum carrying the two magnetic heads, mounted diametrically opposite, rotates at a circumferential velocity of 5.08 m/s. In the VR2020 recorder the two halves of the 12.7 mm (half-inch) wide tape are written one after the other (reversible cassette, maximum playing time 480 min). In fig. 5 half of the width of the tape can therefore be seen. To prevent luminance-signal crosstalk between two successive tracks, the write gaps of the magnetic heads are not located perpendicular to the direction of travel (the gaps are $0.5 \mu\text{m}$ long [6] and $24 \mu\text{m}$ wide). The gaps are mounted at an angle of 15° (the azimuth angle) to the normal to the direction of travel, with the two gaps rotated in opposite directions (see fig. 4). The tracks can therefore be written on the tape without an unused area between them and have a width of slightly less than $24 \mu\text{m}$. Fig. 4 also shows that the width of two successive tracks is practically identical, indicating that the head servosystem for writing operates correctly (this will be briefly discussed later).

Each oblique track contains the information relating to the field formed by the odd (or even) lines of a television picture. In each second the electron beam traces out 50 fields or 25 complete frames on the screen of the picture tube. The field consisting of the even lines starts with a half-line, the field with the odd lines ends with a half-line. In the VIDEO2000 system the tracks are written on the tape in such a way that the start of each track is always displaced by one and a half lines with respect to the previous one; see fig. 8.

[4] R. Theile, *Fernsehtechnik*; Springer, Berlin 1973.

[5] H. Bahr, *Alles über Video*; Philips Fachbücher, 1980.

[6] It is customary to call the dimension of the gap in the direction of the tape travel the 'length', even though the other dimension of the gap is much larger.

The result is that the line-synchronization pulses (seen from the writing head) are always side by side and an adjacent synchronization pulse cannot interfere with the picture information. In figs 5 and 6 the line-synchronization pulses are red areas, grouped along an approximately vertical line. The straightness of these lines shows that the other servosystems controlling the mechanical operation of the VR 2020 during the writing process operate satisfactorily.

In our experimental arrangement (see fig. 1) the incident light beam is given a direction perpendicular to the lines of the Bitter patterns relating to one magnetic head; see fig. 4. The light is reflected by the Bitter patterns from the other head in such a way that it does not enter the entrance diaphragm of the microscope objective. In figs 5 and 6 the Bitter patterns from the other head thus form black lines, so that these figures only show the tracks relating to one magnetic head.

The relation between the pictures in figs 5 and 6 and the corresponding test pattern in fig. 7 should now be clearer. Since identical video signals from the stationary test-pattern picture are repeated at a frequency of 25 Hz, the information relating to one test frame is also present in a direction perpendicular to the track, i.e. approximately in the vertical direction in the photographs. Since the adjacent tracks are displaced by one and a half lines, the Bitter patterns in this direction show (for example) the information relating to the 1st, 7th, 13th, 19th, . . . , 625th line; the 3rd, 9th, 15th, . . . , 621st line; the 5th, 11th, 17th, . . . , 623rd line, and so on; see fig. 8. The test pattern can thus be recognized in the vertical direction in fig. 5, although greatly distorted. As stated at the beginning, the colours do not correspond to those in the test pattern, but to the magnitude of the luminance signal.

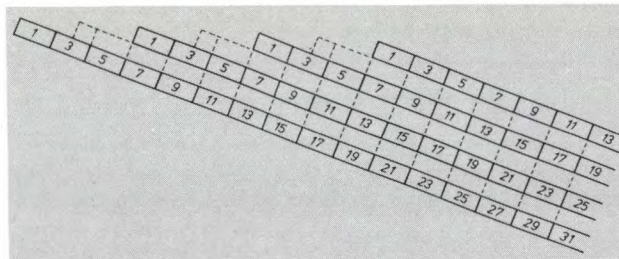


Fig. 8. The start of the oblique video tracks, with the numbers of the odd lines of successive television frames. The order of writing the tracks on the tape is from right to left in this figure (and also in figs 4, 5 and 6). The heads, however, move from left to right in relation to these figures.

In the upper part of fig. 5 patterns can be seen that consist of a series of five red areas corresponding to a total duration of $2\frac{1}{2}$ lines. Together these successive areas correspond to the field-synchronization pulse, which is used for synchronizing the flyback of the electron beam in the picture tube after writing a complete field. Below these pulses there are greenish areas with a duration of $1\frac{1}{2}$ lines, followed by a black area of the same length. The green areas correspond to a signal of 223 kHz, superimposed on a frequency of 3.6 MHz (the level *B* in fig. 3), which is used for the head servosystem for writing mentioned earlier. If no information is written on the tape (the black area), the magnetic head reads the crosstalk originating from the 223 kHz signal on the previous track, and stores its amplitude temporarily in a memory. When the next track is written the same procedure takes place, and then the two crosstalk signals are compared. The position of one of the magnetic heads is then corrected by the appropriate piezoelectric element in such a way that the 223 kHz crosstalk signals of the successive tracks are equal. As noted earlier, this produces tracks on the tape that are virtually constant in width.



Diamond die

The photograph shows a model of a single-crystal diamond die, about 4 cm in diameter. The real diamond, a natural product, is about 40 times smaller. Models like the one shown are used in the Philips diamond die factory at Valkenswaard (in the Netherlands) for demonstration and instructional purposes. The passage for drawing the wire can be seen. On both sides it has an accurately defined 'bell' shape, which determines the ratio of the diameter of the wire before passing through the die to the diameter after it emerges. In the situation shown the wire would pass through the die from top to bottom. The diamonds used for the dies have weights between 0.05 and

3 carats. The corresponding wire diameters after drawing are 8 μm to more than 2 mm. (For making wires of diameter greater than 0.12 mm synthetic diamonds are often used today, in a sintered and hence polycrystalline form.) The polished surface visible in the right foreground of the photograph is used as an inspection window during the production of the passage. In recent years laser drilling and microspark machining have been increasingly used in the manufacture of diamond dies.

We shall shortly publish an article on a number of special applications of spark machining, including the 'sparking' of diamond.

Manipulation of speech sounds

J. 't Hart, S. G. Nootboom, L. L. M. Vogten and L. F. Willems

With the rapid development of microelectronics it is now possible to buy machines that talk. The speech capabilities of such machines are in general still very limited, but it is reasonable to expect that 'synthetic speech' will come into much wider use within the next few years, e.g. in automatic telephone information services, home computers or reading machines for the blind. The development of good synthetic speech will depend on the manipulation of speech sounds by electronic methods. It is therefore necessary first of all to find out which of the many and varied properties of a sound wave are essential to the perception of speech. In the article below the authors discuss some aspects of their research in this field, and deal with some results that should make it possible to produce natural-sounding synthetic speech.

Introduction

Speech is the simplest — and often the best — means for human communication. For solving business and technical problems it is often a much more effective means of communication than writing or using a keyboard and display. And the telephone now offers oral communication over virtually any distance.

Since speech is such a good means of human communication, and since there is an increasing need for information exchange between man and machine, it is only natural to ask whether speech might not also be a suitable medium for such man-machine communication. The rapid development of microelectronics has now made this a real possibility. To answer the question it is necessary to investigate the physical characteristics of speech sounds, the possibilities of control of these characteristics, and their relation to the perception of speech. Research along these lines is currently taking place in many centres around the world. It also constitutes a major part of the programme of the Instituut voor Perceptie Onderzoek (IPO, The Institute for Perception Research) in Eindhoven.

Speech research at the Institute is carried out with the aid of a system (essentially a computer program) called SPARX, standing for SPeech Analysis and Resynthesis eXperiments. This system, using a computer and peripherals for recording and reproducing sound, is capable of analysing natural speech sounds into thirteen parameters that change relatively slowly, i.e. no faster than the speed at which the pharynx and mouth cavity change shape. From these parameters

the speech sounds can then be resynthesized. The parameters represent physical quantities that are directly responsible for distinct elements of speech perception, such as pitch, loudness and the various speech sounds. The coding into thirteen slowly varying parameters first of all allows the *memory capacity* required for storing the speech sounds to be *reduced* very substantially. Secondly, it makes the speech sounds capable of *manipulation*, since the parameters can be selectively processed before resynthesis, so that characteristics such as intonation, i.e. the rise and fall in the pitch of the voice in speech, can be altered without affecting the other characteristics. This ability to manipulate speech sounds is of great importance in speech-perception research and is particularly important in the development of 'speaking machines'. Their counterpart, the 'listening machine' (automatic speech recognition) lies outside the province of this article.

Any machine that plays back tape may be regarded as a speaking machine. Familiar examples are telephone-answering devices and recorded time and weather information by telephone. In this article, however, we think of a 'speaking machine' as equipment that is much more versatile, which can compose its own answers from a vocabulary of words spoken into it. Possible applications of such machines include reading aids for the blind, spoken instructions on how to use equipment, complex telephone systems that provide information automatically, and computers that give a spoken output. The 'vocabulary' can also be thought of as built up from smaller spoken-in units than words, e.g. syllables or phonemes.

J. 't Hart, Prof. Dr S. G. Nootboom, Ir L. L. M. Vogten and Ir L. F. Willems are with the Institute for Perception Research, Eindhoven.

A primary aim in the development of speaking machines is maximum economy in the storage of speech signals. With a system like SPARX a speech signal of 120 kbit per second is reduced via the thirteen parameters to 16 kbit/s. It is also possible, however, at the expense of the quality of the reproduced speech, to make do with fewer bits per second, by coarsening the parameter descriptions. The effect of such coarsening on the quality of speech reproduction can be studied with SPARX.

In the second place the aim is to generate fluent and readily intelligible speech utterances with the small units of the vocabulary. If the machine reproduces the right words in the right sequence but does *nothing more*, the result is useless, since in natural speech the sound of a word depends closely on its context. To obtain fluent, intelligible speech the physical characteristics of each word must be matched to its context. This is a point where there is a particular need for research on speech perception and for the possibility of manipulating speech sounds.

The present article deals with research of this nature with the SPARX system. In the first section we recapitulate the physical structure of natural speech, showing how it comes about and presenting a generally accepted model of speech production, which is also at the base of SPARX. Next we deal with the system itself, discussing some of the possibilities it offers for 'playing with speech sounds'. Finally we discuss the rules for generating sentence intonation, as an example of what is needed to synthesize speech utterances that are reasonably intelligible.

Natural speech

Fig. 1 shows a cross-section of the organs of speech. The physical principle of speech production is simple. There is a variable sound source and there is a variable acoustic filter that alters the sound from the source.

For the vowels and the voiced consonants (e.g. m, n, l, r, b, d) the source sound is produced by vibration of the vocal cords (2), which transmit air-pressure pulses from the trachea (1) at a particular frequency (the 'source frequency'). The source frequency determines the perceived pitch of the sound. A speaker regulates the source frequency, and hence the pitch of his speech, by varying the tension in his vocal cords. The average pitch varies from one speaker to another, and is generally higher for women and children than for men. The energy content of the pressure pulses, and hence the loudness, is determined by the pressure drop across the glottis and the tension in the vocal cords. The sound volume varies very considerably: speech has a wide dynamic range (fig. 2).

The acoustic filter for the voiced sounds consists of the mouth cavity and the pharynx, and when the soft palate (4) does not shut off the nasal cavity — e.g. for the utterance of the nasal consonants m and n — it also includes the nasal cavity. While speaking the speaker continuously changes the shape of the pharynx and mouth cavity by movements of the tongue, the lower jaw and the lips, so that the filter action and hence the physical structure of speech sounds are continually changing. This is the main cause of the differences between the individual vowels and consonants.

For the unvoiced fricatives f, s and sh the source sound is noise produced by turbulence of the air stream from the lungs in a narrowing of the mouth cavity. For f this narrowing takes place between lower lip and upper teeth, for s between the tip of the tongue and the teeth-ridge and for sh between the tongue and the hard palate. The sounds of the voiced fricatives v and z have two sources: vibrations of the vocal cords and air turbulences. For the unvoiced plosives p, t and k the mechanism is again somewhat different. To produce these sounds the passage of air through the mouth cavity is stopped for a short time (about 50 to

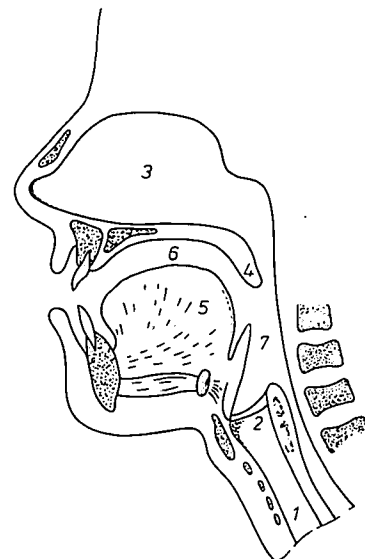


Fig. 1. Cross-section of the organs of speech. 1 Trachea. 2 Vocal cords. 3 Nasal cavity. 4 Soft palate. 5 Tongue. 6 Mouth cavity. 7 Pharynx.

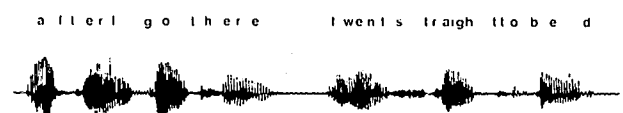


Fig. 2. Recording of the sound wave of the spoken sentence: 'After I got here, I went straight to bed'. The diagram gives the output voltage from the microphone (vertical axis) as a function of time (horizontal axis).

150 ms), and when the air pressure built up during this time is suddenly released by opening the closure, a short burst of sound is produced. The closure for p takes place between the lips, for t between the tip of the tongue and the teeth-ridge and for k between the back of the tongue and the soft palate. For the voiced plosives b, d and g there is again a sudden opening of a closure of the mouth cavity. These differ from the unvoiced plosives, however, in that the vocal cords vibrate during the closure.

For these fricative and plosive sounds the acoustic filter is confined to the space between the narrowing that produces the turbulences and the mouth opening. While the acoustic filter for g, for example, is formed by a large part of the mouth cavity, only a small part is used for producing s, while the source sound of f leaves the mouth practically unfiltered.

The source-filter theory of speech production

A generally accepted theory of speech production is G. Fant's 'source-filter theory' [1]; see fig. 3a. In this model the sound originating from the source *U* as a series of pulses, or as acoustic noise, passes through two filters, consisting of the filter *O* formed by the pharynx, the mouth cavity and the nasal cavity, and the filter *R* that represents the sound radiation at the mouth opening. An essential feature of voiced sound is that the spectrum contains a large number of overtones in addition to the fundamental tone (of frequency *F*₀); it is only through these overtones that the filter *O* can produce such a marked effect. In the spectrum of a normal voiced source the amplitude of the overtones decreases by about 12 dB per octave. If the source sound is acoustic noise, it is usually assumed that its spectrum is practically flat ('white noise').

The amplitude spectrum *S(f)* of a steady-state speech signal is now equal to the spectrum *U(f)* of the source, multiplied by the transfer functions *O(f)* and *R(f)* of the filters *O* and *R*:

$$S(f) = U(f) \cdot O(f) \cdot R(f). \quad (1)$$

The transfer function *R(f)* represents the increasing directivity of the mouth opening to sound as the frequency increases and is equivalent to an increase of the on-axis sound radiation by 6 dB per octave.

The transfer function *O(f)* mainly determines the 'nature' of the speech sound. The mouth cavity together with the pharynx may be regarded as a somewhat irregularly shaped tube, which is almost closed at one end and open at the other. Such a tube has a number of resonant frequencies that correspond to peaks, known as formants, in the transfer function *O(f)* (fig. 3b). Each formant is characterized by a centre frequency and a bandwidth. For the perception

of speech no more than five formants are generally required, in the frequency range from 100 Hz to 5 kHz. In order of increasing frequency they are denoted by *F*₁ to *F*₅. The intrinsic feature of vowels that distinguishes them from each other appears to depend primarily on the first three formants. The formants *F*₄ and *F*₅ do not essentially contribute to the recognizability of speech sounds, but they do largely determine the naturalness of the speech and the recognizability of the speaker.

In speech research as carried out with systems like SPARX a simpler model is used, the 'synthesis model' [2], as illustrated in fig. 4. It contains only one filter *O*, which combines the function of the filter *O* in fig. 3a with that of *R* (6 dB increase per octave), and also includes the 12 dB decrease per octave of the voiced source. The voiced source *V* therefore has a 'flat' spectrum here, consisting of a series of frequency

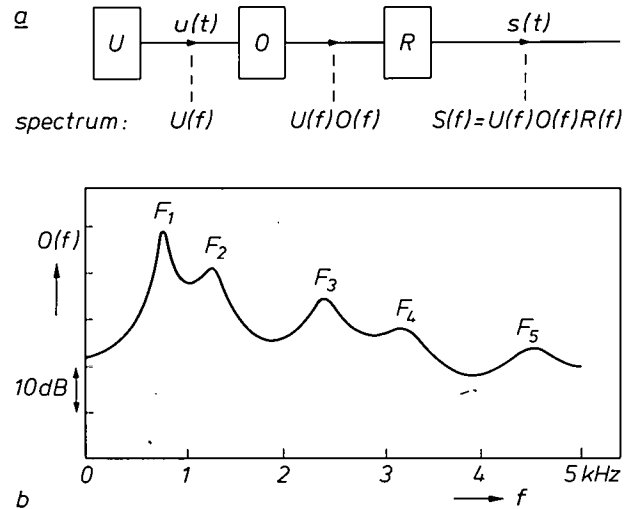


Fig. 3. a) Block diagram of G. Fant's source-filter model [1] for the speech organs. The source sound *u(t)* from the source *U* (lungs and vocal cords) passes through the filter *O* (pharynx and mouth cavity) and the filter *R* (mouth opening) to become the speech sound *s(t)*. The signals in the 'frequency domain' are indicated beneath the diagram. b) Transfer function *O(f)* of the filter *O* for a particular shape of the pharynx-mouth channel; the peaks are called formants.

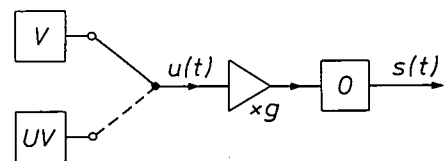


Fig. 4. The synthesis model. The filter *O* contains the functions of *O* and *R* from fig. 3 together with the spectral structure of the source (decrease of 12 dB per octave). The spectrum of the source is therefore 'flat' here: it is either a series of equidistant frequency components of equal amplitude (*V*, voiced) or white noise (*UV*, unvoiced). The amplitude (*g*) of the sound is controlled by a variable amplifier between source and filter.

components of equal amplitude at the frequencies F_0 , $2F_0$, $3F_0$ and so on (fig. 5b). The corresponding source signal consists of a series of delta pulses with a period of $1/F_0$ (fig. 5a). The unvoiced source UV is again white noise. The spectrum $S(f)$ of the 'synthesized' speech signal $s(t)$ (see fig. 5c) is given by:

$$S(f) = U(f) \cdot O(f). \quad (2)$$

The envelope of this spectrum is not flat, and consequently the discrete pulses of fig. 5a are spread out in time (fig. 5d). The level of the speech signal can be adjusted by an amplifier between source and filter.

In this model speech is described by the following parameters: the binary parameter (V/UV) that determines whether the source is voiced or unvoiced; the source frequency F_0 — for the case of a voiced source; the amplitude g ; and the frequencies F_1 to F_5 and the bandwidths B_1 to B_5 of the formants that characterize O . Anticipating the operation of SPARX, fig. 6 gives an example of a speech utterance analysed into these thirteen parameters. Since these parameters represent distinct elements of the speech perception, they can be used as the starting points for reproducing and manipulating speech sounds. Neither this model nor

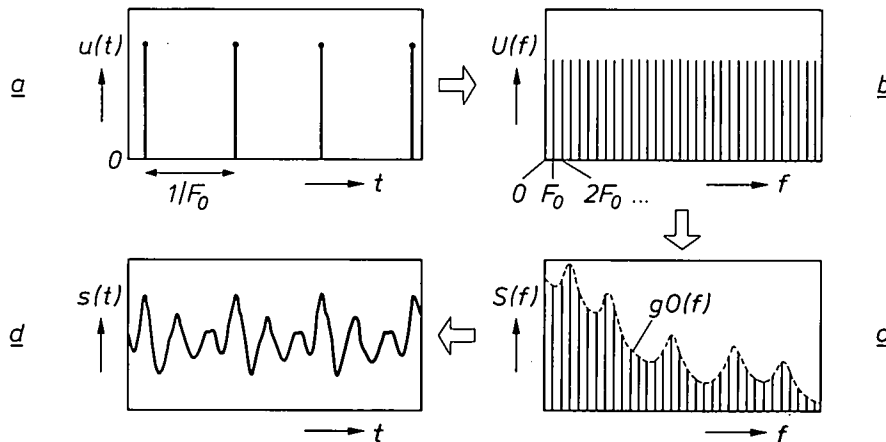


Fig. 5. Signals and spectra in the synthesis model with voiced source. a) The source signal $u(t)$, a series of delta pulses with period $1/F_0$. b) The spectrum $U(f)$ of $u(t)$, a series of components of equal amplitude at a spacing F_0 (fundamental with overtones). c) The spectrum $S(f)$ of the speech signal, which is the product of $U(f)$, the amplitude gain g and the transfer function $O(f)$ of the filter O in fig. 4. d) The resulting speech signal $s(t)$.

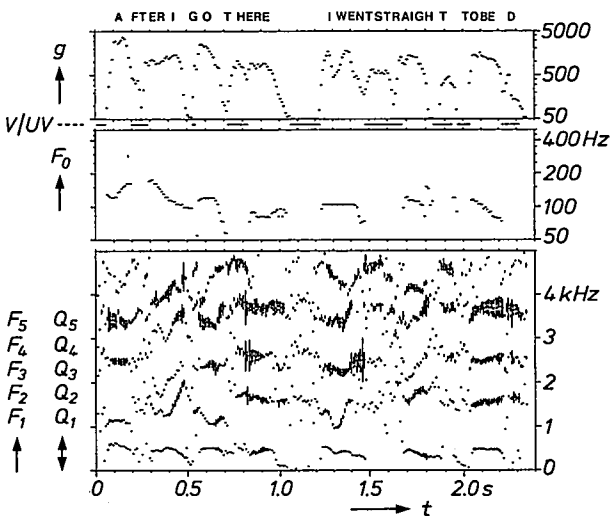


Fig. 6. Analysis of the speech utterance: 'After I got here, I went straight to bed' into the thirteen parameters, as functions of time. Logarithmic scales are used for the amplitude g and the source frequency F_0 . For V/UV , voiced (V) is indicated by white, and unvoiced (UV) by black. In the lower part, the formant frequencies and bandwidths are marked every 10 ms by five vertical bars. The centres of the bars represent F_1 to F_5 , the lengths give the quality factors Q_1 to Q_5 ; the bandwidths are derived from these by the relation $B = F/Q$. A long bar thus corresponds to a higher peak in the spectrum. The parameters are stored at a rate of 16 kbit/s.

SPARX take account of the possibility that the source sound may simultaneously be periodic and noisy, like the sounds 'v' and 'z'.

The SPARX system

Analysis of a natural speech signal $s(t)$ (fig. 7a) into the thirteen parameters is really a matter of resolving the spectrum $S(f)$ of the signal (fig. 7b) into the two factors of the model (fig. 7c,d): a relatively smooth transfer function $O(f)$ of the filter O and the spectrum of the source $U(f)$, which is either a series of equidistant frequency components of equal amplitude or a white-noise spectrum. The analysis, and also the resynthesis, are not carried out in the 'frequency domain' (Fourier analysis and Fourier synthesis), however, but entirely in the 'time domain' (operations on pulse series); this is done with a digital computer. We shall now discuss this method.

[1] G. Fant, Acoustic theory of speech production, Mouton, The Hague 1970.
 [2] See for example J. D. Markel and A. H. Gray Jr., Linear prediction of speech, Springer, Berlin 1976.

In the analysis the analog speech signal first has to be digitized. It is therefore sampled at a frequency (in SPARX) of 10 kHz, i.e. at a sampling rate that makes proper allowance for frequencies in the spectrum up to nearly 5 kHz (Nyquist's theorem^[3]). Each sample is then quantized to an integer value between -2048 and +2047 (12 bits). One second of speech has thus now been coded into 10 000 × 12 bits, i.e. into 120 kbit. For the computer calculations we can now treat the signals as time series of signal values (samples).

The *filter* and *source characteristics* are extracted from the speech signal independently of each other (fig. 8). We shall first consider the extraction of the filter characteristics ('formant extraction'). This is done in an analysis window of 250 samples (25 ms of speech), which shifts 100 samples (10 ms) at a time. The overlap between the successive windows is thus 15 ms. The window is made large enough for it to always contain more than one period of the source sound, yet kept small enough not to smooth out too much of the variation of the speech parameters with

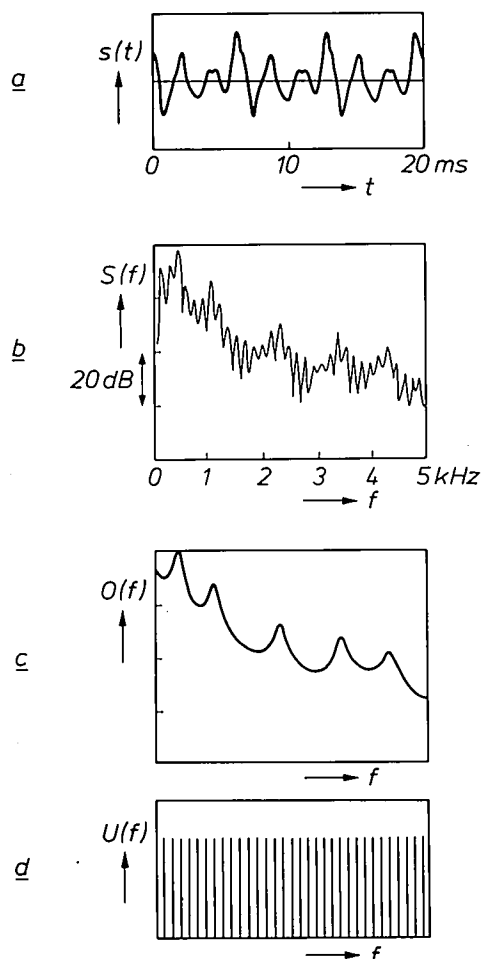


Fig. 7. The analysis problem 'in the frequency domain' is to analyse the spectrum (b) of a natural (steady-state) speech signal (a) into the factors of the synthesis model (c,d). The vertical scale is linear in (a) and logarithmic in (b) and (c). The analysis is in fact performed with a digital computer in the time domain.

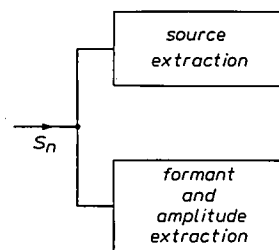


Fig. 8. Diagram of the analysis. The extraction of the source characteristics ($V/UV, F_0$) and the extraction of the formants (F_1, \dots, B_5) and the amplitude (g) from the time sequence of speech samples s_n are carried out 'in parallel', independently of each other.

time. Before the analysis of the signal segment in a window, it is multiplied by a 'Hamming window'^[3] to avoid adverse transient effects due to the abrupt start and finish. The changes in the physical characteristics of the speech sounds can readily be followed with the analysis period of 10 ms of speech.

The *formant extraction* — without knowledge of the source — is performed by *linear prediction* of the speech signal. This is done using the model of the filter O in fig. 9. The filter action is obtained by feedback of the output signal s_n to the input via the 'predictor' P (a transversal filter). In fig. 9 the indication z^{-1} represents the operator that delays the signal by one sampling period. The signal \hat{s}_n that appears at the output of P is therefore a linear combination of $s_{n-1}, s_{n-2}, \dots, s_{n-10}$:

$$\hat{s}_n = \sum_{j=1}^{10} a_j s_{n-j}. \tag{3}$$

For the output signal in fig. 9 we now have:

$$s_n = g u_n + \hat{s}_n. \tag{4}$$

This *linear* relation enables us to *predict* an output signal value from the instantaneous input value and the ten preceding output values; hence the name 'linear prediction'. Now u_n is only seldom $\neq 0$ (the sampling frequency is very much higher than the source frequency), so that 'nearly all' speech samples are predictable just from the ten preceding samples^[4]. This feature is utilized for making the 'best possible' determination of the coefficients a_1, \dots, a_{10} of P , without knowing anything about the source. The method of least squares is used to calculate the values of a_1 to a_{10} that give the best match of \hat{s}_n to the actual signal value s_n (see fig. 10), for all sample predictions \hat{s}_n in the window. Since \hat{s}_{11} is the first prediction, we look for the values for which the expression

$$E = \sum_{n=11}^{250} (s_n - \hat{s}_n)^2 \tag{5}$$

has a minimum.

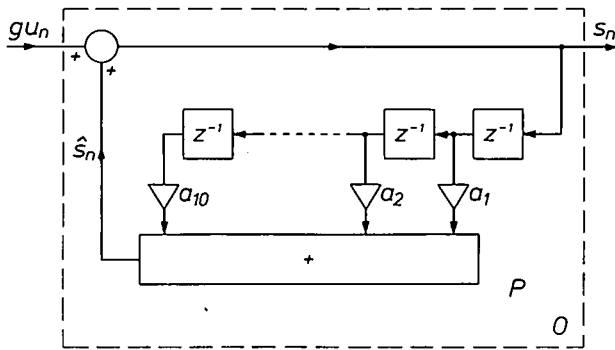


Fig. 9. Representation of the filter O as a 10th-order digital filter. The filtering action is produced by feedback of the output signal s_n to the input via the 'predictor' P . The predictor delays the signal in ten steps, multiplies it after each step by a coefficient (a_j), and adds the results. Expressed mathematically $P(z) = \sum_{j=1}^{10} a_j z^{-j}$, where z^{-j} is the operator that delays the signal by j sampling periods. At every instant the filter O is thus characterized by the ten coefficients a_1, \dots, a_{10} .

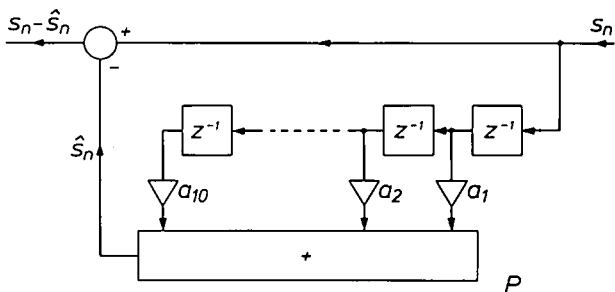


Fig. 10. Determination of the coefficients of the predictor P for each analysis window of 25 ms from natural speech. The (sampled) speech signal s_n is applied to the predictor and the signal \hat{s}_n predicted by P is subtracted from it. The desired values of the coefficients are the values for which $(s_n - \hat{s}_n)^2$, summed over the entire analysis window, has a minimum.

The values found for a_1 to a_{10} are then converted into formant frequencies and formant bandwidths. The conversion is based on the model in *fig. 11* for the filter O , which is equivalent to five resonators in cascade, each characterized by a frequency and a bandwidth. This conversion completes the formant extraction.

The procedure adopted also provides a good measure for the amplitude g of the speech signal: the square root of the minimum E_{\min} found for E (see eq. 5). In resynthesis this gives good results. This is to be expected, since eq. (4) indicates that the mean of $(s_n - \hat{s}_n)^2$ taken over a segment can never be zero, but must be equal to the mean of $(gu_n)^2$.

The method outlined here for determining the coefficients a_1 to a_{10} of the filter O is called 'Linear Predictive Coding' (LPC)^[4]; it is also referred to as 'Wiener Filtering'^[6] or 'Inverse Filtering'^[7]. The method is heuristic: it does not follow strictly logically from the problem posed, but can ultimately be justified because it leads to manageable equations for a_1 to a_{10} and, as appears from the resynthesis, it gives good analysis results^[2].

The model in *fig. 9* has its limitations. It is an 'all-pole filter', which implies that it can represent 'resonances' but not 'antiresonances'. Since the nasal cavity operates as an 'antiresonator' during nasal sounds, the model does not do full justice to such sounds.

For the *extraction of the source characteristics* (see *fig. 8*) the first procedure is to count the number of zero crossings of the speech signal in unit time. If this number exceeds a critical threshold, the source is classified as unvoiced (UV), in the other case as voiced (V). This extremely simple procedure serves well in practice.

For determining the source frequency of voiced sound the analysis window is increased to 40 ms, to ensure that at least two periods fall within the window even at the lowest source frequencies; this is necessary for the recognition of the periodic structure. The determination is made by using a modification of the autocorrelation method, which is based on the strong correlation between samples spaced by one period (if there are any periods at all)^[8]. *Fig. 12* gives the auto-

correlation function $R(k) (= \sum_{n=1}^{N-k} s_n s_{n+k})$ as an example for a speech segment of 40 ms ($N = 400$). The time shift τ of the largest peak away from the origin gives the pitch period.

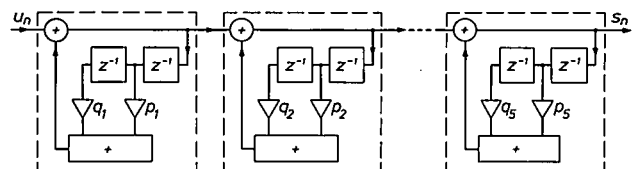


Fig. 11. Representation of the filter O as five second-order digital filters in cascade. If the coefficients p_1, \dots, q_5 satisfy the relation

$$1 + \sum_{i=1}^{10} a_i z^{-i} = \prod_{k=1}^5 (1 + p_k z^{-1} + q_k z^{-2})$$

the filter is equivalent to the 10th-order filter in *fig. 9*. The a 's are converted into the p 's and q 's using the Bairstow routine^[5]. Each of the second-order filters is equivalent to a resonator, and the filter O is thus equivalent to five resonators in cascade. The resonant frequency F and the bandwidth B of a resonator are connected with the coefficients p and q by the relations

$$p = -2 \exp(-\pi BT) \cos 2\pi FT, \\ q = \exp(-2\pi BT),$$

where T is the sampling period.

[3] See for example L. R. Rabiner and R.W. Schafer, *Digital processing of speech signals*, Prentice-Hall, Englewood Cliffs 1978.

The Nyquist theorem is also discussed in: F. W. de Vrijer, *Philips tech. Rev.* 36, 305, 1976, on page 343.

[4] See for example B. S. Atal and S. L. Hanauer, *J. Acoust. Soc. Amer.* 50, 637, 1971.

[5] See for example C.-E. Fröberg, *Introduction to numerical analysis*, 2nd edition, Addison-Wesley, Reading, Mass., 1969.

[6] N. Levinson, *J. Math. and Phys.* 25, 261, 1947.

[7] J. D. Markel, *IEEE Trans. AU-20*, 129, 1972.

[8] L. R. Rabiner, *IEEE Trans. ASSP-25*, 24, 1977.

Determining an autocorrelation function is very time-consuming, because of the many multiplications. For this reason SPARX does not use the autocorrelation function but a modification of it, in which s_n is not multiplied by the displaced sampling value s_{n+k} itself but by the sign of s_{n+k} . Only additions and subtractions are then required.

To save more computer time, the search for the peak of the autocorrelation function in each analysis window is limited to a small range on the τ -scale (see fig. 12) concentrated around the peak found in the preceding correlation diagram. This method is based on the knowledge that the source frequency in speech sound varies rather slowly. It guarantees a certain continuity in the measured F_0 -value, even through the analysis windows in which the periodicity of the signal is not very clear, e.g. because of low intensity. The disadvantage of the method is that once a serious error has been made, for example an octave jump, the error remains. For the first analysis window of a series of 'voiced windows' the area on the τ -scale must be made wide enough for it to include the possible source periods of the speaker.

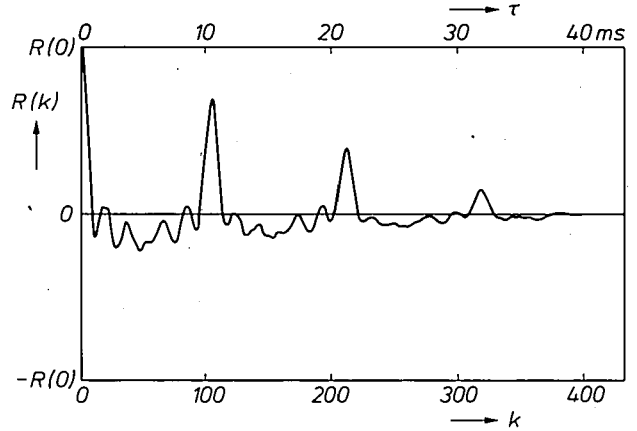


Fig. 12. The autocorrelation function $R(k)$ of a speech segment 40 ms long. This function gives the correlation between samples separated by k sampling periods (τ seconds, where $\tau = kT$, and T is the sampling period. If the signal has periodicity, the periods appear as maxima in the autocorrelation diagram.

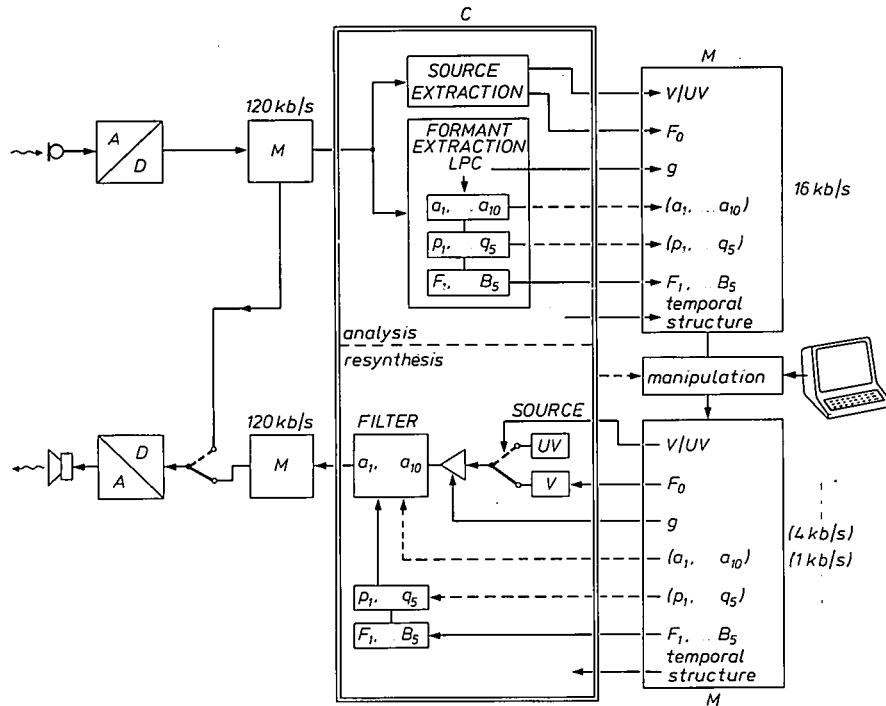


Fig. 13. Block diagram of SPARX. The analysis takes place in the upper half, the resynthesis in the lower half. A/D analog-to-digital conversion, D/A digital-to-analog conversion. M disk memory; this is used at four places in the diagram: for storing the speech signal in digital form both before the analysis and after the resynthesis, and for storing the parameters obtained from the analysis and as used for the resynthesis. A 'manipulation' can be made between both sets of parameters (coarsening of the description or a change in the variation of some of the parameters). C computer, in which the actual analysis and resynthesis take place. The blocks in C only represent software, of course, and not hardware. The manipulation can be made either by the intervention of the experimenter or by means of a program in the computer. In the resynthesis the parameters a_1, \dots, a_{10} or p_1, \dots, q_5 obtained in the formant extraction can be used instead of F_1, \dots, B_5 .

Fig. 13 gives a block diagram of the entire SPARX system. The upper part represents the analysis of a speech signal, as described above. To summarize, the spoken analog signal is converted to 120 kbit/s, and stored in the disk memory. It is then available for analysis into the parameters

$V/UV, \dots, B_5$. These are stored in M for each 10 ms of speech.

The *resynthesis* (lower part in fig. 13) takes place as indicated in the diagram given in fig. 4, and essentially amounts to a reversal of the analysis. The parameters from the memory are used to control a 'source' and a

'filter' (not hardware, but software). They are updated after each 10 ms of speech; the output signal is again stored in M , as samples at a rate of 120 kbit/s, and after conversion it can be made audible through a loudspeaker.

Speech research with SPARX

With SPARX the set of parameters obtained by analysis can be modified before the resynthesis, and the effect of the modification on the speech sounds can be studied. We shall consider here two kinds of modification: coarsening the parameters ('bit-rate reduction') and altering the behaviour of the parameters ('playing with sound').

To save memory space in the computer a coarser form of storage can be used. The result of such a coarsening can be monitored with SPARX. The coarsening takes place in two ways. In the first place the model parameters are quantized more coarsely, and in the second place the number of analysis steps per second is reduced.

This reduction is sometimes permissible because the changes in the speech sounds usually take place fairly slowly. Rapid changes in speech effects, such as the sudden increase of energy in plosives (p, t, k, b, d) are not therefore reproduced so well when the number of analysis steps is reduced.

The permissible quantization is not the same for all parameters. The human ear is fairly sensitive to

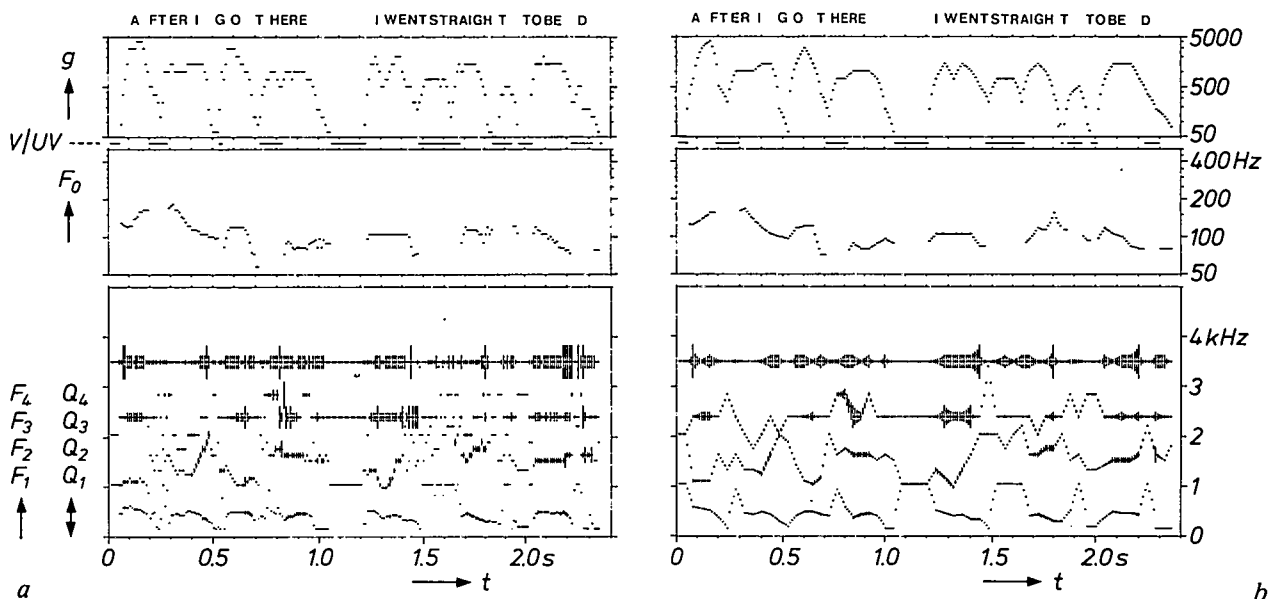


Fig. 14. Analysis data from the same speech utterance as in fig. 6 but with a reduced number of bits for the parameters: a) 4 kbit/s, b) 1 kbit/s. So as to reduce the bit rate, the number of formants is limited to 4, and only the quality factor Q_4 of the fourth formant can be varied; F_4 is fixed. In (b) the parameters are only given for every 40 milliseconds; in between there is automatic linear interpolation.

Bit-rate reduction

As we said in the introduction, it is very important in various applications of speaking machines to be able to store speech utterances or speech signals as economically as possible.

After analysis a speech signal in SPARX is stored at a rate of 16 kbit/s (each second of speech contains 100 analysis frames, each frame produces 13 parameters, and each parameter is on average described by some 12 bits). This is a substantial reduction compared with the original 120 kbit/s (see page 138). In direct resynthesis the process does give some loss of quality, but the result is highly intelligible and the sound quality is good.

changes in the frequencies of the lowest three formants, but is much less sensitive to changes in the frequencies of the fourth and fifth formants. These can therefore be quantized very roughly or even taken as fixed. Nor do the bandwidths of the formants have to be very exactly preserved.

If careful attention is paid to the quantization, speech is still readily understandable after a reduction to 1 kbit/s. Below this level the intelligibility of speech rapidly decreases. In practical applications 4 kbit/s seems to be a good compromise. Fig. 14 gives the parameters for the same speech utterance as in fig. 6, but now described at a rate of 4 kbit/s and 1 kbit/s, respectively.

Playing with speech sounds

SPARX allows the experimenter to 'play' with the variations in pitch, in amplitude or in one or more of the other parameters, without affecting the others, and to study the effect of such changes on the speech sounds. In addition to the thirteen parameters there is a 'hidden' parameter, which has already been included in fig. 13: the time structure. Fig. 15 gives an example in which both the pitch and the time structure are manipulated, in such a way that the utterance not only sounds different, but has changed its meaning.

and sounds reasonably natural. The chance of finding such rules is based on the possibility of replacing the capricious variation of the parameters of natural speech by a strongly simplified, 'stylized' variation, with not too great a degradation of the intelligibility and sound quality.

Rules of adaptation will have to be made for the duration, the sound volume, the source frequency and the formant frequencies of the word sounds in the vocabulary. Our research on this topic at our Institute, particularly on the source frequency (the pitch),

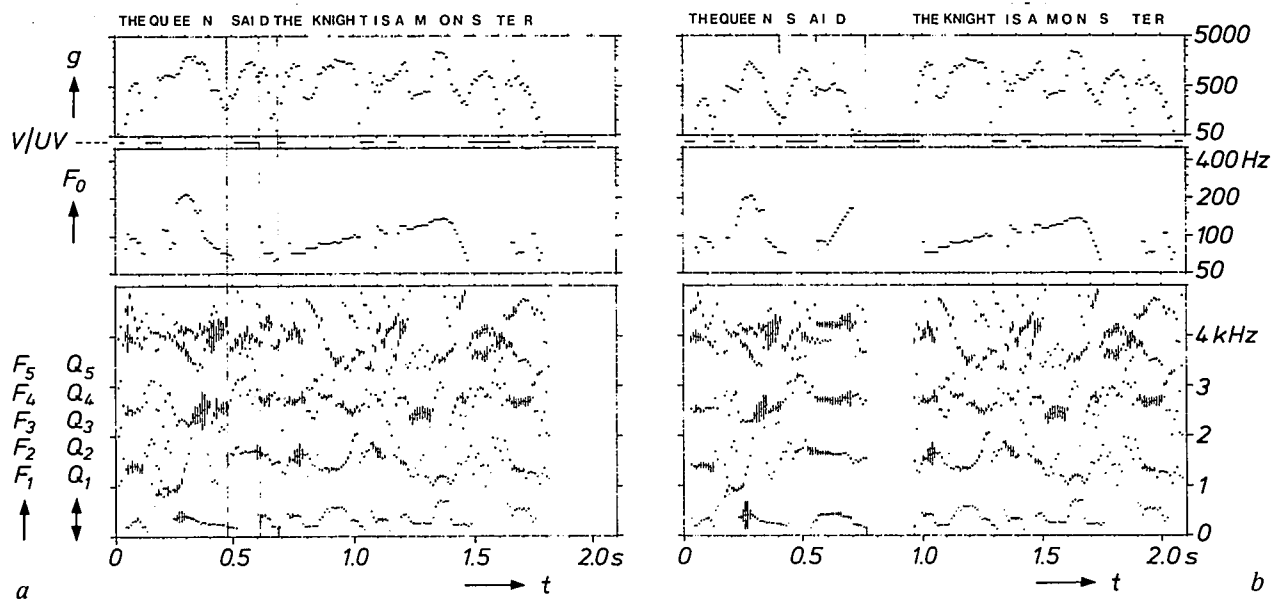


Fig. 15. *a*) Data from an analysis of the speech utterance: "The Queen", said the Knight, "is a monster". *b*) The same data after the following manipulations: change of the pitch contour of 'The Queen said' by a factor of 0.8; lengthening of the sound 'ai' in 'said' by a factor of 2; shortening of 'The Queen' by a factor of 0.8; insertion of a 200ms pause after 'said'. These manipulations change the meaning to: 'The Queen said "The Knight is a monster"'.

The ability to manipulate speech sounds is essential to the development of speaking machines. To make this clear in the present context, we shall consider a 'text-to-speech' system that converts text (presented in digital form) into speech sound. The system recognizes the words of the text and calls the appropriate speech sounds, coded in the thirteen parameters, from its vocabulary^[9]. As we have seen (page 135) the word sounds called must then be matched to their sound context in the speech utterances. Many workers in many countries are trying to find appropriate rules for this matching process. The rules must be simple enough for them to be translated into algorithms for the machine, and at the same time they must produce adaptations that lead to speech that is understandable

has resulted in a set of useful rules that we call an 'intonation grammar'. Investigations into rules for the other parameters are in full progress.

To conclude this article we shall discuss the intonation grammar we have devised for British English. We shall consider only a very simplified version, since we are only concerned here with illustrating our research on this subject. We have deduced this grammar from a large number of stylized pitch contours of natural speech. The stylization can be quite drastic without noticeable degradation of the resultant speech. Fig. 16 gives an example of a stylized pitch contour that trained listeners found indistinguishable from the

^[9] Systems of this kind already exist for American English; see for example D. H. Klatt, IEEE Trans. ASSP-24, 391, 1976.

original contour. We tested the usefulness of our intonation grammar in a large number of listening tests, in which the original intonation of natural speech utterances was replaced via SPARX by intonation constructed from the grammar. This showed that the great majority of pitch contours can be simulated very satisfactorily with the simple version.

An intonation grammar for British English (simple version)

Our grammar (in its simplest version) indicates that the character of the pitch variation is always as shown in *fig. 17*: the pitch goes up and down between three slowly falling lines, the low, middle and high 'declination line' (1, 2 and 3). A sentence always begins at the middle declination line and always ends at the lowest one. The pitch variation of a sentence always forms a *coherent* entity, which runs 'virtually' over the unvoiced ('pitchless') parts. The distance between the lines 1 and 2 and between 2 and 3 amounts to six semitones. For long sentences ($t > 5$ s) the difference between the initial and final pitch of (say) the lowest declination line is constant, but for short sentences ($t < 5$ s) the difference decreases with the length of the sentence; see *fig. 18*. A quantitative formulation of this rule is given in the caption. During silent intervals longer than 250 ms the declination is stopped.

Between the declination lines there are four kinds of transition: an 'accentuating' rise (4) from 2 to 3, two kinds of 'non-accentuating' rise (5) and (6) from 1 to 2, and an 'accentuating' fall (7) from 3 to 1. The 'small' rises 4, 5 and 6 last for 80 ms, the 'large' fall 7 lasts for 160 ms. Two things are now necessary for the positioning of the transitions: a) *the presented text* must be provided with marks indicating intervals and stressed words; b) *the coded speech sound of each of the words in the memory* must have marks relating to the stressed and the last syllable.

We look first at the markings in the presented text. These consist of *brackets*, which mark the beginning or the end of a sentence, *strokes* marking other important grammatical boundaries and *underlinings* to indicate the words that have to be stressed. The following is given as an example:

[no cricketer / has ever been accused / of taking drugs before the match / and no-one throws bottles at the players / as in football]

The brackets correspond to the full stops in the text. Where there are commas there should certainly be strokes, but in general there will be far more strokes than commas (see example). The strokes, like the underlinings, will therefore either have to be separ-

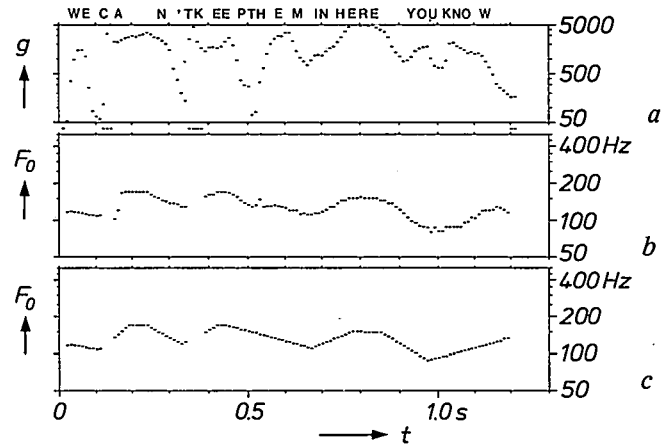


Fig. 16. Variation of the amplitude (a) and the source frequency (b) of the sentence: 'We can't keep them in here you know', in the analysis of natural speech. c) A stylized pitch contour, which the ear can hardly distinguish from that of (b).

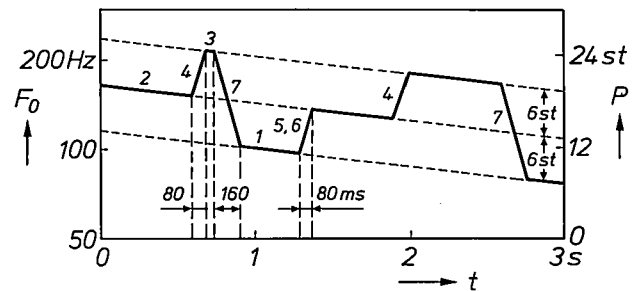


Fig. 17. General picture of pitch movements in a sentence. Variations in the pitch P are plotted on a linear scale in semitones (st), which corresponds to a logarithmic scale for the frequency F_0 as used for the analysis data. The pitch goes up and down between three straight declination lines (1, 2 and 3) at distances of 6 st, via three kinds of rise (4, 5, 6) and one kind of fall (7). The rises last 80 ms, the fall 160 ms.

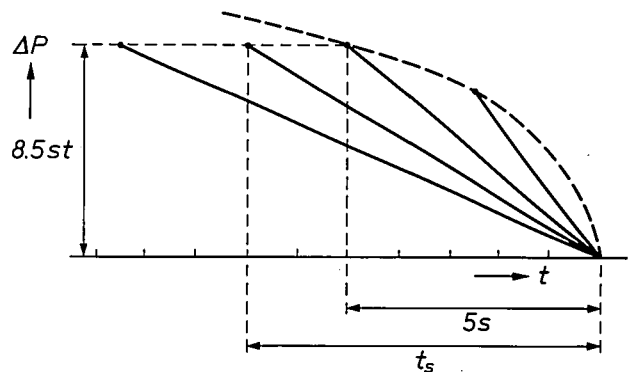


Fig. 18. The low declination line for different lengths t_s of the speech utterance. The lines are made to coincide at their end-points; this gives a good picture because for one speaker the pitch at the end of an utterance is always found to be roughly the same. For the slope D (in semitones per second) and the difference between the pitch at the beginning and the end $\Delta P = -Dt_s$ (in semitones) the following empirical rules apply:

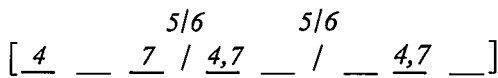
$$D = -11/(t_s + 1.5), \text{ hence } \Delta P = 11t_s/(t_s + 1.5) \text{ for } t_s \leq 5 \text{ s,}$$

$$D = -8.5/t_s, \text{ hence } \Delta P = 8.5 \text{ for } t_s > 5 \text{ s.}$$

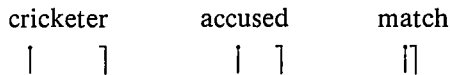
When t_s is large, ΔP is therefore constant, and when it is small ΔP decreases as shown in the figure.

ately marked by hand, or will have to be determined by text-analytical rules. These are by no means simple and we shall not consider them here. The brackets and strokes divide the text into 'blocks'.

The transitions in the sentence are distributed as follows. At the strokes, and only there, there are transitions 5 or 6. Each block therefore contains one rise 4 and one fall 7 (see fig. 17). If there is one stressed word, both rise and fall occur in that word; if there are two, then the 4 occurs in the first and the 7 in the second. Cases with more than two stressed words per block are not considered here. We illustrate these rules in the following example; thin dashes represent unstressed words and thick dashes represent stressed words:



To permit a decision as to where the rise 4 or the fall 7 (or both) occurs in or near a stressed word, the speech sound of that word must be provided with a 'flag' | in the memory, which marks the *vowel onset of the syllable with the lexical stress*. For the decision on the choice between 5 and 6 and their precise location, *the end of the last tonal part of the word sound* must also be marked (|). Examples:



(The speech sounds here are replaced for convenience by the words in ordinary spelling.)

Where the 'parsing of the sentence' now requires a fall 7, this always starts 30 ms after the flag |. If the previous rise 4 occurs in the same word, it starts 80 ms before this flag (and therefore ends there; *fig. 19a*); if, on the other hand, it occurs in a previous word, it starts 30 ms before the flag | of that word (*fig. 19b*).

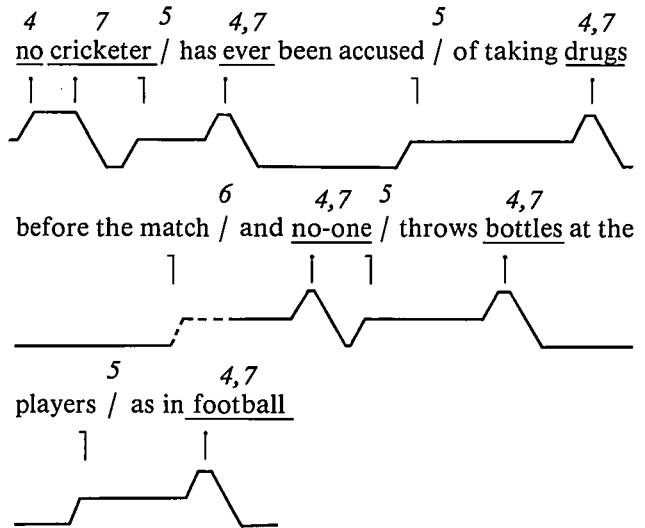
The difference between 5 and 6 is a difference between the audibility or non-audibility of the rise before the associated stroke. In the first case we have a rise 5; this ends at the flag | before the stroke (*fig. 19c*). In the other case we have a rise 6; this starts at that flag (*fig. 19d*). The rise 6 is therefore always virtual or partly virtual.

There still remains the choice between 5 and 6. The rise 6 is chosen if one or more of the following situations occurs (see *fig. 20*):

- *The last syllable before the stroke* has the lexical stress (whether that word is stressed or not) or it only has a 'short' tonal part (*fig. 20a* and *b*).
- *The first word after the boundary* is stressed and its first syllable has the lexical stress, or it is the word 'and' or 'or' (*fig. 20c* and *d*).

If neither of these situations occurs, the choice falls on 5. To exemplify *fig. 20b*, the caption to *fig. 20* gives a few words with a 'short' and a 'long' tonal part in the last syllable.

These rules enable the pitch contour to be established. The result for the example is:



The stroke after 'match' receives a 6 for two reasons: the last tonal part before the stroke is short, and after the stroke comes 'and'. The four other strokes receive a 5: none of the situations in *fig. 20* occurs here. In the example only the word flags relevant to the sentence analysis are noted. The declina-

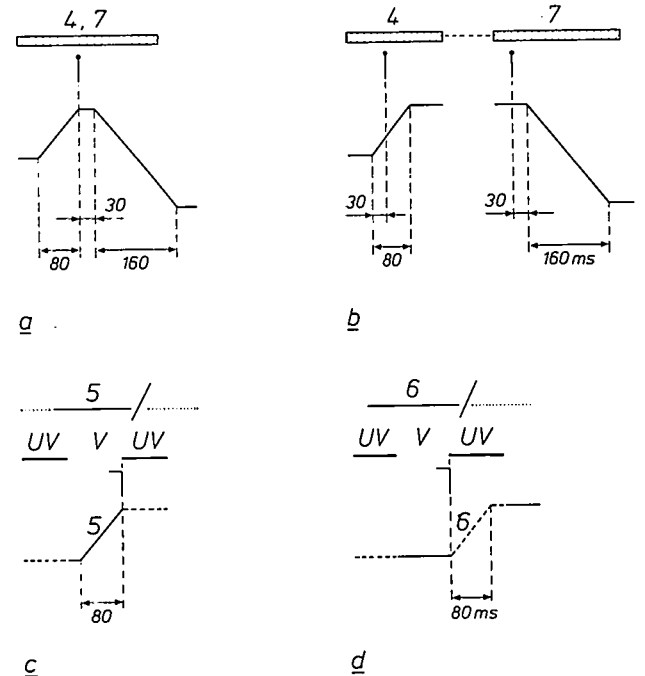


Fig. 19. *a, b*) Location of the transitions 4 and 7 relative to |, *a*) for one stressed word, *b*) for two stressed words in a block. *c, d*) Location of the rises 5 and 6 relative to the flag | at the last transition V to UV before a block boundary; 5 is audible, 6 (at least partly) is inaudible (virtual).

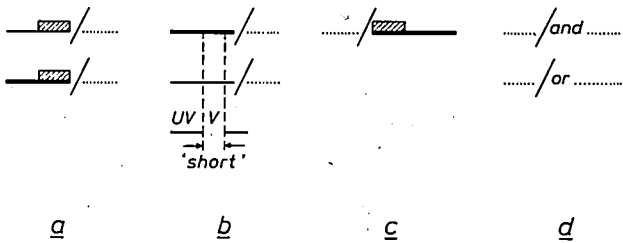


Fig. 20. Cases where the rise δ should be chosen at a boundary. The horizontal lines represent words, the thick lines stressed words, and a block represents a syllable with lexical stress. a) Last syllable before the stroke has the lexical stress. b) Last syllable before the stroke has a 'short' tonal part. c) First word after the stroke is accentuated and its first syllable has the lexical stress. d) First syllable after the stroke is 'and' or 'or'. Examples of words with a 'short' or 'long' tonal part in the last syllable:

short: match, it, bat.
 long: cricketer, really, bad.

tion is not included in the contour. Since the speech sounds here are also replaced by the words in ordinary spelling, the time structure is probably distorted.

A contour of this type only gives *changes* in pitch; the pitch itself still has to be established at one point. We have found that for a single speaker all the sentences end at approximately the same pitch (see fig. 18); the choice of that pitch establishes the pitch of the entire speech utterance. For a simulated male voice 75 Hz is a suitable final value of the source frequency, and for a simulated female voice 150 Hz.

In this article we have confined ourselves to a simple version of an intonation grammar with seven pitch movements. A refinement of our method enables us to distinguish at least 18 pitch movements for British English [10]. The simpler version discussed in the article represents the most common type of intonation pattern in British English speech, however, and therefore seems to us to be very suitable for pitch regulation in synthetic speech.

[10] J. R. de Pijper, IPO Annual Progress Report 15, 54, 1980.

Summary. At the Institute for Perception Research (IPO) speech is studied with the aid of a system called 'SPARX', for SPeech Analysis and Resynthesis eXperiments. It is based on the commonly accepted 'synthesis model' for speech production, consisting of a sound source that produces a 'voiced' pulse series or 'unvoiced' noise both with a flat spectral envelope, plus a variable filter that brings about all spectral effects of the speech organs. In SPARX an incoming speech signal is digitized and then analysed — by means of linear predictive coding (LPC) and a modification of the autocorrelation method — into the thirteen parameters of the model. These are the binary parameter that indicates whether the source is voiced or unvoiced, the source frequency, the amplitude, and the five frequencies and bandwidths of the 'formants', which characterize the filter. From these parameters the speech sound can be resynthesized, and manipulated by operations on the parameters. Examples discussed are coarsening of the parameters for the purpose of storing speech signals more economically, and manipulation of the pitch and time structure in such a way as to alter the meaning of a speech utterance. Finally an intonation grammar is discussed that seems suitable for regulating the pitch of synthetic speech. SPARX is used here for determining the extent to which the pitch variation can be 'stylized' without degrading the perception.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- H. A. Algra & J. M. Robertson:** Inhomogeneities in horizontally dipped LPE (La,Ga):YIG films studied by spin wave resonance.
J. appl. Phys. **50**, 4295-4301, 1979 (No. 6). *E*
- P. M. Asbeck, D. A. Cammack, J. J. Daniele & V. Klebanoff:** Lateral mode behavior in narrow stripe lasers. IEEE J. QE-15, 727-733, 1979 (No. 8). *N*
- H. M. J. M. van Ass:** Diffusion kinetics of some glass systems used for the production of optical fibres. J. non-cryst. Solids **33**, 325-334, 1979 (No. 3). *E*
- R. N. Bates & M. D. Coleman:** Millimetre wave finline balanced mixers. Conf. Proc. 9th Eur. Microwave Conf., Brighton 1979, pp. 721-725. *R*
- M. Berth & C. Venger:** L'implantation ionique dans GaAs pour transistors à effet de champ et circuits intégrés. Acta Electronica **23**, 23-35, 1980 (No. 1). *L*
- M. Binet:** Mesure rapide des profils de concentration et de mobilité des couches minces de GaAs. Acta Electronica **23**, 53-61, 1980 (No. 1). *L*
- H. Bouma** (Institute for Perception Research, Eindhoven): Introduction to section 'Letter and word recognition'. Processing of visible language, Vol. 1, ed. P. A. Kolars, M. E. Wrolstad & H. Bouma, pp. 221-225; Plenum Press, New York 1979.
- H. Bouma** (Institute for Perception Research, Eindhoven): Introduction to section 'Technological media for visual presentation'. Processing of visible language, Vol. 1, ed. P. A. Kolars, M. E. Wrolstad & H. Bouma, pp. 447-449; Plenum Press, New York 1979.
- D. G. Bouwhuis** (Institute for Perception Research, Eindhoven): Word knowledge and letter recognition as determinants of word recognition. Processing of visible language, Vol. 1, ed. P. A. Kolars, M. E. Wrolstad & H. Bouma, pp. 269-281; Plenum Press, New York 1979.
- J. J. M. Braat & P. F. Greve:** Aplanatic optical system containing two aspheric surfaces. Appl. Optics **18**, 2187-2191, 1979 (No. 13). *E*
- F. J. A. den Broeder & J. van der Borst:** Fe₈₀B_{20-x}Si_x glasses: A study of some physical properties as depending on metalloid content. J. appl. Phys. **50**, 4279-4282, 1979 (No. 6). *E*
- J. J. van den Broek, H. Donkersloot, G. van Tendeloo* & J. van Landuyt*** (* Rijksuniversitair Centrum Antwerpen): Phase transformations in pure and carbon-doped Al₄₅Mn₅₅ alloys. Acta metall. **27**, 1497-1504, 1979 (No. 9). *E*
- E. Bruninx, A. van Eenbergen & A. Schouten:** The determination of zinc, copper, lead and manganese in phosphate-containing matrices by coprecipitation on iron(III) hydroxide and x-ray fluorescence spectrometry. Anal. chim. Acta **109**, 419-423, 1979 (No. 2). *E*
- K. H. J. Buschow:** Intermetallic compounds of rare earths and non-magnetic metals. Rep. Prog. Phys. **42**, 1373-1477, 1979 (No. 8). *E*
- K. H. J. Buschow, U. Goebel** (T.H. Darmstadt) & **E. Dormann** (Univ. Bayreuth): Valence instabilities in CeSn₃ and YbAl₃. Phys. Stat. sol. (b) **93**, 607-615, 1979 (No. 2). *E*
- J. P. Chané & J. Hallais:** La croissance épitaxiale de GaAs pour transistors à effet de champ. Acta Electronica **23**, 11-21, 1980 (No. 1). *L*
- T. A. C. M. Claasen & W. F. G. Mecklenbräuker:** The Wigner distribution — a tool for time-frequency signal analysis, Part II: Discrete-time signals. Philips J. Res. **35**, 276-300, 1980 (No. 4/5). *E*
- T. A. C. M. Claasen, W. F. G. Mecklenbräuker, J. B. H. Peek & N. van Hurck:** Signal processing method for improving the dynamic range of A/D and D/A converters. Proc. 1979 Int. Symp. on Circuits and systems (ISCAS), Tokyo, pp. 193-196. *E*

- T. E. G. Daenen:** Cyclic reaction mechanism in the electrodeposition of aluminium.
Nature **280**, 378-380, 1979 (No. 5721). *E*
- Ph. Delsarte:** A refined version of Khachian's algorithm.
Philips J. Res. **35**, 307-319, 1980 (No. 4/5). *B*
- J. G. Dil & B. A. J. Jacobs:** Apparent size of reflecting polygonal obstacles of the order of one wavelength.
J. Opt. Soc. Amer. **69**, 950-960, 1979 (No. 7). *E*
- P. Eckerlin & S. Garbe:** Analysis of tungsten compounds in the wall region of halogen lamps.
Philips J. Res. **35**, 320-325, 1980 (No. 4/5). *A*
- W. G. Essers & R. Walter:** Some aspects of the penetration mechanisms in metal-inert-gas (MIG) welding.
Arc physics and weld pool behaviour, Int. Conf., London 1979, Vol. 1, pp. 289-300; 1980. *E*
- S. Garbe:** Morphological effects of tungsten filaments in halogen lamps: their influence on the temperature distribution and evidence for hot spot growth due to different faceting.
Philips J. Res. **35**, 326-336, 1980 (No. 4/5). *A*
- P. J. Gibson:** The Vivaldi Aerial.
Conf. Proc. 9th Eur. Microwave Conf., Brighton 1979, pp. 101-105. *R*
- J.-M. Goethals & C. Couvreur:** A cryptanalytic attack on the Lu-Lee public-key cryptosystem.
Philips J. Res. **35**, 301-306, 1980 (No. 4/5). *B*
- G. Harding:** Dose rate control in tomography — a study of image quality from a transfer function standpoint.
Radiol. diagn. **20**, 581-586, 1979 (No. 4). *H*
- M. Helmig:** Microdensitometrie en beeldanalyse.
Fotonica-meded. **5**, No. 3, 11-17, 1979 (July). *E*
- R. N. Jackson:** Flat television display — the shape of things to come.
Electronics and Power **25**, 615-621, 1979 (Sept.). *R*
- B. A. Joyce:** Present status and future directions for MBE.
Surface Sci. **86**, 92-101, 1979. *R*
- E. Klotz, R. Linde, U. Tiemens & H. Weiss:** Synthetische Tomogramme durch codierte Abbildung.
Radiol. diagn. **20**, 587-603, 1979 (No. 4). *H*
- A. J. Linssen & H. L. Peek:** Effects of trichloroethane on generation and annihilation of stacking faults during oxidation of (100) silicon.
Philips J. Res. **35**, 263-275, 1980 (No. 4/5). *E*
- J. Lohstroh:** ISL, a fast and dense low-power logic, made in a standard Schottky process.
IEEE J. SC-14, 585-590, 1979 (No. 3). *E*
- J. Lohstroh:** Static and dynamic noise margins of logic circuits.
IEEE J. SC-14, 591-598, 1979 (No. 3). *E*
- G. M. Martin, G. Jacob & G. Poiblaud** (RTC La Radiotechnique-Compelec, Caen): Les matériaux GaAs semi-isolants: paramètres essentiels et méthodes de caractérisation.
Acta Electronica **23**, 37-51, 1980 (No. 1). *L*
- D. Meignant & A. Mitonneau:** Caractérisation de pièges dans la couche active de transistors à effet de champ au GaAs.
Acta Electronica **23**, 81-90, 1980 (No. 1). *L*
- R. Memming:** Charge transfer processes at semiconductor electrodes.
Electroanalytical chemistry **11**, ed. A. J. Bard, pp. 1-84; Dekker, New York 1979. *H*
- P. C. Müräu, R. Liebert & B. M. Singer:** An electrophoretic X-ray imaging device.
IEEE Trans. ED-26, 1153-1155, 1979 (No. 8). *N*
- M. Rocchi:** Post-caractérisation des transistors à effet de champ et des circuits intégrés en GaAs sur motifs de test.
Acta Electronica **23**, 63-80, 1980 (No. 1). *L*
- R. Schäfer:** Direct solution of the radiative transfer equation for plane-parallel atmospheres.
J. quant. Spectrosc. rad. Transfer **23**, 455-466, 1980 (No. 5). *A*
- J. M. Shannon:** Hot-electron camel transistor.
IEE J. Solid-St. Electron Dev. **3**, 142-144, 1979 (No. 5). *R*
- M. Sintzoff:** Certification et programmation par approximation.
Actes des Journées Francophones sur la Certification du Logiciel, Genève 1979, pp. 102-113. *B*
- E. T. J. M. Smeets & J. Politiek:** Very-low-noise silicon avalanche photodiodes made by the channeling of aluminium in $\langle 110 \rangle$ silicon.
Appl. Phys. Letters **35**, 112-113, 1979 (No. 2). *E*
- J. L. Teszner:** Introduction (*to issue on Gallium arsenide for field effect transistors and integrated circuits*).
Acta Electronica **23**, 7-9, 1980 (No. 1). (*In English and in French.*) *L*
- M. Urner-Wille & K. Witter:** Compensation point switching in homogeneous amorphous GdFe-films.
J. Magn. magn. Mat. **13**, 77-80, 1979 (No. 1/2). *H*
- H. J. Verbeek** (Philips Centre for Technology, Eindhoven): Tribological systems and wear factors.
Wear **56**, 81-92, 1979 (No. 1).
- H. Verweij & J. H. J. M. Buster** (Philips Semiconductor Devices Factory, Nijmegen): The structure of lithium, sodium and potassium germanate glasses, studied by Raman scattering.
J. non-cryst. Solids **34**, 81-99, 1979 (No. 1). *E*
- J. H. Waszink & G. J. P. M. van den Heuvel:** Measurements and calculations of the resistance of the wire extension in arc welding.
Arc physics and weld pool behaviour, Int. Conf., London 1979, Vol. 1, pp. 227-239; 1980. *E*

Contents of Philips Telecommunication Review 38, No. 3, 1980:

- W. G. Bax, Ph. Uythoven & J. Wagenmakers:** Field trial of a 140 Mb/s coaxial line system (pp. 93-103).
D. W. Rollema: A maritime traffic control centre (pp. 104-120).
D. W. Rollema: Coast Guard radar for four lighthouses on Dutch coast (p. 121).
S. T. Soames & R. A. Mulkerrin: HERALD, a small business telephone communication system (pp. 122-130).
G. Schouten: A figure of merit for solderability (pp. 131-138).

Contents of Philips Telecommunication Review 38, No. 4, 1980:

- W. M. Pannell:** A quasi-sync radio paging system (pp. 141-148).
A. C. Steenhuisen: Man-machine language for digital PRX telephone systems (pp. 149-167).
F. P. van Enk & T. Ryan: Remote control of base stations in mobile radio (pp. 168-175).
R. Baird & F. P. van Enk: The DS-1002 high-speed data system for mobile radio (pp. 176-186).
J. Noordanus & K. Everaarts: End-to-end supervision of digital radio relay sections (pp. 187-194).

Contents of Electronic Components and Applications 2, No. 3, 1980:

- F. J. Burgum & E. B. G. Nijhof:** Inverter circuit for a PWM motor speed control system (pp. 130-142).
H. W. Evers: Mains pollution caused by domestic appliances, Part 3 — Voltage fluctuation and flicker (pp. 143-149).
A. J. Rees & E. I. Várszegi: Electret microphone for telephony (pp. 150-164).
Asymmetric J-FET improves radio performance (pp. 165-169).
W. Hesse & U. Schillhof: Digital control of radio and audio equipment, Part 7 — RTS tuning controls and the microcomputer (pp. 170-174).
J. Fasser & A. M. L. Hodemaekers: Interrupted current-loop dialling for pushbutton telephones (pp. 175-189).

Contents of Electronic Components and Applications 2, No. 4, 1980:

- F. Burgum, E. B. G. Nijhof & A. Woodworth:** Gate turn-off switch (pp. 194-202).
J. A. A. den Ouden: Development of digital filters using the 8X300 microprocessor (pp. 203-214).
D. J. van der Wal: Interference suppression for FM radios (pp. 215-218).
B. G. Starr & J. C. F. van Loon: LSI circuit for AC motor speed control (pp. 219-229).
R. J. van de Plassche: Monolithic 14-bit DAC with 85 dB S/N ratio (pp. 235-241).
D. J. G. Janssen & L. van de Meeberg: PCM codec with on-chip digital filters (pp. 242-250).

Contents of Electronic Components and Applications 3, No. 1, 1980:

- CQL10 semiconductor laser for information readout (pp. 2-5).
W. B. Rosink: Analogue control system for a.c. motor with PWM variable speed drive (pp. 6-15).
A. Franken & W. Lohuis: Highlight handling with diode-gun 'Plumbicon' tubes (pp. 17-20).
H. J. H. van Heffen: Ceramic permanent magnets for d.c. motors, Part 1 — Performance equations (pp. 22-30).
B. H. A. Goddijn: New hybrid stepping motor design (pp. 31-37).
A. P. M. Moelands: Serial I/O with the MAB8400 series microcomputers (pp. 38-46).
P. R. Brennand & B. Murray: Frequency synthesiser using LSI devices (pp. 47-61).

Recent United States Patents

Abstracts from patents that describe inventions from the following research laboratories, which form part of or cooperate with the Philips group of companies:

Philips Research Laboratories, Eindhoven, The Netherlands	E
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	R
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	A
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	B
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N

4 272 574

Optically readable information disc

G. J. M. Lippits

A. J. M. van den Broek

R. Dijkstra

The invention relates to an information disc having a laminated structure which can be read optically. The information disc comprises a transparent substrate which is preferably manufactured from a synthetic resin, for example plexiglass, having thereon a radiation-cured lacquer layer in which the information track is present. The lacquer layer used comprises a radiation cross-linkable protic compound which after curing is aprotic. The lacquer layer preferably comprises a polythiol compound as well as a polyene compound in an equivalent ratio of 1:1.

E

4 275 091

Method of duplicating plastic information carriers

G. J. M. Lippits

A. J. M. van den Broek

A. J. G. Op het Veld

R. Dijkstra

J. de Jonge

The invention relates to a method of reproducing plastic record carriers, in particular duplicating video records. According to the invention, a metal die is used which is provided with a thin-liquid molding resin of a particular composition which can be polymerized by radiation. A radiation-pervious substrate which is manufactured from synthetic material, for example polymethylmethacrylate, is provided on the molding resin. The molding resin is exposed to light via the substrate after which the cured molding resin together with the substrate connected thereto is removed from the die. The molding resin used in the process comprises low-molecular monomers or oligomers which contain on an average 25-70% by weight of hydrocarbon groups and/or phenyl groups. The molding resin is aprotic and has a functionality as regards unsaturation which is between the values 2 and 6. A suitable molding resin contains mono-, tri- or tetra esters of acyclic acid. The molding resin preferably has a swelling capacity with respect to the substrate and for that purpose preferably comprises a vinylmonomer. The metal die used in the method preferably is a quite flat die which is obtained by providing the master disk which is a flat glass plate with information track with a nickel layer, gluing hereon a flat stiffening plate and then removing the master disk. The resulting father disk may be used as a die. Alternatively, further metal copies may be made herefrom which are provided in a simpler manner with a flat stiffening plate. The invention also extends to the molding resin, substrate and die used in the method, as well as to the resulting plastic record carriers.

E

4 272 776

Semiconductor device and method of manufacturing same

B. H. Weijland

W. H. C. G. Verkuijden

An inset oxide isolated integrated circuit, with multiple levels of inset oxide, polycrystalline regions, and channel stops.

E

4 272 995

Ionization flow meter

M. P. Weistra

An ionization flow meter which can determine with high accuracy the flow of a gas expressed as a gas velocity, volume flow, or mass flow and substantially independently of pressure, temperature and, as the case may be, moisture content. By applying the theory of corona discharges, data relating to the flow-determining factors such as ion mobility and gas density can be derived from the known measured voltage and current values using electronic means. For example, the slope S of the I-V characteristic curve can be determined from the measured values of voltage V and current I .

E

4 276 494

Cathode ray tube with transversely supported electrode and conductive wall coating

J. H. T. van Roosmalen

G. A. H. M. Vrijssen

In a cathode-ray tube, in particular a camera tube, the inner wall of the glass envelope is coated with an electrically conductive material interrupted in the proximity of electrodes extending transversely to

E



the wall coating and supported by transversely extending supporting surfaces. At the area of each of the interruptions the envelope has a stepwise decrease of the inside diameter in two steps. In the direction of decreasing diameter the first of these steps constitutes the supporting surface for the transverse electrode and the interruption in the wall coating is provided on a wall portion of the second of these steps. The interruptions provided in this manner do not exert any disturbing influence on the electron beam in the tube.

4 276 529

Magnet coil arrangement for generating a homogeneous magnetic field for magnetic resonance arrangements

J. Heinzerling
R. Rieckeheer

H

The invention relates to a magnet coil arrangement for generating a magnetic field which is homogeneous at least in its center, preferably for magnetic resonance spectroscopy. The arrangement consists of four flat ring coils whose coil planes extend perpendicularly to an axis of examination which extends through the coil centers. The coils are symmetrically arranged with respect to a point situated on this axis. By means of a magnetic coil arrangement of this kind, a magnetic field can be generated which extends in the direction of the axis of examination and rotationally-symmetrically with respect thereto and which, in comparison with the magnetic fields generated by means of known coil arrangements, has an improved homogeneity in the direction perpendicular to the axis of examination over a larger range which extends in the direction perpendicular to the axis of examination in the center of the magnetic coil arrangement.

4 276 611

Device for the control of data flows

P. G. Jansen
J. L. W. Kessels

E

A commutation device for the selective control of data transport. At least two data inputs and data outputs, each of the latter having a buffer for storing a data word. A number of possibilities of data transport can be selectively controlled, four for a single connection and two different ones for pair-wise connection. Seven input control lines are provided, two lines for receiving a signal which indicates whether information is present on the associated input line, two lines for indicating the selected output buffer, two erase lines for making a data buffer freely accessible after output of data from the data buffer, and one priority line for granting priority to one of the two input lines if both lines select the same data buffer. There are four output control lines, two lines which indicate that the data present on the input lines have been taken up in the selected output buffer, and two lines which indicate whether an output buffer contains data. The commutation device can effect the data transport itself and can be grouped in specific arrangements to form a buffer in which the data partly determine their own path.

4 276 649

Receiver for digital signals in line code

G. C. Groenendaal
E. A. Aagaard

E

Receiver for a digital line code signal. This receiver comprises a line code decoder and a digital-to-analog converter. To reduce the audibility of bursts this receiver also comprises a line-code violation detector detecting whether the received signal deviates from the line code; as well as a pulse generator. Each time the line-code violation detector detects that the received signal deviates from the line code, the output signal of the pulse generator is applied to the digital-to-analog converter instead of the output signal of the line code decoder.

4 276 650

Method of synchronizing a quadphase receiver and clock synchronization device for carrying out the method

F. de Jager
R. A. van Doorn
J. J. Verboom
M. G. Carasso

E

The invention relates to a method for the clock synchronization of a receiver for demodulating a quadphase coded data signal and to a clock synchronization device for carrying out the method. In the method according to the invention the first bit is compared (correlated) with the third bit and the second bit with the fourth bit: a high degree of correlation indicates that synchronization has been obtained and a low degree indicates absence of synchronization.

4 277 138

Diffraction grating and system for the formation of color components

H. Dammann

H

A device for spatially separating specific spectral regions, preferably of color components from a wideband spectrum which is actively and/or passively radiated by objects. The spectral regions, or color components, are derived from the diffraction orders of a diffraction grating (phase grating), which is disposed in the pupil of an imaging lens and whose groove profile consists of several steps, which produce path length differences which are integral multiples of a specific wavelength.

4 277 542

Resistance material

A. H. Boonstra
C. A. H. A. Mutsaers
F. N. G. R. van der Kruijs

E

Resistance material consisting of a mixture of metal oxidic compounds, metal oxides, a permanent binder and a temporary binder, the resistance-determining component consisting of barium-rhodate $\text{BaRh}_6\text{O}_{12}$. This component has a linear positive temperature coefficient of the resistance (TCR) and enables the production of a resistor having a very low TCR by combining the material with a material having a negative TCR. The resistor is obtained by firing this resistance material after it has been applied onto a substrate.

4 277 686

Device for determining internal body structures by means of scattered radiation

G. Harding

H

The invention relates to a device for measuring a scatter coefficient distribution in a plane of a body. The plane is irradiated in different directions by a primary radiation beam along beam paths which are each time situated in parallel in a direction. Scattered radiation which is generated by a primary radiation beam along its path is measured by detectors which are situated on both sides of the plane and which enclose the body as completely as possible. The scatter coefficient distribution is determined by iteration by calculating a scatter value for each beam path from an assumed distribution and by comparing this scatter value with the associated measured scattered radiation. From the difference between calculated and measured values a correction is determined and taken up in the calculated value.

4 277 713

Low-pressure gas discharge lamp and method for making

J. Hasker E
J. C. G. Vervest
C. Peters
L. C. J. Vroomen

Low-pressure discharge lamp having an elongate discharge vessel which contains a thinly distributed filamentary body permeable to the gas discharge, said body comprising a helical support filament which is supported by the inner surface of the discharge vessel and is at least one further filament, supported by the support filament and extending therefrom towards the axis of the discharge vessel.

4 278 888

Apparatus for determining the spatial distribution of the absorption of radiation in a body

W. Wagner H

A computed tomography device wherein a body contour outside an examination area is determined by measurements made with the aid of an auxiliary radiation source (for example light or ultrasound).

4 278 912

Electric discharge tube having a glass-sealed electric leadthrough and method of manufacturing such an electric leadthrough

G. A. H. M. Vrijssen E
J. P. T. Franssen

An electric discharge tube is provided with a hermetically sealed leadthrough which electrically connects electrodes on the inner and outer walls of the envelope. The leadthrough consists of an aperture in the envelope having a conductive layer provided on the wall of the aperture. The aperture is hermetically sealed by means of a plug of thermally devitrified glass which is provided in the form of a suspension of a devitrified glass powder in an organic binder. To manufacture the leadthrough, the envelope of the tube is subjected to temperature treatments in which at a first temperature range the binder is fired from the suspension in an oxygen-containing atmosphere, and at a second temperature range the devitrifiable glass is devitrified in a non-oxidizing atmosphere. A hermetically sealed leadthrough results, without excessive oxidation of the electrodes, while the deformation of the glass envelope at the area of the leadthrough is avoided.

4 279 157

Method of and device for determining the internal structure of a body by means of acoustic beams

H. Schomberg H
M. Tasto

A method and device for determining the internal structure of a body by means of acoustic beams. Transit times and intensities of acoustic beams passing through the body in different spatial directions are measured to establish the refractive index distribution and the acoustic absorption coefficient distribution, respectively at the points of a point matrix associated with the body. The non-rectilinear course of the acoustic beams is taken into account in this respect. This results in reconstructed images of higher quality.

4 279 253

Epilation apparatus

F. Haes E
G. M. P. G. Hermes
C. M. Reijnhout

There is provided an epilation apparatus comprising a drivable member having a hair-gripping wall, and a stationary complementary member having a confronting wall spaced from the hair-gripping wall of the drivable member to provide a hair gap therebetween.

4 280 049

X-ray spectrometer

H. W. Werner E
A. W. Witmer
W. F. Knippenberg

An X-ray spectrometer which is arranged inside an evacuable housing and which comprises a wavelength dependent X-ray detection system and, for irradiating the specimen to be examined, an electron source with an electron deflection system for generating an electron beam and an X-ray source for generating an X-ray beam. The X-ray source consists of an anticathode on which the electron beam can be directed by the electron deflection system in order to generate the X-ray beam.

4 280 068

Bulk channel charge coupled device having improved input linearity

P. J. Snijder E

In bulk channel charge coupled devices the nonlinearity in the input characteristic caused by varactor effects is removed by moving the potential well in which the charge packets are generated below the input electrode to the surface where the center of electrical charge is substantially independent of the value of the charge. Said shift can be obtained by external means, for example an extra d.c. voltage at the input electrode, or by internal means, for example a thicker oxide below the input electrode.

4 280 089

Automatic incrementing attenuation arrangement

R. J. van de Plassche E
E. C. Dijkmans

Attenuation arrangement comprising a step attenuator arranged in cascade with a controllable voltage divider via first and second voltage terminals, the step attenuator comprising a series arrangement of attenuation elements for dividing a voltage applied across said series arrangement into a plurality of voltage increments, which voltage increments are individually switchable between the two voltage terminals for varying the output voltage of the controllable voltage divider for the voltage range of the relevant voltage element, the direction of the polarity of the voltage between the two voltage terminals changing at a switch-over from one voltage increment to an adjacent voltage increment.

4 280 158

Magneto-resistive reading head

E. de Niet E

A magnetic reading head having a magneto-resistive element which is connected to a reading amplifier. In order to reduce the modulation noise (Barkhausen effect) when making the relationship between the resistance variation and the strength of the signal field linear in a negative feedback loop of the reading amplifier an electric turn is present which turn is positioned relative to the magneto-resistive element in such manner that a negative feedback field (H_t) can be generated with it which causes a magnetic flux in the element which is directed oppositely to the magnetic flux caused in the element by a magnetic field (H_y) to be detected.

4 280 858

Method of manufacturing a semiconductor device by retarding the diffusion of zinc or cadmium into a device region

C. J. M. van Oopdorp E
H. Veenliet

A semiconductor device and a method for manufacturing the semiconductor device are disclosed for forming an abrupt and accurately positioned p-n junction between a substrate and a substrate-adjointing region. This is achieved in accordance with the present

invention by diffusing zinc or cadmium from a surface of the substrate-adjointing region to the substrate, and abruptly limiting or retarding the diffusion of the zinc or cadmium into the substrate near a junction between the substrate and the region. This is accomplished in accordance with the present invention by selecting the net donor concentration in the substrate near the junction to be higher than the concentration of zinc or cadmium at the surface of the substrate-adjointing region.

4 281 396

Magnetic strip domain memory system

J. Roos

E

A magnetic memory device in which information is stored in the form of strip domains in a layer of magnetic material supported by a layer of ferromagnetic material. The ferromagnetic material contains a pattern of alternately magnetized strips for sustaining a magnetic field periodically varying in a first coordinate direction and directed transverse to the domain layer. The device also includes a generator for receiving and converting data into configurations of the strip domains in the plate.

4 283 226

Method of preparing titanium iron-containing material for hydrogen storage

H. H. van Mal

E

H. A. van Esveld

J. S. van Wieringen

K. H. J. Buschow

A material for storing hydrogen consisting of a titanium-iron alloy having 5-30 at. % of one or more metals of the group chromium, zirconium, manganese and vanadium.

4 283 689

Microwave oscillator circuit with improved efficiency

H. Tjassens

E

A microwave oscillator circuit, suitable for use as the local oscillator in beam transmitters, radar systems and satellite TV receivers, comprises an active element (IMPATT, Gunn diode) at one end of a coaxial transmission line which is terminated at its other end by a matched load. At a suitable distance from the diode a high Q transmission cavity resonator is coupled to the transmission line via a first coupling hole. A drawback of such a circuit is that a portion of the oscillator power at the required oscillator frequency f_0 is dissipated in the terminal impedance Z_0 . This is obviated by coupling the transmission resonant cavity to the transmission line via a second coupling hole. The distance between the first and the second coupling hole is $\frac{1}{4}\lambda$. As a result the terminal impedance, which has been transformed very frequency-selective to a very high value in situ of the second coupling hole, is transformed to a very low value at the first coupling hole and very little power is dissipated in this low impedance, so that a considerable improvement of the circuit efficiency has been achieved.

4 283 837

Semiconductor device and method of manufacturing same

A. Slob

E

A semiconductor device includes a silicon substrate having an insulating layer with a window. A silicon layer is deposited on the insulating layer and on the silicon substrate surface in the window. This silicon layer has n-type and p-type conductive layer parts which adjoin each other within the window and which each serve as both a connection conductor and an electrode of an active zone of the device. Semiconductor devices in accordance with the invention feature very small surface areas, and are thus particularly suitable for high frequency operation.

4 284 069

Wall element comprising a solar collector which is disposed between two transparent panes

H. Hörster

A

W. Hermann

K. Klinkenberg

A wall element, comprising a solar collector which is arranged between two panes and which comprises a number of rotatable absorber plates. One side of the absorber plates is provided with a non-selective black coating and the other side is provided with a selective, heat-reflective layer. The absorber plates are accommodated in evacuated, transparent tubes.

4 286 156

Device for determining the spatial absorption distribution in a plane of examination

W. Wagner

H

The device in accordance with the invention comprises detectors, a first part of which is not struck by radiation during measurement of the useful signal, whilst a second part which is struck directly by the radiation during the measurement of the useful radiation, is shielded during a next measurement. During the last measurement, the first detector part is struck by scattered radiation. From the two signals thus formed, a signal free from scattered radiation can be obtained by subtraction.

4 286 177

Integrated injection logic circuits

C. M. Hart

E

A. Slob

An "Integrated Injection Logic" integrated circuit in which bias currents are supplied by means of a current injector. The current injector is a multi-layer structure in which current is supplied by means of injection and collection of charge carriers via rectifying junctions, to predetermined zones of the circuit to be biased. Such zones are preferably biased by charge carriers which are collected by such zones from one of the layers of the current injector. The circuit also preferably includes a region for reducing carrier injection from a predetermined zone.

4 286 266

Display device

M. de Zwart

E

J. L. A. M. Heldens

In a liquid crystal display device which upon controlling with direct voltage shows a memory effect that can be erased by alternating voltage, the written area can expand beyond the edge of the electrode so that edge zones can no longer be erased readily. This disadvantage is avoided by covering the edges of the electrodes with an insulating layer.

4 286 318

Control loop

K. A. Immink

E

A. Hoogendoorn

A control loop provided with a control unit for realizing a transfer characteristic having a number of peaks at a fundamental frequency and harmonics thereof. The control unit comprises a memory device for digitally storing a number of samples of the error signal appearing in the control loop during a cycle period equal to the period corresponding to the fundamental frequency. Furthermore, there are provided means for comparing the sample stored in the memory device with the value of the error signal one cycle period later and, depending on this comparison, correcting the memory content of the relevant memory location. The variation of the error signal stored in the memory device is furthermore cyclically employed as a control signal for the control loop.

COMPACT disc DIGITAL AUDIO

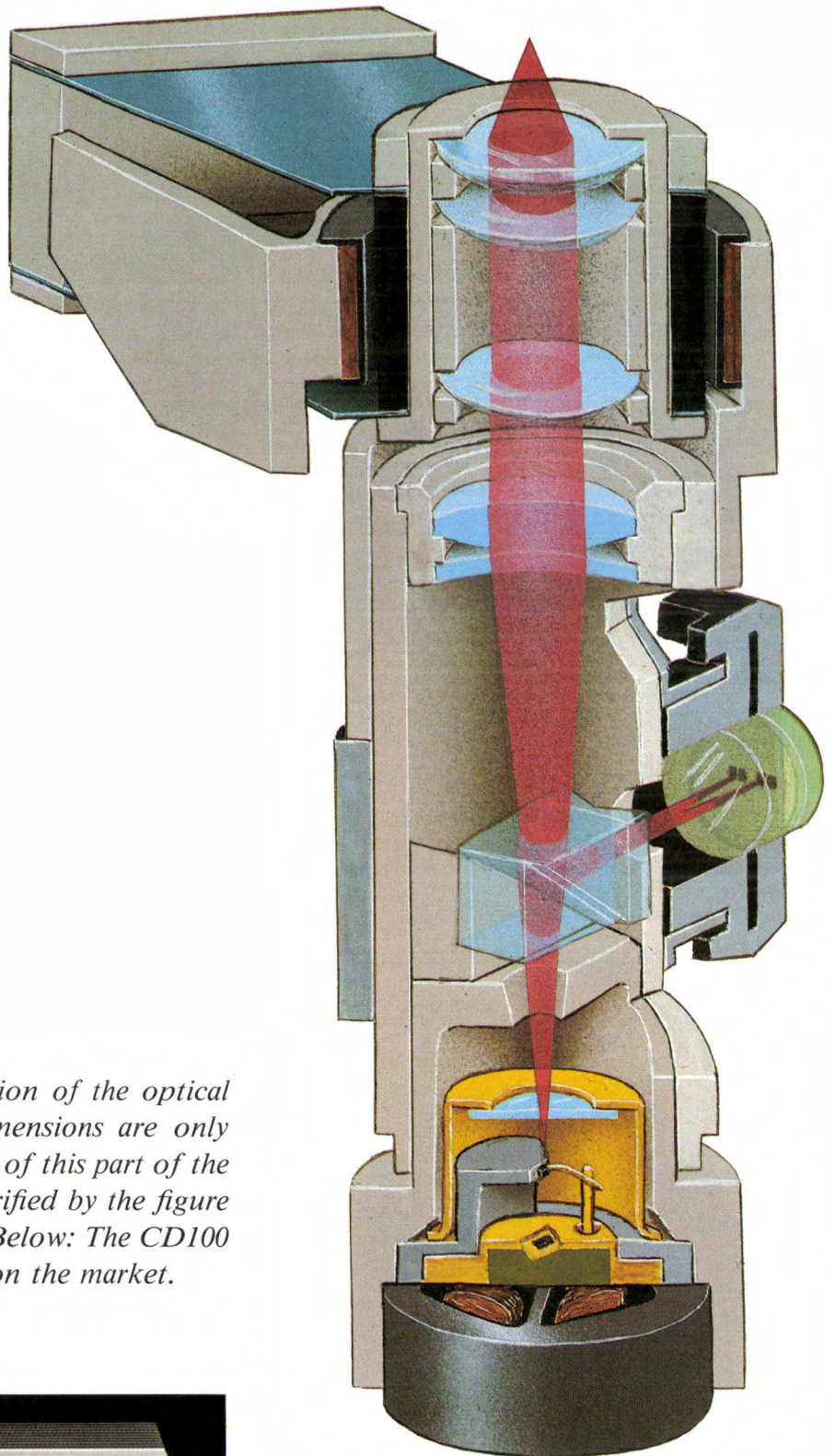
In 1877 Edison's phonograph played the nursery rhyme 'Mary had a little lamb', after he had recorded it on the wax cylinder in his own voice: the human voice had been reproduced for the first time in history. Then came Berliner's wax disc, followed by the 78-turns-per-minute shellac disc, and eventually the modern long-play record. Now when we enjoy the music from our 'LPs' at home it is almost perfectly reproduced by our hi-fi equipment. However, the record player itself is a weak link in the chain, since damage to the vulnerable disc often introduces an unwanted accompaniment of undesirable sounds to the music. This cannot happen with the Compact Disc. It is scanned optically, so playing it cannot produce any damage, and dust and fingermarks have far less effect — because errors can in fact be corrected. Another way in which the Compact Disc differs from the conventional long-play record is that the sound is recorded on the disc in digital form. The digital processing of the audio signals will perhaps be a more far-reaching development in the history of record-playing equipment than the change from acoustic to electrical reproduction was in its time. As we shall see, digital processing brings many advantages. To make the best use of it, it is necessary to build up a complete system that extends from the record-manufacturer's equipment to the record player at home. It is clear that so extensive a system only has a chance of success if

records and playing equipment from different manufacturers are all absolutely compatible.

The study of the possibility of recording audio signals optically on a disc was started in 1974 at Philips Research Laboratories, in close cooperation with the Philips Audio Division. It soon became clear that a different method would have to be used from that in the LaserVision system^[]: digital signal coding, instead of analog modulation methods. More and more people from the Research Laboratories and the Audio Division gradually became involved in the project. After an agreement 'in principle' had been reached with the Sony company in 1979, extensive technical discussions were started, mainly about the signal processing. Eventually a system standard was produced, including contributions from both companies. Then licensing agreements were made with a number of other manufacturers of audio equipment and gramophone records.*

This issue contains four articles written by staff from Philips Research Laboratories and the Compact Disc Development Laboratory of the Audio Division. They give an account of various aspects of the revolutionary Compact Disc system: the complete system, the modulation of the digital signal, the method of error correction and the conversion of the digital signal into the analog signal.

[*] Formerly called the VLP system.



Right: An artist's impression of the optical pick-up, whose actual dimensions are only 45×12 mm. The operation of this part of the playback equipment is clarified by the figure and caption on page 153. Below: The CD100 player, the first to be put on the market.



The Compact Disc Digital Audio system

M. G. Carasso, J.-B. H. Peek and J. P. Sinjou

Introduction

During the many years of its development the gramophone has reached a certain maturity. The availability of long-play records of high quality has made it possible to achieve very much better sound reproduction in our homes than could be obtained with the machine that first reproduced the sound of the human voice in 1877. A serious drawback of these records is that they have to be very carefully handled if their quality is to be preserved. The mechanical tracking of the grooves in the record causes wear, and damage due to operating errors cannot always be avoided. Because of the analog recording and reproduction of the sound signal the signal-to-noise ratio may sometimes be poor (< 60 dB), and the separation between the stereo channels (< 30 dB) leaves something to be desired.

For these and other problems the Compact Disc system offers a solution. The digital processing of the signal has resulted in signal-to-noise ratios and a channel separation that are both better than 90 dB. Since the signal information on the disc is protected by a 1.2 mm transparent layer, dust and surface damage do not lie in the focal plane of the laser beam that scans the disc, and therefore have relatively little effect. Optical scanning as compared with mechanical tracking means that the disc is not susceptible to damage and wear. The digital signal processing makes it possible to correct the great majority of any errors that may nevertheless occur. This can be done because error-correction bits are added to the information present on the disc. If correction is not possible because there are too many defects, the errors can still be detected and 'masked' by means of a special procedure. When a Compact Disc is played there is virtu-

ally no chance of hearing the 'tick' so familiar from conventional records.

With its high information density and a playing time of an hour, the outside diameter of the disc is only 120 mm. Because the disc is so compact, the dimensions of the player can also be small. The way in which the digital information is derived from the analog music signal gives a frequency characteristic that is flat from 20 to 20 000 Hz. With this system the well-known wow and flutter of conventional players are a thing of the past.

Another special feature is that 'control and display' information is recorded, as 'C&D' bits. This includes first of all 'information for the listener', such as playing time, composer and title of the piece of music. The number of a piece of music on the disc is included as well. The C&D bits also contain information that indicates whether the audio signal has been recorded with pre-emphasis and should be reproduced with de-emphasis^[1]. In the Compact Disc system a pre-emphasis characteristic has been adopted as standard with time constants of 15 and 50 μ s. In some of the versions of the player the 'information for the listener' can be presented on a display and the different sections of the music on the disc can be played in the order selected by the user.

In the first article of a series of four on the Compact Disc system we shall deal with the complete system, without going into detail. We shall consider the disc, the processing of the audio signal, reading out the signal from the disc and the reconstitution of the audio signal. The articles that follow will examine the system aspects and modulation, error correction and the digital-to-analog conversion.

Drs M. G. Carasso and Dr Ir J. B. H. Peek are with Philips Research Laboratories, Eindhoven; J. P. Sinjou is with the Philips Audio Division, Eindhoven.

^[1] See F. W. de Vrijer, Modulation, Philips tech. Rev. 36, 305-362, 1976, in particular pages 323 and 324.

The disc

In the LaserVision system^[2], which records video information, the signal is recorded on the disc in the form of a spiral track that consists of a succession of pits. The intervals between the pits are known as 'lands'. The information is present in the track in analog form. Each transition from land to pit and vice versa marks a zero crossing of the modulated video signal. On the Compact Disc the signal is recorded in a similar manner, but the information is present in the track in digital form. Each pit and each land represents a series of bits called channel bits. After each land/pit or pit/land transition there is a '1', and all the channel bits in between are '0'; see *fig. 1*.

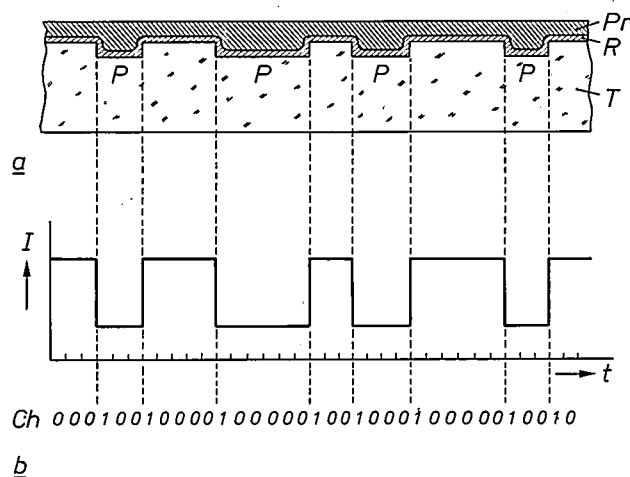


Fig. 1. *a)* Cross-section through a Compact Disc in the direction of the spiral track. *T* transparent substrate material, *R* reflecting layer, *Pr* protective layer. *P* the pits that form the track. *b)* *I* the intensity of the signal read by the optical pick-up (see *fig. 2*), plotted as a function of time. The signal, shown in the form of rectangular pulses, is in reality rounded and has sloping sides^[3]. The digital signal derived from this waveform is indicated as a series of channel bits *Ch*.

The density of the information on the Compact Disc is very high: the smallest unit of audio information (the audio bit) covers an area of $1 \mu\text{m}^2$ on the disc, and the diameter of the scanning light-spot is only $1 \mu\text{m}$. The pitch of the track is $1.6 \mu\text{m}$, the width $0.6 \mu\text{m}$ and the depth $0.12 \mu\text{m}$. The minimum length of a pit or the land between two pits is $0.9 \mu\text{m}$, the maximum length is $3.3 \mu\text{m}$. The side of the transparent carrier material *T* in which the pits *P* are impressed — the upper side during playback if the spindle is vertical — is covered with a reflecting layer *R* and a protective layer *Pr*. The track is optically scanned from below the disc at a constant velocity of 1.25 m/s . The speed of rotation of the disc therefore varies, from about 8 rev/s to about 3.5 rev/s .

Processing of the audio signal

For converting the analog signal from the microphone into a digital signal, pulse-code modulation (PCM) is used. In this system the signal is periodically sampled and each sample is translated into a binary number. From Nyquist's sampling theorem the frequency of sampling should be at least twice as high as the highest frequency to be accounted for in the analog signal. The number of bits per sample determines the signal-to-noise ratio in the subsequent reproduction.

In the Compact Disc system the analog signal is sampled at a rate of 44.1 kHz , which is sufficient for reproduction of the maximum frequency of $20\,000 \text{ Hz}$. The signal is quantized by the method of uniform quantization; the sampled amplitude is divided into equal parts. The number of bits per sample (these are called audio bits) is 32, i.e. 16 for the left and 16 for the right audio channel. This corresponds to a signal-to-noise ratio of more than 90 dB . The net bit rate is thus $44.1 \times 10^3 \times 32 = 1.41 \times 10^6$ audio bits/s. The audio bits are grouped into 'frames', each containing six of the original samples.

Successive blocks of audio bits have blocks of parity bits added to them in accordance with a coding system called CIRC (Cross-Interleaved Reed-Solomon Code)^[4]. This makes it possible to correct errors during the reproduction of the signal. The ratio of the number of bits before and after this operation is 3:4. Each frame then has C&D (Control and Display) bits, as mentioned earlier, added to it; one of the functions of the C&D bits is providing the 'information for the listener'. After the operation the bits are called data bits.

Next the bit stream is modulated, that is to say the data bits are translated into channel bits, which are suitable for storage on the disc; see *fig. 1b*. The EFM code (Eight-to-Fourteen Modulation) is used for this: in EFM code blocks of eight bits are translated into blocks of fourteen bits^[5]. The blocks of fourteen bits are linked by three 'merging bits'. The ratio of the number of bits before and after modulation is thus 8:17.

For the synchronization of the bit stream an identical synchronization pattern consisting of 27 channel bits is added to each frame. The total bit rate after all these manipulations is 4.32×10^6 channel bits/s.

^[2] See Philips tech. Rev. 33, 187-193, 1973.

^[3] See *fig. 3* of the article by J. P. J. Heemskerck and K. A. Schouhamer Immink, on p. 159 of this issue.

^[4] See H. Hoeve, J. Timmermans and L. B. Vries, Error correction and concealment in the Compact Disc system, this issue, p. 166.

^[5] See J. P. J. Heemskerck and K. A. Schouhamer Immink, Compact Disc: system aspects and modulation, this issue, p. 157.

^[6] J. C. J. Finck, H. J. M. van der Laak and J. T. Schrama, Philips tech. Rev. 39, 37, 1980.

Table I. Names of the successive signals, the associated bit rates and operations during the processing of the audio signal.

Name	Bit rate in 10^6 bits/s	Operations
Audio signal		PCM (44.1 kHz)
Audio bit stream	1.41	CIRC (+ parity bits) Addition of C&D bits
Data bit stream	1.94	EFM Addition of merging bits Addition of synchronization patterns
Channel bit stream	4.32	

disc (called the 'master'). A pattern of pits is produced on this disc by means of a photographic developing process. After the surface has been coated with a thin silver layer, an electroplating process is applied to produce a nickel impression, called the 'metal father'. From this 'father disc' impressions called 'mother discs' are produced in a similar manner. The impressions of the mother discs, called 'sons' or 'stampers', are used as tools with which the pits P are impressed into the thermoplastic transparent carrier material T of the disc; see fig. 1.

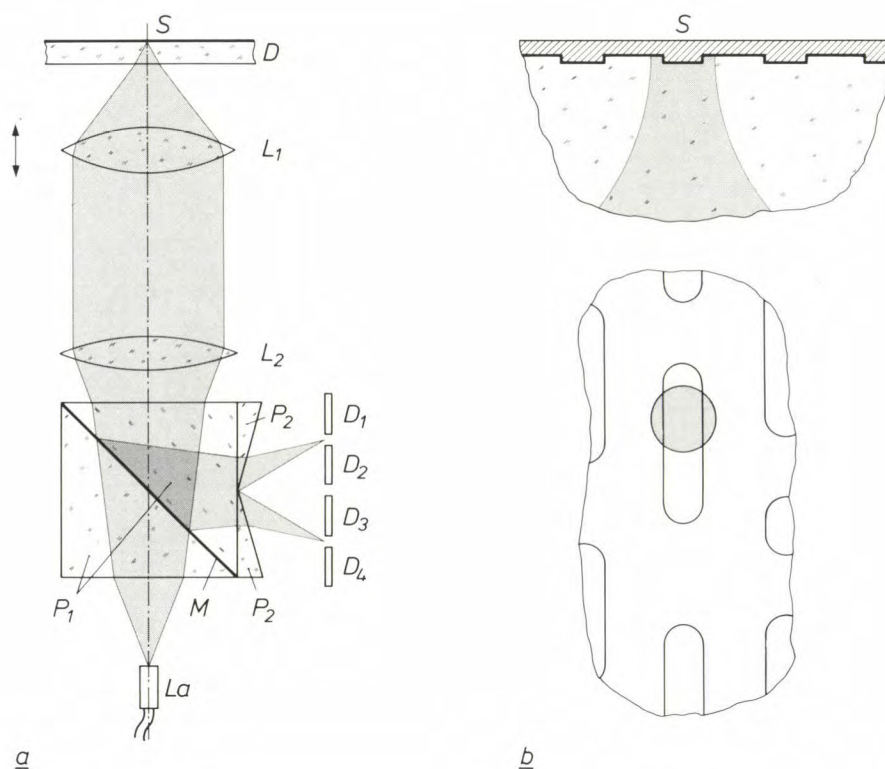


Fig. 2. *a*) Diagram of the optical pick-up. D radial section through the disc. S laser spot, the image on the disc of the light-emitting part of the semiconductor laser La . L_1 objective lens, adjustable for focusing. L_2 lens for making the divergent laser beam parallel. M half-silvered mirror formed by a film evaporated on the dividing surface of the prism combination P_1 . P_2 beam-splitter prisms. D_1 to D_4 photodiodes whose output currents can be combined in various ways to provide the output signal from the pick-up and also the tracking-error signal and the focusing-error signal. (In practice the prisms P_2 and the photodiodes D_1 to D_4 are rotated by 90° and the reflection at the mirror M does not take place in a radial plane but in a tangential plane.) *b*) A magnified view of the light spot S and its immediate surroundings, with a plan view. It can clearly be seen that the diameter of the spot (about $1 \mu\text{m}$) is larger than the width of the pit ($0.6 \mu\text{m}$).

Table I gives a survey of the successive operations with the associated bit rates, with their names. From the magnitude of the channel bit rate and the scanning speed of 1.25 m/s it follows that the length of a channel bit on the disc is approximately $0.3 \mu\text{m}$.

The signal produced in this way is used by the disc manufacturer to switch on and off the laser beam that illuminates the light-sensitive layer on a rotating glass

Read-out from the disc

As we have seen, the disc is optically scanned in the player. This is done by the AlGaAs semiconductor laser described in an earlier article in this journal [6]. Fig. 2 shows the optical part of the 'pick-up'. The light from the laser La (wavelength 800 nm) is focused through the lenses L_2 and L_1 on to the reflecting layer of the disc. The diameter of the light spot S is about

1 μm . When the spot falls on an interval between two pits, the light is almost totally reflected and reaches the four photodiodes D_1 - D_4 via the half-silvered mirror M . When the spot lands on a pit — the depth of a pit is about $\frac{1}{4}$ of the wavelength in the transparent substrate material — interference causes less light to be reflected and an appreciably smaller amount reaches the photodiodes. When the output signals from the four photodiodes are added together the result is a fairly rough approximation^[3] to the rectangular pulse pattern present on the disc in the form of pits and intervals.

The optical pick-up shown in fig. 2 is very small (about 45×12 mm) and is mounted in a pivoting arm that enables the pick-up to describe a radial arc across the disc, so that it can scan the complete spiral track. Around the pivotal point of the arm is mounted a 'linear' motor that consists of a combination of a coil and a permanent magnet. When the coil is energized the pick-up can be directed to any required part of the track, the locational information being provided by the C&D bits added to each frame on the disc. The pick-up is thus able to find independently any particular passage of music indicated by the listener. When it has been found, the pick-up must then follow the track accurately — to within $\pm 0.1 \mu\text{m}$ — without being affected by the next or previous track. Since the track on the disc may have some slight eccentricity, and since also the suspension of the turntable is not perfect, the track may have a maximum side-to-side swing of $300 \mu\text{m}$. A tracking servosystem is therefore necessary to ensure that the deviation between pick-up and track is smaller than the permitted value of $\pm 0.1 \mu\text{m}$ and in addition to absorb the consequences of small vibrations of the player.

The tracking-error signal is delivered by the four photodiodes D_1 to D_4 . When the spot S , seen in the radial direction, is situated in the centre of the track, a symmetrical beam is reflected. If the spot lies slightly to one side of the track, however, interference effects cause asymmetry in the reflected beam. This asymmetry is detected by the prisms P_2 , which split the beam into two components. Beyond the prisms one component has a higher mean intensity than the other. The signal obtained by coupling the photodiodes as $(D_1 + D_2) - (D_3 + D_4)$ can therefore be used as a tracking-error signal.

As a result of ageing or soiling of the optical system, the reflected beam may acquire a slowly increasing, more or less constant asymmetry. Owing to a d.c. component in the tracking-error signal, the spot will then always be slightly off-centre of the track. To compensate for this effect a second tracking-error signal is generated. The coil that controls the pick-up arm

is therefore supplied with an alternating voltage at 600 Hz, with an amplitude that corresponds to a radial displacement of the spot by $\pm 0.05 \mu\text{m}$. The output sum signal from the four photodiodes — which is at a maximum when the spot is in the centre of the track — is thus modulated by an alternating voltage of 600 Hz. The amplitude of this 600 Hz signal increases as the spot moves off-centre. In addition the sign of the 600 Hz error signal changes if the spot moves to the other side of the track. This second tracking-error signal is therefore used to correct the error signal mentioned earlier with a direct voltage. The output sum signal from the photodiodes, which is processed in the player to become the audio signal, is thus returned to its maximum value.

The depth of focus of the optical pick-up at the position of S (see fig. 2) is about $4 \mu\text{m}$. The axial deviation of the disc, owing to various mechanical effects, can have a maximum of 1 mm. It is evident that a servosystem is also necessary to give correct focusing of the pick-up on the reflecting layer. The objective lens L_1 can therefore be displaced in the direction of its optical axis by a combination of a coil and a permanent magnet, in the same way as in a loud-speaker. The focusing-error signal is also provided by the row of photodiodes D_1 to D_4 . If the spot is sharply focused on the disc, two sharp images are precisely located between D_1 and D_2 and between D_3 and D_4 . If the spot is not sharply focused on the disc, the two images on the photodiodes are not sharp either, and have also moved closer together or further apart. The signal obtained by connecting the photodiodes as $(D_1 + D_4) - (D_2 + D_3)$ can therefore be used for controlling the focusing servosystem. The deviation in focusing then remains limited to $\pm 1 \mu\text{m}$.

Reconstitution of the audio signal

The signal read from the disc by the optical pick-up has to be reconstituted to form the analog audio signal.

Fig. 3 shows the block diagram of the signal processing in the player. In *DEMODO* the demodulation follows the same rules that were applied to the EFM modulation, but now in the opposite sense. The information is then temporarily stored in a buffer memory and then reaches the error-detection and correction circuit *ERCO*. The parity bits can be used here to correct errors, or just to detect errors if correction is found to be impossible^[4]. These errors may originate from defects in the manufacturing process, damage during use, or fingermarks or dust on the disc. Since the information with the CIRC code is 'interleaved' in time, errors that occur at the input of *ERCO* in one

frame are spread over a large number of frames during decoding in *ERCO*. This increases the probability that the maximum number of correctable errors per frame will not be exceeded. A flaw such as a scratch can often produce a train of errors, called an error burst. The error-correction code used in *ERCO* can correct a burst of up to 4000 data bits, largely because the errors are spread out in this way.

ERCO are synchronized by a clock generator *C* controlled by a quartz crystal.

Fig. 3 also illustrates the control of the disc speed n_D . The bit stream leaves the buffer memory at a rate synchronized by the clock generator. The bit stream enters the buffer memory, however, at a rate that depends on the speed of revolution of the disc. The extent to which n_D and the sampling rate are matched

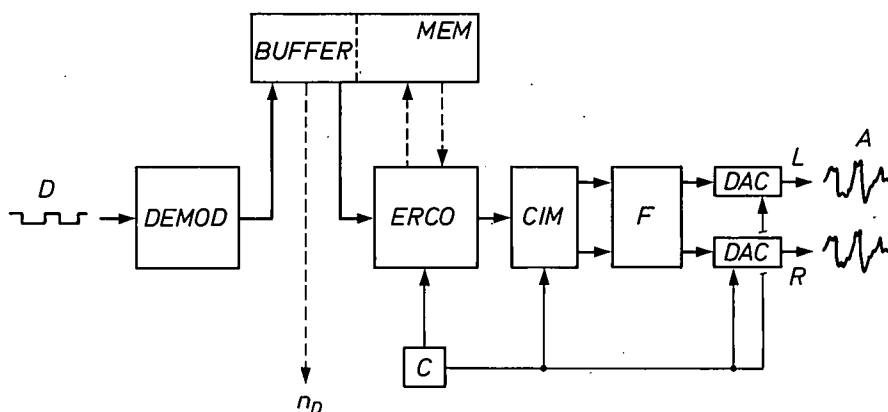


Fig. 3. Block diagram of the signal processing in the player. *D* input signal read by the optical pick-up; see fig. 2. *A* the two output analog audio signals from the left (*L*) and the right (*R*) audio channels. *DEMOD* demodulation circuit. *ERCO* error-correction circuit. *BUFFER* buffer memory, forming part of the main memory *MEM* associated with *ERCO*. *CIM* (Concealment: Interpolation and Muting) circuit in which errors that are only detected since they cannot be corrected are masked or 'concealed'. *F* filters for interpolation. *DAC* digital-to-analog conversion circuits. Each of the blocks mentioned here are fabricated in VLSI technology. *C* clock generator controlled by a quartz crystal. The degree to which the buffer memory capacity is filled serves as a criterion in controlling the speed of the disc.

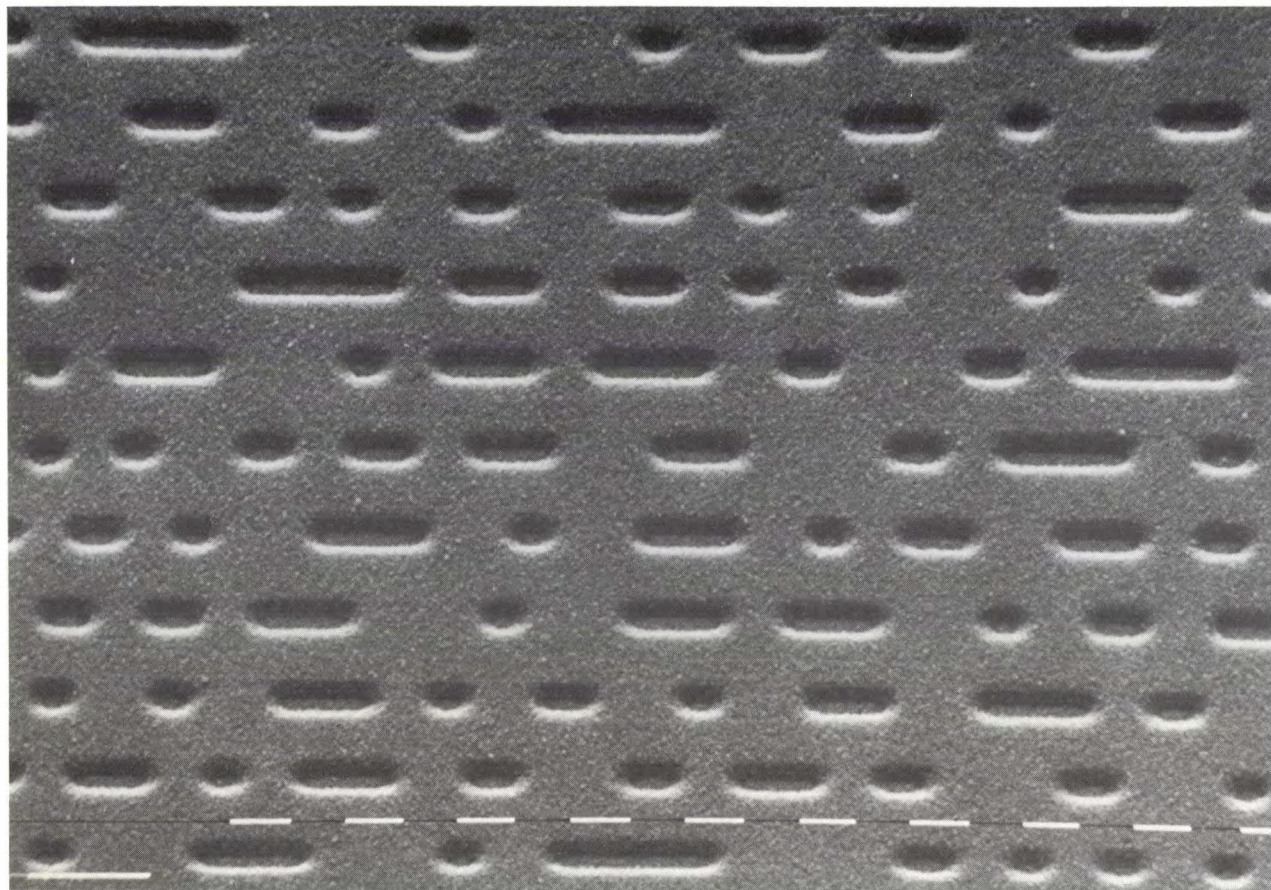
If more errors than the permitted maximum occur, they can only be detected. In the *CIM* block (Concealment: Interpolation and Muting) the errors detected are then masked. If the value of a sample indicates an error, a new value is determined by linear interpolation between the preceding value and the next one. If two or more successive sample values indicate an error, they are made equal to zero (muting). At the same time a gradual transition is created to the values preceding and succeeding it by causing a number of values before the error and after it to decrease to zero in a particular pattern.

In the digital-to-analog converters *DAC*^[7] the 16 bit samples first pass through interpolation filters *F* and are then translated and recombined to recreate the original analog audio signal *A* from the two audio channels *L* and *R*. Since samples must be recombined at exactly the same rate as they are taken from the analog audio signal, the *DACs* and also *CIM* and

determines the 'filling degree' of the buffer memory. The control is so arranged as to ensure that the buffer memory is at all times filled to 50% of its capacity. The analog signal from the player is thus completely free from wow and flutter, yet with only moderate requirements for the speed control of the disc.

[7] See D. Goedhart, R. J. van de Plassche and E. F. Stikvoort, Digital-to-analog conversion in playing a Compact Disc, this issue, p. 174.

Summary. Digital processing of the audio signal and optical scanning in the Compact Disc system yield significant advantages: insensitivity to surface damage of the disc, compactness of disc and player, excellent signal-to-noise ratio and channel separation (both 90 dB) and a flat response over a wide range of frequencies (up to 20000 Hz). The Compact Disc, with a diameter of only 120 mm, gives a continuous playing time of an hour or more. The analog audio signal is converted into a digital signal suitable for transcription on the disc. After the digital signal has been read from the disc by an optical 'pick-up' the original audio signal is recreated in the player.



The information on the Compact Disc is recorded in digital form as a spiral track consisting of a succession of pits. The pitch of the track is $1.6\ \mu\text{m}$, the width $0.6\ \mu\text{m}$ and the depth of the pits $0.12\ \mu\text{m}$. The length of a pit or the land between two pits has a minimum value of $0.9\ \mu\text{m}$ and a maximum value of $3.3\ \mu\text{m}$. The scale at the bottom indicates intervals of $1\ \mu\text{m}$.

Compact Disc: system aspects and modulation

J. P. J. Heemskerk and K. A. Schouhamer Immink

In this article we shall deal in more detail with the various factors that had to be weighed one against the other in the design of the Compact Disc system. In particular we shall discuss the EFM modulation system ('Eight-to-Fourteen Modulation'), which helps to produce the desired high information density on the disc.

tical to B_1 — from the disc and reconverts it to the orchestral sound. The system between *COD* and *DECOD* is the actual *transmission channel*; B_1 and B_0 consist of 'channel bits'.

Fig. 2 shows the encoding system in more detail. The audio signal is first converted into a stream B_1 of 'audio bits' by means of pulse-code modulation. A

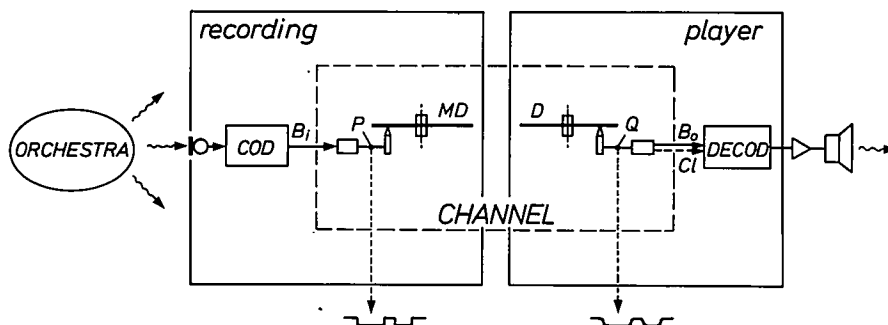


Fig. 1. The Compact Disc system, considered as a transmission system that brings sound from the studio into the living room. The transmission channel between the encoding system (*COD*) at the recording end and the decoding system (*DECOD*) in the player, 'transmits' the bit stream B_1 to *DECOD* via the write laser, the master disc (*MD*), the disc manufacture, the disc (*D*) in the player and the optical pick-up; in the ideal case B_0 is the same as B_1 . The bits of B_0 , as well as the clock signal (*Cl*) for further digital operations, have to be detected from the output signal of the pick-up unit at *Q*.

Fig. 1 represents the complete Compact Disc system as a 'transmission system' that brings the sound of an orchestra into the living room. The orchestral sound is converted at the recording end into a *bit stream* B_1 , which is recorded on the master disc. The master disc is used as the 'pattern' for making the discs for the user. The player in the living room derives the bit stream B_0 — which in the ideal case should be iden-

number of bits for 'control and display' (C&D) and the parity bits for error correction are then added to the bit stream^{[1][2]}. This results in the 'data bit stream' B_2 . The modulator converts this into channel bits (B_3). The bit stream B_1 is obtained by adding a synchronization signal.

Dr J. P. J. Heemskerk is with the Philips Audio Division, Eindhoven; Ir K. A. Schouhamer Immink is with Philips Research Laboratories, Eindhoven.

[1] M. G. Carasso, J. B. H. Peek and J. P. Sinjou, The Compact Disc Digital Audio system, this issue, p. 151.

[2] H. Hoeve, J. Timmermans and L. B. Vries, Error correction and concealment in the Compact Disc system, this issue, p. 166.

The number of data bits n that can be stored on the disc is given by:

$$n = \eta A/d^2,$$

where A is the useful area of the disc surface, d is the diameter of the laser light spot on the disc and η is the 'number of data bits per spot' (the number of data bits that can be resolved per length d of track). A/d^2 is the number of spots that can be accommodated side by side on the disc. The information density n/A is thus given by:

$$n/A = \eta/d^2. \quad (1)$$

The spot diameter d is one of the most important parameters of the channel. The modulation can give a higher value of η . We shall now briefly discuss some of the aspects of the channel that determine the specification for the modulation system.

We shall consider one example here to illustrate the way in which such tolerances affect the design: the choice of the 'spot diameter' d . We define d as the half-value diameter for the light intensity; we have

$$d = 0.6 \lambda/NA,$$

where λ is the wavelength of the laser light and NA is the numerical aperture of the objective. To achieve a high information density (1) d must be as small as possible. The laser chosen for this system is the small CQL10 [31], which is inexpensive and only requires a low voltage; the wavelength is thus fixed; $\lambda \approx 800$ nm. This means that we must make the numerical aperture as large as possible. With increasing NA , however, the manufacturing tolerances of the player and the disc rapidly become smaller. For example, the tolerance in the local 'skew' of the disc (the 'disc tilt') relative to the objective-lens axis is proportional to NA^{-3} . The

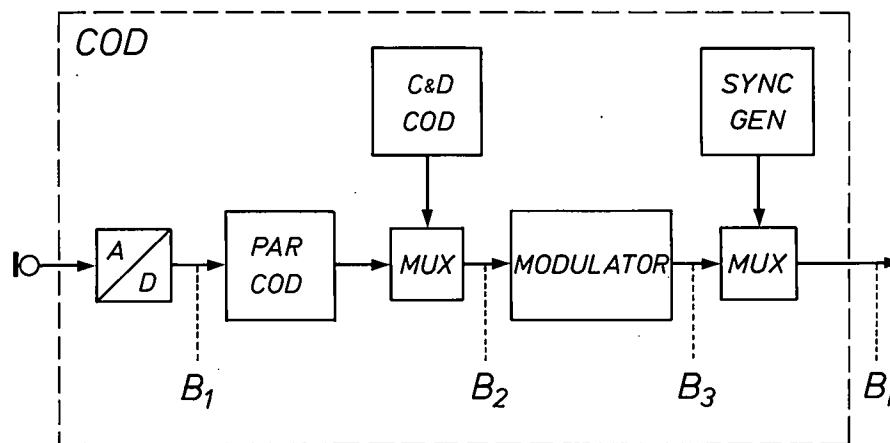


Fig. 2. The encoding system (COD in fig. 1). The system is highly simplified here; in practice for example there are two audio channels for stereo recording at the input, which together supply the bit stream B_1 by means of PCM, and the various digital operations are controlled by a 'clock', which is not shown. The bit stream B_1 is supplemented by parity and C&D (control and display) bits (B_2), modulated (B_3), and provided with synchronization signals (B_i). MUX: multiplexers. Fig. 9 gives the various bit streams in more detail.

The channel

The bit stream B_i in fig. 1 is converted into a signal at P that switches the light beam from the write laser on and off. The channel should be of high enough quality to allow the bit stream B_i to be reconstituted from the read signal at Q .

To achieve this quality all the stages in the transmission path must meet exacting requirements, from the recording on the master disc, through the disc manufacture, to the actual playing of the disc. The quality of the channel is determined by the player and the disc: these are mass-produced and the tolerances cannot be made unacceptably small.

tolerance for the disc thickness is proportional to NA^{-4} , and the depth of focus, which determines the focusing tolerance, is proportional to NA^{-2} . After considering all these factors in relation to one another, we arrived at a value of 0.45 for NA . We thus find a value of $1 \mu\text{m}$ for the spot diameter d .

The quality of the channel is evaluated by means of an 'eye pattern', which is obtained by connecting the point Q in fig. 1 to an oscilloscope synchronized with the clock for the bit stream B_0 ; see fig. 3a. The signals originating from different pits and lands are superimposed on the screen; they are strongly rounded,

mainly because the spot diameter is not zero and the pit walls are not vertical. If the transmission quality is adequate, however, it is always possible to determine whether the signal is positive or negative at the 'clock times' (the dashes in fig. 3a), and hence to reconstitute the bit stream. The lozenge pattern around a dash in this case is called the 'eye'. Owing to channel imperfections the eye can become obscured; owing

important parameters, both for the player and for the disc. The list is far from complete, of course.

With properly manufactured players and discs the channel quality can still be impaired by dirt and scratches forming on the discs during use. By its nature the system is fairly insensitive to these^[1], and any errors they may introduce can nearly always be corrected or masked^[2]. In the following we shall see that the modulation system also helps to reduce the sensitivity to imperfections.

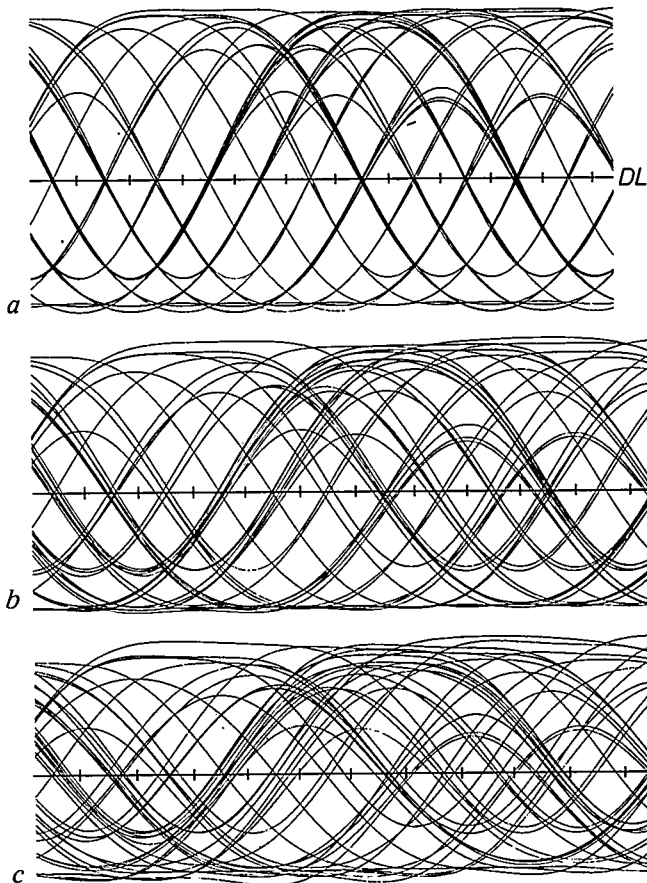


Fig. 3. Eye pattern. The figures give the read signal (at Q in fig. 1) on an oscilloscope synchronized with the bit clock. At the decision times (marked by dashes) it must be possible to determine whether the signal is above or below the decision level (DL). The curves have been calculated for *a*) an ideal optical system, *b*) a defocusing of $2\ \mu\text{m}$, *c*) a defocusing of $2\ \mu\text{m}$ and a disc tilt of 1.2° . The curves give a good picture of experimental results.

to 'phase jitter' of the signal relative to the clock an eye becomes narrower, and noise reduces its height. The signals in fig. 3a were calculated for a perfect optical system. Fig. 3b shows the effect of defocusing by $2\ \mu\text{m}$ and fig. 3c shows the effect of a radial tilt of 1.2° in addition to the defocusing. In fig. 3b a correct decision is still possible, but not in fig. 3c.

This example also gives some idea of the exacting requirements that the equipment has to meet. A more general picture can be obtained from Table I, which gives the manufacturing tolerances of a number of

Table I. Manufacturing tolerances.

Player	Objective-lens tilt $\pm 0.2^\circ$
	Tracking $\pm 0.1\ \mu\text{m}$
	Focusing $\pm 0.5\ \mu\text{m}$
	R.M.S. wavefront noise of read laser beam $0.05\ \lambda$ (40 nm)
Disc	Thickness $1.2 \pm 0.1\ \text{mm}$
	Flatness $\pm 0.6^\circ$ (at the rim corresponding to a sag 0.5 mm)
	Pit-edge positioning $\pm 50\ \text{nm}$
	Pit depth $120 \pm 10\ \text{nm}$

Bit-stream modulation

The playing time of a disc is equal to the track length divided by the track velocity v . For a given disc size the playing time therefore increases if we decrease the track velocity in the system (the track velocity of the master disc and of the user disc). However, if we do this the channel becomes 'worse': the eye height decreases and the system becomes more sensitive to perturbations. There is therefore a lower limit to the track velocity if a minimum value has been established for the eye height because of the expected level of noise and perturbation. We shall now show that we can decrease this lower limit by an appropriate bit-stream modulation.

We first consider the situation without modulation. The incoming data bit stream is an arbitrary sequence of ones and zeros. We consider a group of 8 data bits in which the change of bit value is fastest (fig. 4a). Un-coded recording (1: pit; 0: land, or vice versa) then gives the pattern of fig. 4b. This results in the rounded-off signal of fig. 4c at Q in fig. 1; fig. 4d gives the eye pattern. The signal in fig. 4c represents the highest frequency (f_{m1}) for this mode of transmission, and we have $f_{m1} = \frac{1}{2} f_d$, where f_d is the data bit rate. The half eye height a_1 is equal to the amplitude A_1 of the highest-frequency signal.

[1] 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.

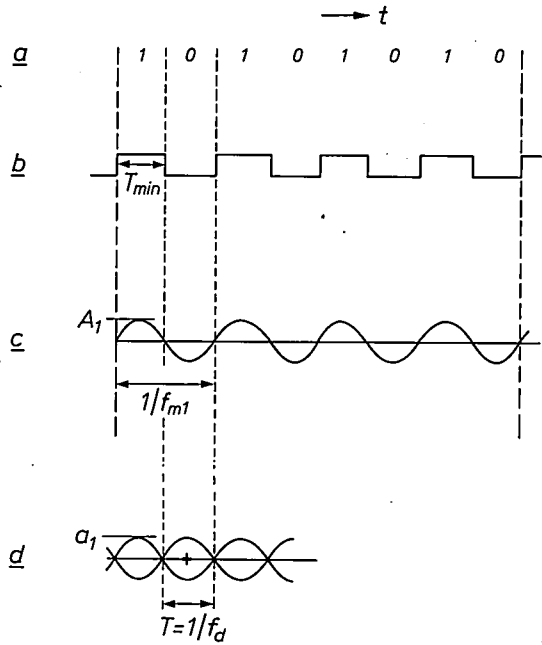


Fig. 4. Direct recording of the data bit stream on the disc. a) Data bit stream of the highest frequency that can occur. b) Direct translation of the bit stream into a pattern of pits. c) The corresponding output signal (at Q in fig. 1); its amplitude A_1 is found with the aid of fig. 5. d) The eye pattern that follows from (c). T_{min} minimum pit or land length; f_{m1} highest frequency; T data bit length; f_d data bit rate. We have $T_{min} = T$; $f_{m1} = \frac{1}{2} f_d$.

The relation between the eye height and the track velocity now follows indirectly from the 'amplitude-frequency characteristic' of the channel; see fig. 5. In this diagram A is the amplitude of the sinusoidal signal at Q in fig. 1 when a sinusoidal unit signal of frequency f is presented at P. With the aid of Fourier

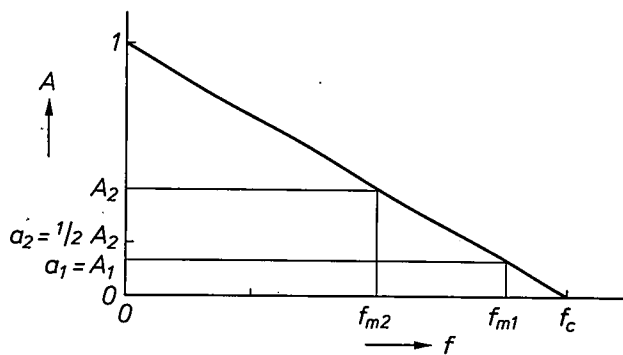


Fig. 5. Amplitude-frequency characteristic of the channel. The diagram gives the amplitude A of the sinusoidal signal at Q (fig. 1) when a sinusoidal unit signal is presented at P as a function of the frequency f . The transfer is 'cut off' at the frequency f_c , which is given by $f_c = (2NA/\lambda)v$. The line shown applies to an ideal optical system; in reality A is always somewhat lower; the cut-off frequency is then effectively lower. The 'maximum frequencies' f_{m1} , f_{m2} , the amplitudes A_1 , A_2 and the 'half eye heights' a_1 , a_2 relate to the 'direct' and 'modulated' writing of the data bits on the disc; see figs 4 and 6.

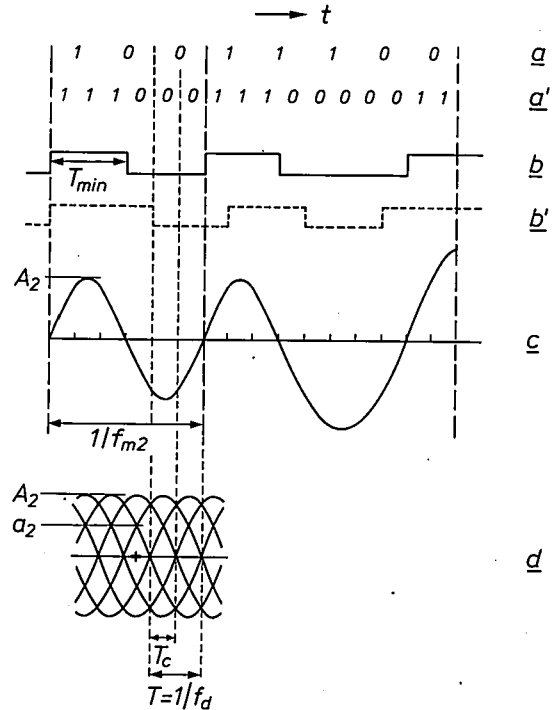


Fig. 6. Eight-to-sixteen modulation. Each group of 8 data bits (a) is translated with the aid of a dictionary into 16 channel bits (a'), in such a way that the run length is equal to at least three channel bits. b) Pattern of pits produced from the bit stream (a'). b') pattern of pits obtained with a different input signal. c) The read signal corresponding to (b); its amplitude is again determined from fig. 5. d) The resultant eye pattern. The half eye height (a_2) here is only half the amplitude (A_2) of the approximately sinusoidal signal of maximum frequency (f_{m2}).

analysis and synthesis the output signal can be calculated from $A(f)$ for any input signal. The line in the diagram represents a channel with a perfect optical system. In the first part of this section we shall take this for granted. The true situation will always be less favourable. The 'cut-off frequency' is determined by the spot diameter and the track velocity v ; in the ideal case $f_c = (2NA/\lambda)v$.

For a given track velocity we now obtain the half eye height a_1 in fig. 4 directly from fig. 5: it is equal to the amplitude A_1 at the frequency f_{m1} . If v , and hence f_c , is varied, the line in fig. 5 rotates about the point 1 on the A -axis. For a given minimum value of a_1 , the figure indicates how far f_c can be decreased; this establishes the lower limit for v . In particular, if the minimum value for a_1 is very small, f_c can be decreased to a value slightly above $f_{m1} (= \frac{1}{2} f_d)$.

Fig. 6 gives the situation with modulation: an imaginary 8→16 modulation, which is very close to EFM, however. Each group of 8 incoming data bits (fig. 6a) is converted into 16 channel bits (fig. 6a'). This is done by using a 'dictionary' that assigns unambiguously but otherwise arbitrarily to each word of 8 bits a word of 16 bits, but in such a way that the resultant channel bit stream only produces pits and lands that

are at least three channel bits long (fig. 6*b*). On the time scale the minimum pit and land lengths ('the minimum run length' T_{\min}) have become $1\frac{1}{2}$ times as long as in fig. 4, but a simple calculation shows that about as much information can nevertheless be transmitted as in fig. 4 (256 combinations for 8 data bits), because there is a greater choice of pit-edge positions per unit length (see fig. 6*b* and *b'*); the 'channel bit length' T_c has decreased by a half.

With the modulation we have managed to reduce the highest frequency (f_{m2}) in the signal (see fig. 6*c*, left; $f_{m2} = \frac{1}{3}f_d = \frac{2}{3}f_{m1}$). Therefore f_c and v can be reduced by a factor of $1\frac{1}{2}$ for the case in which a very small eye height is tolerable (see fig. 5); this represents an increase of 50% in playing time.

The modulation also has its disadvantages. In the first place the half eye height (a_2) in this case is only half of the amplitude (A_2) of the signal at the highest frequency (see fig. 6*d*). This has consequences if the minimum eye height is not very small. For example, the modulation becomes completely unusable if the half eye height in fig. 5 has to remain larger than $\frac{1}{2}$ ($a_2 > \frac{1}{2}$ implies $A_2 > 1$); uncoded recording is then still possible ($A_1 = a_1$). In the second place, the tolerance for time errors and for the positioning of pit edges, together with the eye width (T_c), has decreased by a half. In designing a system, the various factors have to be carefully weighed against one another.

To show qualitatively how a choice can be made, we have plotted the half eye height in fig. 7 as a function of the 'linear information density' σ (the number of incoming data bits per unit length of the track; $\sigma = f_d/v$) for three systems: '8→8 modulation' (i.e. uncoded recording), 8→16 modulation, and a system that also has about the same information capacity (256 combinations for 8 data bits) in which, however, the minimum run length has been increased still further, again at the expense of eye width of course ('8→24 modulation', $T_{\min} = 2T$, $T_c = \frac{1}{3}T$). The figure is a direct consequence of the reasoning above, with the assumption that the cut-off frequency is 20% lower than the ideal value $(2NA/\lambda)v$, as a first rough adjustment to what we find in practice for the function $A(f)$.

In qualitative terms, the 8→16 system has been chosen because the nature of the noise and perturbations is such that the eye can be smaller than at *A* in fig. 7, but becomes too small at *C*. An improvement is therefore possible with 8→16 modulation, but not with 8→24 modulation.

For our Compact Disc system we have $\sigma = 1.55$ data bits/ μm ($f_d = 1.94$ Mb/s, $v = 1.25$ m/s^[1]); the operating point would therefore be at *P* in fig. 7. The model used is however rather crude and in better models *A*, *B* and *C* lie more to the left, so that *P* ap-

proaches *C*. But 8→16 modulation is still preferable to 8→24 modulation, even close to *C*, since the eye width is $1\frac{1}{2}$ times as large as for 8→24 modulation.

EFM is a refinement of 8→16 modulation. It has been chosen on the basis of more detailed models and many experiments. At the eye height used, it gives a gain of 25% in information density, compared with uncoded recording.

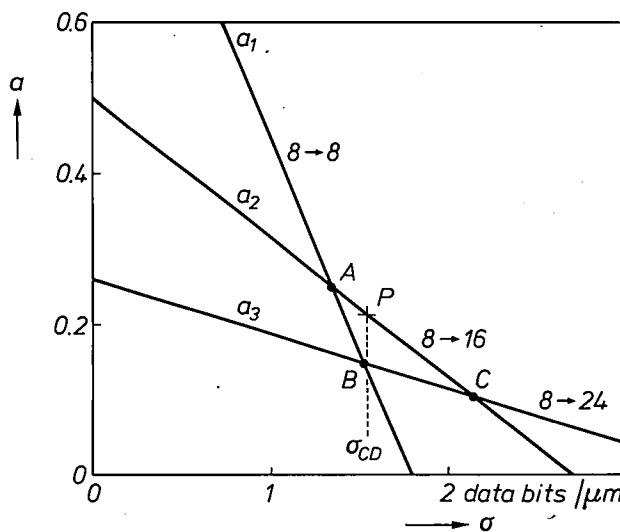


Fig. 7. Half eye height a as a function of the linear information density σ , for 8→8, 8→16 and 8→24 modulation. These systems are characterized by the following values for the channel bit length T_c and the minimum run length T_{\min} :
8→8: $T_c = T$, $T_{\min} = T$ (fig. 4),
8→16: $T_c = \frac{1}{2}T$, $T_{\min} = \frac{3}{2}T$ (fig. 6),
8→24: $T_c = \frac{1}{3}T$, $T_{\min} = 2T$,
where T is the data bit length. The straight lines give the relations that follow from fig. 5:

$$a_1 = c_1(1 - f_{m1}/f_c) \rightarrow a_1 = 1 - \sigma/1.8,$$

$$a_2 = c_2(1 - f_{m2}/f_c) \rightarrow a_2 = 0.5(1 - \sigma/2.7),$$

$$a_3 = c_3(1 - f_{m3}/f_c) \rightarrow a_3 = 0.26(1 - \sigma/3.6),$$

where σ is the numerical value of the linear information density, expressed in data bits per μm . The c 's are the ratios of the half eye height to the amplitude, and the f_m 's the maximum frequencies for the three systems ($c_1 = 1$, $c_2 = \sin 30^\circ = 0.5$, $c_3 = \sin 15^\circ = 0.26$, $f_{m1} = \frac{1}{2}f_d$, $f_{m2} = \frac{1}{3}f_d$, $f_{m3} = \frac{1}{4}f_d$; f_d is the data bit rate). The second set of equations follows from the first set by substituting $0.8 \times (2NA/\lambda)v$ for f_c , with $NA = 0.45$, $\lambda = 0.8 \mu\text{m}$, $v = f_d/\sigma$. The factor 0.8 is introduced as a rough first-order correction to the 'ideal' amplitude characteristic.

Further requirements for the modulation system

In developing the modulation system further we still had two more requirements to take into account.

In the first place it must be possible to regenerate the *bit clock* in the player from the read-out signal (the signal at *Q* in fig. 1). To permit this the number of pit edges per second must be sufficiently large, and in particular the 'maximum run length' T_{\max} must be as small as possible.

The second requirement relates to the 'low-frequency content' of the read signal. This has to be as

small as possible. There are two reasons for this. In the first place, the servosystems for track following and focusing^[1] are controlled by low-frequency signals, so that low-frequency components of the information signal could interfere with the servosystems. The second reason is illustrated in *fig. 8*, in which the read signal is shown for a clean disc (*a*) and for a disc that has been soiled, e.g. by fingermarks (*b*). This causes the amplitude and average level of the signal to fall. The fall in level causes a completely

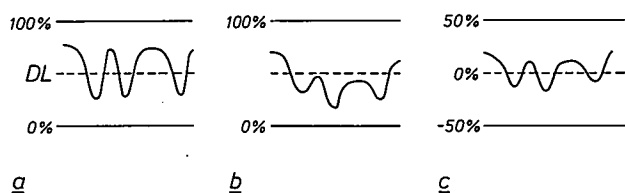


Fig. 8. The read-out signal for six pit edges on the disc, *a*) for a clean disc, *b*) for a soiled disc, *c*) for a soiled disc after the low frequencies have been filtered out. *DL* decision level. Because of the soiling, both the amplitude and the signal level decrease; the decision errors that this would cause are eliminated by the filter.

wrong read-out if the signal falls below the decision level. Errors of this type are avoided by eliminating the low-frequency components with a filter (*c*), but the use of such a filter is only permissible provided the information signal itself contains no low-frequency components. In the Compact Disc system the frequency range from 20 kHz to 1.5 MHz is used for information transmission; the servosystems operate on signals in the range 0-20 kHz.

The EFM modulation system

Fig. 9 gives a schematic general picture of the bit streams in the encoding system. The information is divided into 'frames'. One frame contains 6 sampling periods, each of 32 audio bits (16 bits for each of the two audio channels). These are divided into symbols of 8 bits. The bit stream B_1 thus contains 24 symbols per frame. In B_2 eight parity symbols have been added and one C&D symbol, resulting in 33 'data symbols'. The modulator translates each symbol into a new symbol of 14 bits. Added to these are three 'merging bits', for reasons that will appear shortly. After the addition of a synchronization symbol of 27 bits to the frame, the bit stream B_i is obtained. B_i therefore contains $33 \times 17 + 27 = 588$ channel bits per frame. Finally, B_i is converted into a control signal for the write laser. It should be noted that in B_i '1' or '0' does not mean 'pit' or 'land', as we assumed for simplicity in *fig. 6*, but a '1' indicates a pit edge. The informa-

tion is thus completely recorded by the positions of the pit edges; it therefore makes no difference to the decoding system if 'pit' and 'land' are interchanged on the disc.

Opting for the translation of series of 8 bits following the division into symbols in the parity coding has the effect of avoiding error propagation. This is because in the error-correction system an entire symbol is always either 'wrong' or 'not wrong'. One channel-bit error that occurs in the transmission spoils an entire symbol, but — because of the correspondence between modulation symbols and data symbols — never more than one symbol. If a different modulation system is used, in which the data bits are not translated in groups of 8, but in groups of 6 or 10, say, then the bit stream B_2 is in fact first divided up into 6 or 10 bit 'modulation symbols'. Although one channel-bit error then spoils only one modulation symbol, it usually spoils two of the original 8 bit symbols.

In EFM the data bits are translated 8 at a time into 14 channel bits, with a T_{\min} of 3 and a T_{\max} of 11 channel bits (this means at least 2 and at the most 10 successive zeros in B_i). This choice came about more or less as follows. We have already seen that the choice of about $1\frac{1}{2}$ data bits for T_{\min} , with about 16 channel bits on 8 data bits, is about the optimum for the Compact Disc system^[4]. A simple calculation shows that at least 14 channel bits are necessary for the reproduction of all the 256 possible symbols of 8 data bits under the conditions $T_{\min} = 3$, $T_{\max} = 11$ channel bits. The choice of T_{\max} was dictated by the fact that a larger choice does not make things very much easier, whereas a smaller choice does create far more difficulties.

With 14 channel bits it is possible to make up 267 symbols that satisfy the run-length conditions. Since we only require 256, we omitted 10 that would have introduced difficulties with the 'merging' of symbols under these conditions, and one other chosen at random. The dictionary was compiled with the aid of computer optimization in such a way that the translation in the player can be carried out with the simplest possible circuit, i.e. a circuit that contains the minimum of logic gates.

The merging bits are primarily intended to ensure that the run-length conditions continue to be satisfied when the symbols are 'merged'. If the run length is in danger of becoming too short we choose '0's for the merging bits; if it is too long we choose a '1' for one of them. If we do this we still retain a large measure of freedom in the choice of the merging bits, and we use this freedom to minimize the low-frequency content of the signal. In itself, two merging bits would be sufficient for continuing to satisfy the run-length con-

ditions. A third is necessary, however, to give sufficient freedom for effective suppression of the low-frequency content, even though it means a loss of 6% of the information density on the disc. The merging bits

are shown two data symbols of B_2 and their translation from the dictionary into channel symbols (B_3). From the T_{min} rule the first of the merging bits in this case must be a zero; this position is marked 'X'. In

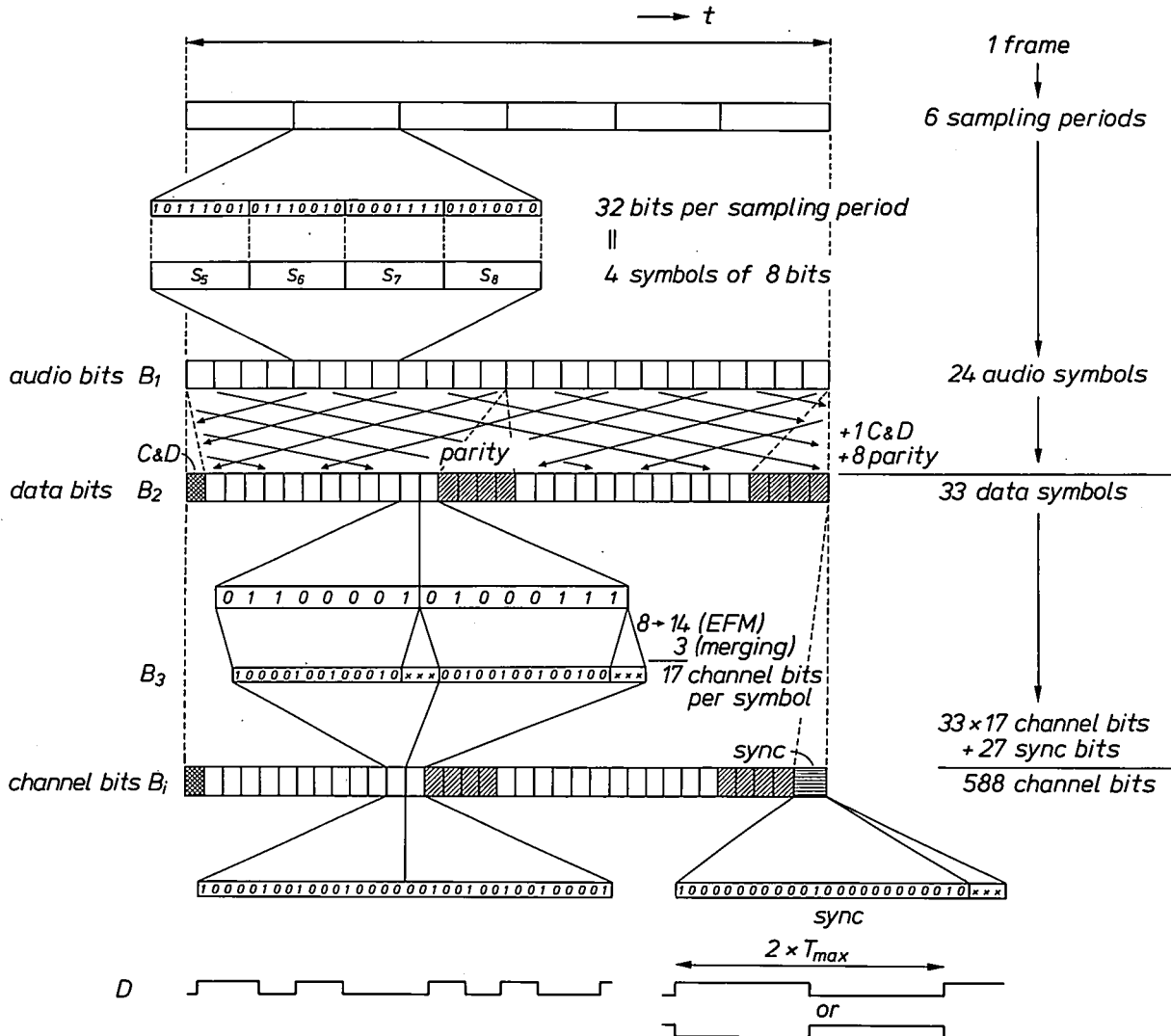


Fig. 9. Bit streams in the encoding system (fig. 2). The information is divided into frames; the figure gives one frame of the successive bit streams. There are six sampling periods for one frame, each sampling period giving 32 bits (16 for each of the two audio channels). These 32 bits are divided to make four symbols in the 'audio bit stream' B_1 . In the 'data bit stream' B_2 eight parity and one C&D symbols have been added to the 24 audio symbols. To scatter possible errors, the symbols of different frames in B_1 are interleaved, so that the audio signals in one frame of B_2 originate from different frames in B_1 . The modulation translates the eight data bits of a symbol of B_2 into fourteen channel bits, to which three 'merging bits' are added (B_3). The frames are marked with a synchronization signal of the form illustrated (bottom right); the final result is the 'channel bit stream' (B_i) used for writing on the master disc, in such a way that each '1' indicates a pit edge (D).

contain no audio information, and they are removed from the bit stream in the demodulator.

Fig. 10 illustrates, finally, how the merging bits are determined. Our measure of the low-frequency content is the 'digital sum value' (DSV); this is the difference between the totals of pit and land lengths accumulated from the beginning of the disc. At the top

the two following positions the choice is free; these are marked 'M'. The three possible choices $XMM = 000, 010$ and 001 would give rise to the patterns of pits as illustrated, and to the indicated waveform of the

[4] A more detailed discussion is given in K. A. Immink, Modulation systems for digital audio discs with optical readout, Proc. IEEE Int. Conf. on Acoustics, speech and signal processing, Atlanta 1981, pp. 587-589.

DSV, on the assumption that the DSV was equal to 0 at the beginning. The system now opts for the merging combination that makes the DSV at the end of the second symbol as small as possible, i.e. 000 in this case. If the initial value had been -3 , the merging combination 001 would have been chosen.

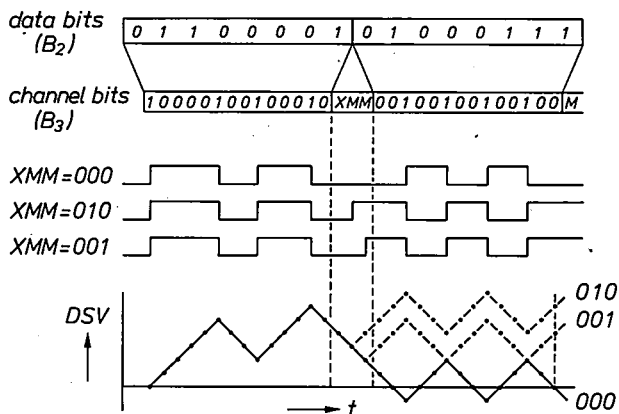


Fig. 10. Strategy for minimizing the digital sum value (DSV). After translation of the data bits into channel bits, the symbols are merged together by means of three extra bits in such a way that the run-length conditions continue to be satisfied and the DSV remains as small as possible. The first run-length rule (at least two zeros one after the other) requires a zero at the first position in the case illustrated here, while the choice remains free for the second and third positions. In this case there are thus three merging alternatives: 000, 010 and 001. These alternatives give the patterns of pits shown in the diagram and the illustrated DSV waveform. The system chooses the alternative that gives the lowest value of DSV at the end of the next symbol. The system looks 'one symbol ahead'; strategies for looking further ahead are also possible in principle.

When this strategy is applied, the noise in the servo-band frequencies (< 20 kHz) is suppressed by about 10 dB. In principle better results can be obtained, within the agreed standard for the Compact Disc system, by looking more than one symbol ahead, since minimization of the DSV in the short term does not always contribute to longer-term minimization. This is not yet done in the present equipment.

Summary. The Compact Disc system can be considered as a transmission system that brings sound from the studio into the living room. The sound encoded into data bits and modulated into channel bits is sent along the 'transmission channel' consisting of write laser — master disc — user disc — optical pick-up. The maximum information density on the disc is determined by the diameter d of the laser light spot on the disc and the 'number of data bits per light spot'. The effect of making d smaller is to greatly reduce the manufacturing tolerances for the player and the disc. The compromise adopted is $d \approx 1 \mu\text{m}$, giving very small tolerances for objective and disc tilt, disc thickness and defocusing. The basic idea of the modulation is that, while maintaining the minimum length for 'pit' and 'land' (the 'minimum run length') required for satisfactory transmission, the information density can be increased by increasing the number of possible positions per unit length for pit edges (the bit density). Because of clock regeneration there is also a maximum run length, and the low-frequency content of the transmission channel must be kept as low as possible. With the EFM modulation system used each 'symbol' of eight data bits is converted into 14 channel bits with a minimum run length of 3 and a maximum run length of 11 bits, plus three merging bits, chosen such that, when the symbols are merged together, the run-length conditions continue to be satisfied and the low-frequency content is kept to the minimum.



This prototype player, which will be put on the market later, will display 'information for the listener' such as title, composer, 'track number' and playing time of the piece of music. The different sections of the music on the disc can also be played in the order selected by the user — the numbers on the far right.

Error correction and concealment in the Compact Disc system

H. Hoeve, J. Timmermans and L. B. Vries

Introduction

When analog signals such as audio signals are transmitted and recorded via an intervening system such as a gramophone record it is difficult to properly correct signal errors that have occurred in the path between the audio source and the receiving end. With suitably coded digital signals, however, a practical means of error correction does exist. We shall demonstrate this with the following example^[1].

Suppose that a message of 12 binary units (bits) has to be transmitted (a stream of digital information can always be divided into groups of a particular size for transmission). The 12 bits x_{ij} are arranged as follows in a matrix, in which all x_{ij} can only have the value 0 or 1:

$$\begin{array}{cccc} x_{11} & x_{12} & x_{13} & x_{14} \\ x_{21} & x_{22} & x_{23} & x_{24} \\ x_{31} & x_{32} & x_{33} & x_{34} \end{array}$$

To discover at the receiving end whether the message read there contains an error, and, if so, what the error is, one extra bit (called a 'parity bit') is added to each row and column: x_{15} , x_{25} , x_{35} and x_{41} , x_{42} , x_{43} , x_{44} respectively. These parity bits provide a check on the correctness of the message received. The values assigned to them are such that x_{i5} ($i = 1, 2, 3$) makes the number of ones in row i even, for example, while x_{4j} ($j = 1, 2, 3, 4$) makes the number of ones in column j even. Next, a further parity bit (x_{45}) is added that has a value such that the number of ones in the block is made even. This results in the following matrix of four rows and five columns:

$$\begin{array}{ccccc} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} \end{array}$$

It is easy to verify that the number of ones in the last row is also even, and so is the number of ones in

the last column. If now a bit, say x_{23} , is incorrectly read at the receiving end, then the number of ones in the second row and the number of ones in the third column will no longer be even, and once this has been ascertained, a 0 at position x_{23} can be changed into a 1, or vice versa, thus correcting the error.

So as to be able in this way to correct one error in 12 information bits, it is necessary to send a total of 20 bits instead of 12: the 'code word' of $n = 20$ bits consists of $k = 12$ information bits and $n - k = 8$ parity bits. The (n, k) code used here, a $(20, 12)$ code, makes it possible to correct single errors and also, as can easily be verified, to detect various multiple bit errors.

The 'rate' of an error-correcting code is taken to be the ratio of the number of information bits to the total number of bits per code word: k/n . The $(20, 12)$ code does not have a high rate, because it requires a relatively large number of parity bits. For the Compact Disc this would entail a considerable reduction in the playing time.

The theory of error-correcting codes^[2] gives design methods that entail a minimal addition of parity bits when certain correction criteria are satisfied. An important concept in this theory is the 'distance' and in particular the 'minimum distance' d_m between two code words of n bits. Distance here is taken to be the number of places in which the bits of the two code words differ from each other. In the above example the minimum distance d_m is equal to 4: if one single bit of the k information bits changes, then the two parity bits of the associated row and column change at the same time, as does the one at the bottom right-hand corner, x_{45} , so that the entire code word has changed at four places. Theory tells us that to correct all the combinations of t errors occurring within one word, the minimum distance must be at least $2t + 1$. To correct single errors, therefore, the minimum distance need be no greater than three. Examples of this are the single-error-correcting Hamming codes^[4].

Ir H. Hoeve and J. Timmermans are with the Philips Audio Division, Eindhoven; Ir L. B. Vries is with Philips Research Laboratories, Eindhoven.

The statement that a code word x , which is received as z because of t errors, can be restored to its original form if the minimum distance is $2t + 1$, can be seen from *fig. 1*. A decoder provided with a list in which all the code words are stored can compare z with each of these code words and thus recover the correct code word unambiguously.

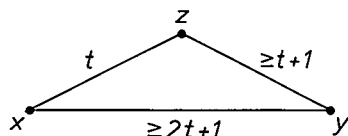


Fig. 1. The original transmitted code word x is received as z owing to t bit errors. Any code word y differing from x lies at a distance $\geq 2t + 1$ from x . To cause z to change into y it is necessary to change at least $t + 1$ bits. It follows that x is the only code word that has a distance t from z .

On the theory of block codes

In the foregoing we have shown with a simple example that it is possible to correct errors. Error-correcting systems do have their limitations, of course. To make this clear we shall consider how error-correcting codes should be designed to guarantee a specific measure of correction, with as few extra bits as possible added to the digital information to be transmitted. It will help if we first say something about the theory of block codes.

So that known and efficient error-correcting codes can be applied, groups of bits are formed by adding together a fixed number s of consecutive bits; these groups are called symbols. With these symbols we now set to work in the same way as with the bits in the foregoing: the information symbols are grouped together to form blocks with a length of k symbols. For error-correction we now add parity symbols to expand each block of k information symbols into a code word of n symbols. The $n - k$ parity symbols to be added are calculated from the k information symbols, and this is done in such a way as to make the error correction as effective as possible. Thus, of the very large number of possibly different words of n symbols only a small fraction, i.e. $2^{(k-n)s}$, become code words (see *fig. 2*). For a given encoding system both n and k are fixed.

As already mentioned in the article on modulation in the Compact Disc system [3], the start of each word is marked by a synchronization symbol. (A word marked by a synchronization symbol is called a 'frame'.) The error-correcting system therefore knows when a new word begins, and the only errors it has to

deal with are errors that occur in the transmission of data.

There are two kinds of errors: those that are distributed at random among the individual bits, the random errors, and errors that occur in groups that may cover a whole symbol or a number of adjacent symbols; these are called 'bursts' of errors. They can occur on a disc as a result of dirt or scratches, which interfere with the read-out of a number of adjacent pits and lands.

The best code for correcting random errors is the one that, for given values of n and k , is able to correct the largest number of independent errors within one code word. In the detection and correction of errors the symbols have to undergo a wide variety of operations. Large k -values (as with the Compact Disc) require extremely complex computing hardware. Practice has shown that the only acceptable solution to this problem is to choose a convenient code. And the only usable codes that enter into consideration, so far as we know at present, are the 'linear codes'.

A code is linear if it obeys the following rule:

If $x = (x_1, \dots, x_n)$ and $y = (y_1, \dots, y_n)$ are code words, then their sum $x + y = (x_1 + y_1, \dots, x_n + y_n)$ is also a code word.

In this sum the symbol $x_i + y_i$ is produced — irrespective of the number of bits s per symbol — by a modulo-2 bit addition. The special feature of the linear code is thus that each sum of code words yields another code word, i.e. a word of n symbols, which also belongs to the small fraction of symbol combinations permitted in the code.

It is this linearity feature that makes it possible to cut down considerably on the extent of the decoding equipment. The *Reed-Solomon codes* [2] are examples of such a linear code. They are also extremely efficient,

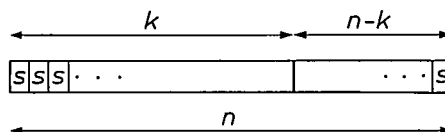


Fig. 2. A code word of length n consists of an information block of k symbols and a parity block of $n - k$ symbols; each symbol comprises s bits. The number of possible words of n symbols is 2^{ns} . The parity bits are fixed for each combination of the ks information bits in accordance with established encoding rules. The number of code words is thus 2^{ks} . It follows that the fraction $2^{(k-n)s}$ of the number of possible words consists of code words.

[1] This example is taken from S. Lin, An introduction to error-correcting codes, Prentice-Hall, Englewood Cliffs 1970.
 [2] See for example F. J. MacWilliams and N. J. A. Sloane, The theory of error-correcting codes, North-Holland, Amsterdam 1978.
 [3] See J. P. J. Heemskerck and K. A. Schouhamer Immink, Compact Disc: system aspects and modulation, this issue, p. 157.

since for every $s > 1$ and $n \leq 2^s - 1$ there exists a Reed-Solomon code with

$$d_m = n - k + 1.$$

Together with the general condition $d_m \geq 2t + 1$ mentioned earlier, which the minimum distance must satisfy for the correction of t errors, this yields $n - k \geq 2t$. Put in another way: to correct t symbol errors it is sufficient to add $2t$ parity symbols. (By 'distance' between two words we mean here the number of positions in which there are different symbols in the two words; it does not matter how many corresponding bits differ from each other within the corresponding symbols.)

In practice a less cumbersome algorithm will generally be used for error correction than the comparison with the aid of a list of all the code words, as described on page 167. We shall not consider the details of the algorithm here. We shall, however, try to give some idea of the manner in which error bursts are tackled with block codes. To do this we must introduce the concept of 'erasure'.

The position (i) of a particular symbol (x_i) in a transmitted code word (x) is called an erasure position if a decoder-independent device signals that the value of x_i is not reliable. This value is then erased, and in the decoding procedure the correct value has to be calculated. The decoding is now simpler and quicker because the positions at which errors can occur are known. (We assume for the moment that no errors occur outside the erasure positions.) The advantage of correcting by means of the erasures is expressed quantitatively by the following proposition:

If a code has a minimum distance d_m , then $d_m - 1$ erasures can be reconstituted.

Since the number of errors that can be corrected without erasure information is $\frac{1}{2}(d_m - 1)$ at most, the advantage of correcting by means of erasures is clear.

The proposition that for a minimum distance d_m , the number of erasures that can be reconstituted is $d_m - 1$, can be proved as follows:

Let x be the code word transmitted and z the word received. Let z be subject to a maximum of $d_m - 1$ erasures, so that it differs from x at a maximum of $d_m - 1$ positions, i.e. for the distance $d(x, z)$ between the code words z and x we have $d(x, z) \leq d_m - 1$. We now replace the symbols at all the erasure positions in z by other symbols, which gives a number of words that we denote by \bar{z} ; we try all the possible substitutions \bar{z} . These words also differ from x at a maximum of $d_m - 1$ positions, i.e. $d(x, \bar{z}) \leq d_m - 1$. Since all the code words different from x have a distance to x that is greater than or equal to the minimum distance d_m , the words \bar{z} include no other code word except x itself. It is therefore only necessary to find out which of the finite number of words \bar{z} represents a code word.

In the Compact Disc system the value of the analog signal to be reproduced is converted at every sampling instant into a binary number of 16 bits per audio channel. For error correction the digital information to be transmitted is divided into groups of eight bits, so that in each sampling operation four information symbols (consisting of audio bits) are generated. In fact, eight parity symbols are added to each block of 24 audio symbols^[4]. The calculation of the parity symbols will not be dealt with here.

Cross-Interleaved Reed-Solomon Code

The error-correcting code used in the Compact Disc system employs not one but two Reed-Solomon codes (C_1, C_2), which are interleaved 'crosswise' (Cross-Interleaved Reed-Solomon Code, CIRC). For code C_1 we have: $n_1 = 32$, $k_1 = 28$, $s = 8$, and for C_2 : $n_2 = 28$, $k_2 = 24$, $s = 8$. The rate of the CIRC we use is $(k_1/n_1)(k_2/n_2) = 3/4$.

For both C_1 and C_2 we have $2t = n - k = 4$, so that for each the minimum distance d_m is equal to $2t + 1 = 5$. This makes it possible to directly correct a maximum of two ($= t$) errors in one code word or to make a maximum of four ($= d_m - 1$) erasure corrections. A combination of both correction methods can also be used.

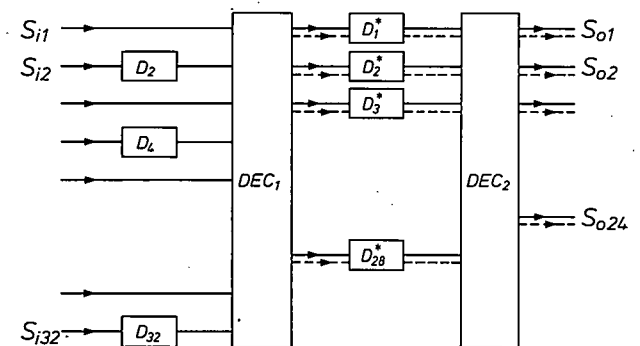


Fig. 3. Schematic representation of the decoding circuit for CIRC. The 32 symbols (s_{11}, \dots, s_{132}) of a frame (24 audio symbols and 8 parity symbols) are applied in parallel to the 32 inputs. The delay lines D_{2i} ($i = 1, \dots, 16$) have a delay equal to the duration of one symbol, so that the information of the 'even' symbols of a frame is cross-interleaved with that of the 'odd' symbols of the next frame. The decoder DEC_1 is designed in accordance with the encoding rules for a Reed-Solomon code with $n_1 = 32$, $k_1 = 28$, $s = 8$. It corrects one error, and if multiple errors occur passes them on unchanged, attaching to all 28 symbols an erasure flag, sent via the dashed lines. Owing to the different lengths of the delay lines D_j^* ($j = 1, \dots, 28$), errors that occur in one word at the output of DEC_1 are 'spread' over a number of words at the input of DEC_2 . This has the effect of reducing the number of errors per DEC_2 word. The decoder DEC_2 is designed in accordance with the encoding rules for a Reed-Solomon code with $n_2 = 28$, $k_2 = 24$, $s = 8$. It can correct a maximum of four errors by means of the erasure-positions method. If there are more than four errors per word, 24 symbol values are passed on unchanged, and the associated positions are given an erasure flag via the dashed lines. S_{01}, \dots, S_{024} outgoing symbols.

Decoding circuit

The error-correction circuit ^[5] is shown schematically in *fig. 3*; *fig. 4* is a photograph of the actual IC. The circuit consists of two decoders, *DEC*, and a number of delay lines, *D* and *D**. The input signal is a sequence of frames ^[6].

pass on 28 symbols unchanged. *DEC*₁ is designed for correcting one error. If it receives a word with a double or triple error, that event is detected with certainty; all the symbols of the received word are passed on unchanged, and all 28 positions are provided with an erasure flag. The same happens in principle for

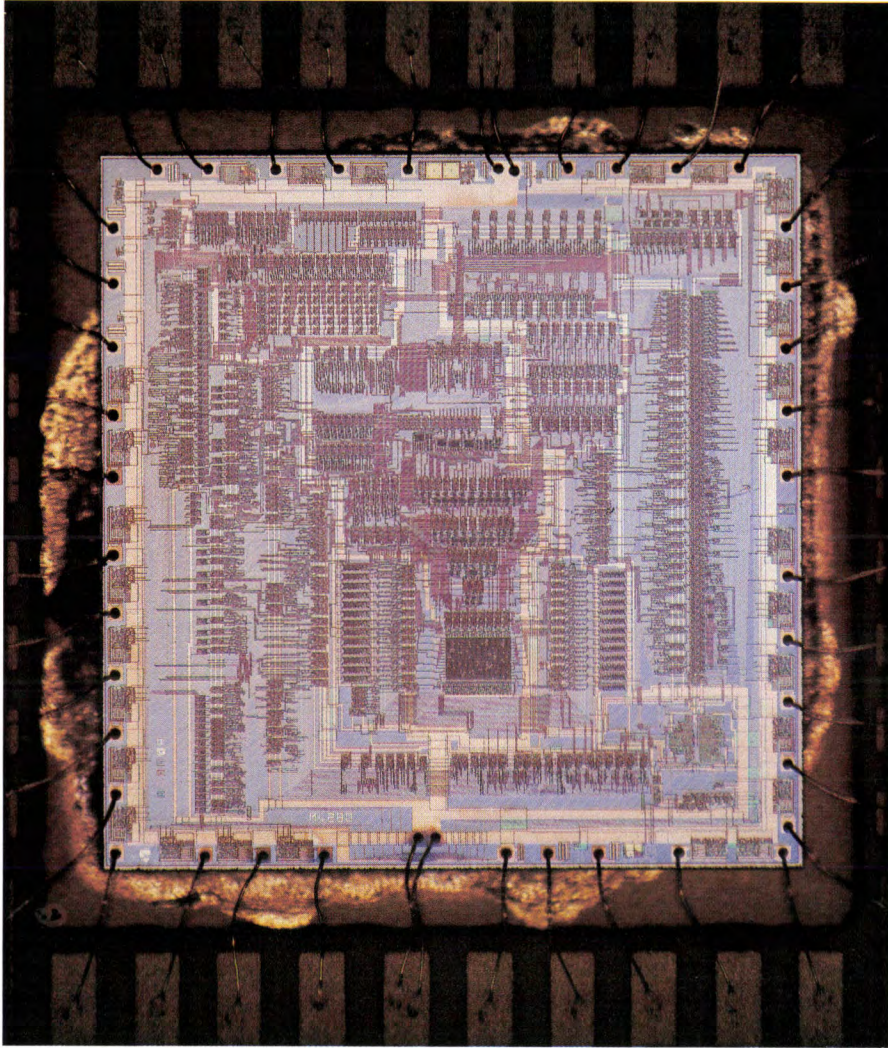


Fig. 4. The integrated circuit for error detection and correction is fabricated in n-channel MOS silicon-gate technology. It has an area of 45 mm² and contains about 12 000 gates.

The 32 symbols of a frame are applied in parallel to the 32 inputs. In passing through the delay lines *D*₂, *D*₄, . . . , *D*₃₂, each of length equal to the duration of one symbol, the even symbols of a frame with the odd symbols of the next frame form the words that are fed to the decoder *DEC*₁. (The symbols of the frames are 'cross-interleaved'. In fact they are 'de-interleaved', because the 'interleaving' ^[4] has already taken place, before the information was recorded on the disc.) If there are no errors in the transmission path, the decoder *DEC*₁ will receive code words that correspond to the encoding rules for *C*₁, and it will

events from 4 to 32 errors, but here there is a small probability ($\approx 2^{-19}$) that this detection will fail. We shall return to this probability later.

The symbols arrive via the delay lines *D*₁^{*}, . . . , *D*₂₈^{*}, which differ from each other in length, at the input of

^[4] The calculation and addition of the parity symbols take place in the encoding circuit *PAR COD* in *fig. 2* of the article of note [3]. Delay lines are used for interleaving the audio and parity symbols.

^[5] This circuit corresponds to the *ERCO* chip in *fig. 3* of the article by M. G. Carasso, J. B. H. Peek and J. P. Sinjou, this issue, p. 151.

^[6] In *fig. 9* of [3] a frame of this kind is represented by the bit stream *B*₂, from which the C&D symbol has already been removed.

DEC_2 in different words. If there are no errors present, DEC_2 will receive words that correspond to the encoding rules for C_2 , and it will pass on 24 audio symbols unchanged. DEC_2 can correct up to four errors, by means of erasure decoding. (In the current Compact Disc system full use is not made of this facility: DEC_2 is arranged in such a way that only two errors are corrected.) If DEC_2 receives a word containing five or more errors with given erasure positions, it will pass on 24 symbols unchanged, but provided with an erasure flag at the appropriate positions; this flag has in fact already been assigned by DEC_1 . A value for the erroneous samples can still be calculated with the aid of a linear interpolation.

As already mentioned, DEC_1 has been designed to allow the correction of single errors, and the detection of double and triple errors. The probability that DEC_1 will not detect quadruple or higher multiple errors is only about 2^{-19} . It may seem strange that the possibility of correcting two random errors is not utilized: in fact it would considerably increase the chance of DEC_1 failing to detect quadruple or higher multiple errors.

The probability P of quadruple or higher multiple errors passing DEC_1 without being detected can be approximated by the expression

$$P = \frac{1 + n_1(2^s - 1)}{2^{s(n_1 - k_1)}} \approx 2^{-19}$$

The numerator contains the number of error patterns with one error or none. (The factor $(2^s - 1)$ is the number of possibilities for one bit error per symbol; such a symbol can occur at n_1 positions. The value 1 is added because zero errors can be achieved in exactly one way.) This complete expression is to be related to the number of possibilities for filling in the parity: $2^{s(n_1 - k_1)}$. For proof of this equation the reader is referred to the literature [7].

When a disc is used for the recording and read-out of digital signals there are few random errors; most errors then occur as bursts. This is because the dimensions of a pit are small in relation to the most common mechanical imperfections such as dirt and scratches. It is therefore very important that multiple errors of this type cannot pass DEC_1 without being indicated with a high degree of certainty.

Since the bursts are 'spread out' over several words at the input of DEC_2 , the number of errors per word hardly ever exceeds the limit value $d_m - 1 = 4$. In this way most error bursts are fully corrected.

Specifications of CIRC

In assessing the quality of our CIRC decoder for Compact Disc applications its ability to correct both error bursts and random errors is of great importance.

The quality characteristics for the correction of bursts are the maximum fully correctable burst length and the maximum interpolation length. The first is determined by the design of the CIRC decoder and in our case amounts to about 4000 data bits, corresponding to a track length on the disc of 2.5 mm. The maximum interpolation length is the maximum burst length at which all erroneous symbols that leave the decoder uncorrected can still be corrected by linear interpolation between adjacent sample values. This 'length' is about 12 000 data bits; see the next section.

Random errors can also introduce multiple errors within one code word now and again; we shall return to this presently. The greater the relative number of errors ('bit error rate', BER) at the receiving end, the greater is the probability of uncorrectable errors. A measure for the performance of this system is the number of sample values that have to be reconstituted by interpolation for a given BER value per unit time. This number of sample values per unit time is called the sample interpolation rate. The lower this rate is at a given BER value, the better the quality of the system for random-error correction.

Table I. Specifications of CIRC.

Aspect	Specification
Maximum <i>completely</i> correctable burst length	≈ 4000 data bits (i.e. ≈ 2.5 mm track length on the disc)
Maximum interpolatable burst length in the <i>worst</i> case	$\approx 12\,300$ data bits (i.e. ≈ 7.7 mm track length)
Sample interpolation rate	One sample every 10 hours at BER = 10^{-4} ; 1000 samples per minute at BER = 10^{-3}
Undetected error samples (clicks)	Less than one every 750 hours at BER = 10^{-3} ; negligible at BER $\leq 10^{-4}$
Code rate	3/4
Structure of decoder	One special LSI chip plus one random-access memory (RAM) for 2048 words of 8 bits
Usefulness for future developments	Decoding circuit can also be used for a four-channel version (quadraphonic reproduction)

An objective assessment of the quality of the error-correcting system also requires an indication of the number of errors that pass through uncorrected and are therefore not corrected by the system. These uncorrected and uncorrected errors may produce a clearly audible 'click' in the reproduction.

The main features of the CIRC system are summarized in Table I. Details of the calculation relating to the quality can be found in the literature [7].

Concealment of residual errors

The purpose of error concealment is to make the errors that have been detected but not corrected by the CIRC decoder virtually inaudible. Depending on the magnitude of the error to be concealed, this is done by interpolation or by muting the audio signal^[8].

Two consecutive 8 bit symbols delivered by the decoder together form a 16 bit sample value. Since a sample value in the case of a detected error carries an erasure flag, the concealment mechanism 'knows' whether a particular value is reliable or not. A reliable sample value undergoes no further processing, but an unreliable one is replaced by a new value obtained by a linear interpolation between the (reliable) immediate neighbours. Sharp 'clicks' are thus avoided; all that happens is a short-lived slight increase in the distortion of the audio signal. With alternate correct and wrong sample values, the bandwidth of the audio signal is halved during the interpolation (10 kHz).

If the decoder delivers a sequence of wrong sample values, a linear interpolation does not help. In that case the concealment mechanism deduces from the configuration of the erasure flags that the signal has to be muted. This is done by rapidly turning the gain down and up again electronically, a procedure that starts 32 sampling intervals before the next erroneous sample values arrive. To achieve this the reliable values are first sent through a delay line with a length of 32 sampling intervals, while the unreliable values are processed immediately. The gain is kept at zero for the duration of the error and then turned up again in 32 sampling intervals. The gain variation follows a cosine curve (from 0 to 180° and from 180 to 360°) to avoid the occurrence of higher-frequency components. This also means that there are no clicks when the audio signal is muted, as in switching the player on and off, during an interval in playing or during the search procedure.

Maximum burst-interpolation length

Two associated 16 bit sample values, one from the left and one from the right audio channel, together form a sample set. If these sets were fed to the concealment circuit in the correct sequence, it would not be possible to interpolate more than one set from their reliable neighbouring sets. This would mean that in the case of an error longer than the maximum correctable burst length signal muting would very soon have to be applied.

By interleaving the sample sets it becomes possible to interpolate new sets for a given length of consecutive erroneous sets. This is done by alternating groups of 'even' sample sets with groups of 'odd' sets. Such a group, odd or even, can be interpolated from

its neighbouring group or groups. The maximum burst-interpolation length is thus equal to the length of such a group. In our system we have grouped the twelve 16 bit sample values of a frame in the way shown in fig. 5. The odd and even groups are separated by the parity values Q . Since these are not necessary for the reconstitution of the original signal and may therefore permissibly be unreliable, they increase the interpolation length. The maximum length with this grouping is certainly seven or even eight sample values, for some error patterns.

The delay lines corresponding to D_i^* (see fig. 3) in the encoder^[4] have placed eight frames between two successive sample values, after interleaving. The maximum burst length that can always be interpolated is therefore 56 frames. This presupposes, of course, that we are working with sample values consisting of two immediately consecutive symbols; the distance between all successive symbols is four frames, however. This is also the work of the delay lines D_i^* .

The delay lines corresponding to D_i (again fig. 3) in the encoder^[4] now ensure, however, that this distance is alternatively three and five frames, after interleaving. The distance of five frames is responsible for a decrease in the maximum interpolation length from 56 to 51 frames. We have tacitly assumed here that the burst also comes within a block of eight frames. If we discount this assumption, there is still a reduction of a length of 1 frame - 2 symbols. The maximum burst length that can be interpolated with certainty has now become 50 frames + 2 symbols.

So far we have taken no account of random errors that can be interpolated; this is the subject of the next and final section. At this point we shall simply mention the effect of the interpolation of such errors on the maximum interpolation length.

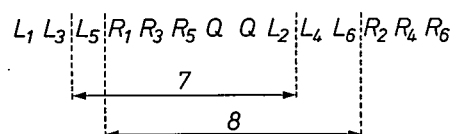


Fig. 5. Grouping of the sample values within a frame; L_i values for the left channel, R_i values for the right channel. For each sequence of seven unreliable values, new values can be calculated with certainty from reliable neighbours (e.g. if L_5 to L_2 are unreliable, the new values are interpolated from R_6 of the preceding frame and from the reliable values of the above frame). Given a favourable situation, new values can in fact be derived for eight consecutive values (e.g. values for R_1 up to L_6 from R_6 of the preceding frame, the reliable values of the above frame and L_1 of the succeeding frame).

[7] L. M. H. E. Driessen and L. B. Vries, Performance calculations of the Compact Disc error correcting code on a memoryless channel, in: 4th Int. Conf. on Video and data recording, Southampton 1982 (IERE Conf. Proc. No. 54), pp. 385-395.

[8] Error concealment takes place in the *CIM* chip in fig. 3 of the article of note [5].

To achieve good results in the treatment of random errors, the symbols are finally sent through a further set of delay lines Δ_i with a length of two frames. These delay lines, which serve purely and simply for 'restoring' uncorrected random errors, cause in their turn a reduction of the interpolation length by two frames. The final maximum burst length that is guaranteed capable of interpolation is thus 48 frames + 2 symbols, which corresponds to 12 304 bits.

Interpolation of random errors

If the symbols S_{oi} (fig. 3) after the decoder DEC_2 were already in the correct sequence, a pattern of errors might arise that would rule out any possibility

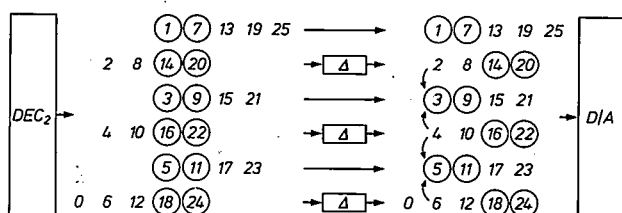


Fig. 6. The effect of the delay lines Δ_i with a length equal to the duration of two frames on the signal from the decoder DEC_2 . Each number represents a sample set, and a circle around a number is an erasure flag. A frame, consisting of 24 symbols or 6 sample sets, is represented by a complete column. The succession of frames on the left in the figure (sample sets that are irrelevant in the present context have been omitted) comes direct from DEC_2 and comprises a pattern of random errors, causing the total rejection of two consecutive frames (1, 14, 3, . . . 11, 24). It can be seen, however, that the chosen grouping enables a new value from reliable neighbours to be interpolated for each unreliable sample set, e.g. a value for 5 from 4 and 6. After passing through the delay lines Δ_i with a length equal to the duration of two frames, the sample sets are applied in the correct sequence to the D/A converter. If a frame in the succession of frames on the right in the figure were to be completely rejected, no interpolation would be possible.

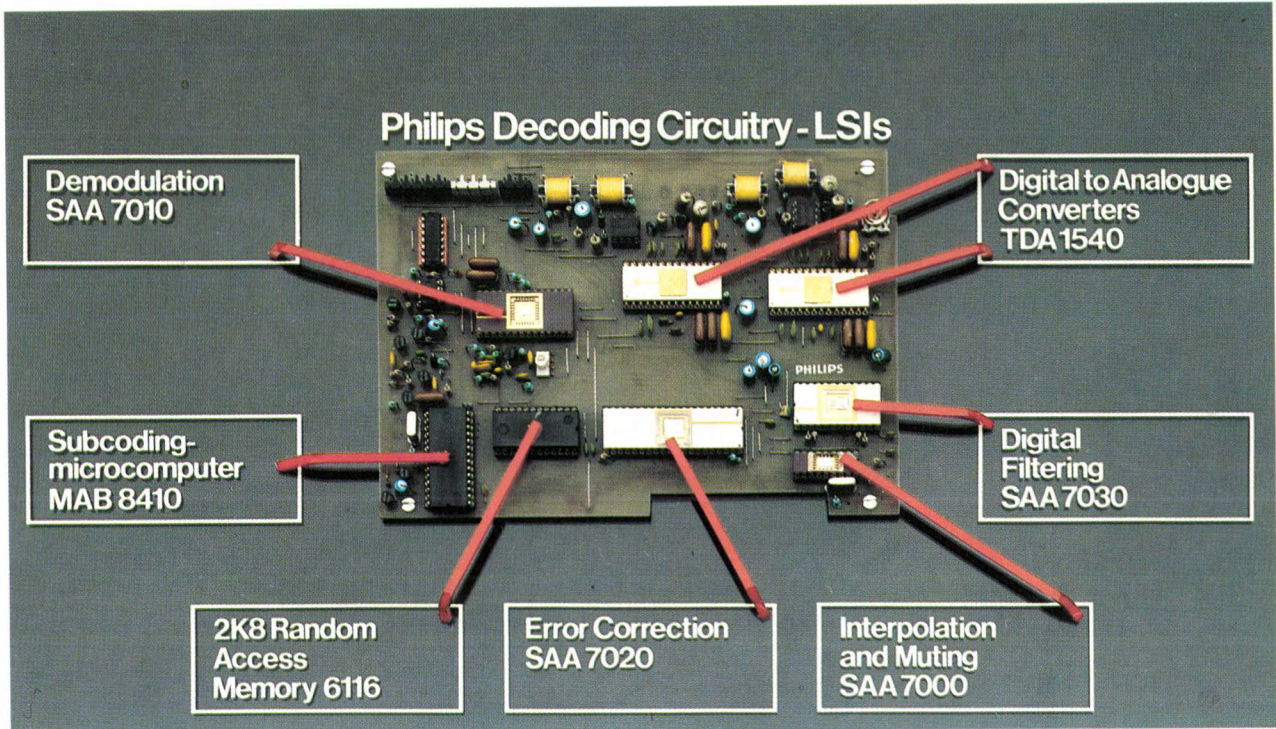
of interpolation, even though there were no long error bursts. This would happen if DEC_1 failed to detect an error but DEC_2 had detected it, resulting in the rejection of the entire frame at the output of DEC_2 . As described on page 170, however, the chance of DEC_1 failing is very small.

Since we prefer not to have to mute the audio signal, the concealment network contains a set of delay lines Δ_i , with a length of two frames, which ensure that the symbols of a single or double completely rejected frame from DEC_2 can still be interpolated from the reliable adjacent frames (see fig. 6). The probability that three completely rejected frames will occur within the interpolatable length determined by Δ_i is negligible.

After the symbols have passed through the delay lines Δ_i , they are in the correct sequence. Most of the errors have been corrected and the signal is ready for the digital-to-analog conversion^[9].

[9] See D. Goedhart, R. J. van de Plassche and E. F. Stikvoort, Digital-to-analog conversion in playing a Compact Disc, this issue, p. 174.

Summary. After an example showing how errors in a digital signal can be corrected, the article deals with the theory of block codes. The treatment of random errors and error bursts is discussed. Error correction in the Compact Disc system uses a Cross-Interleaved Reed-Solomon Code (CIRC), which is a combination of a (32,28) and a (28,24) code. One of the two decoders in the CIRC decoding circuit corrects single errors, the other corrects double errors. The residual errors are interpolated linearly to a length of up to 12 000 bits, and longer errors are muted. The interpolation and the signal muting take place in a separate chip, whose configuration is briefly discussed.



The most important printed circuit in the player contains a number of ICs in VLSI technology for reconstituting the analog audio signal from the digital signal read from the Compact Disc. The chips for error correction and digital-to-analog conversion are discussed in this issue.

Digital-to-analog conversion in playing a Compact Disc

D. Goedhart, R. J. van de Plassche and E. F. Stikvoort

Introduction

The last stage in the series of operations on the signal in the Compact Disc system is the return from the digital code to the analog signal, which has the same shape as the acoustic vibration that was picked up by the microphone.

After decoding and error correction the digital signal has the form of a series of 16 bit words. Each word represents the instantaneous numerical value of the measured sound pressure in binary form, and is therefore a sample of the acoustic signal. There are 44 100 of these words per second.

The digital-to-analog converter in the Compact Disc player generates an electric current of the appropriate magnitude for each word and keeps it constant until the next word arrives. The electric current thus describes a 'staircase' curve that approximates to the shape of the analog signal (*fig. 1a*). In terms of frequency, the steps in the staircase represent high frequencies, which extend beyond the band of the analog audio signal (20 Hz - 20 kHz). These high frequencies have to be suppressed by a lowpass filter; in the Compact Disc player their level should be reduced to at least 50 dB below that of the maximum audio signal.

If this high attenuation of the frequencies above the audio band is to be achieved solely with an analog lowpass filter, the filter must meet a very tight specification. It was decided to avoid this problem in the Philips Compact Disc player by introducing a filter operation, earlier in the digital stages. This was done by 'oversampling' by a factor of four: a digital filter, operating at four times the sampling rate ($4 \times 44.1 \text{ kHz} = 176.4 \text{ kHz}$) delivers signal values at this increased frequency, thus refining the staircase curve (*fig. 1b*) and making it easier to filter out the

high frequencies. As a result it is possible to make do with a relatively simple lowpass filter of the third order after the digital-to-analog conversion.

The conversion of the 16 bit words into an analog signal is performed in the Philips Compact Disc player by a 14 bit digital-to-analog converter available as an integrated circuit and capable of operating at the high sampling rate of 176.4 kHz. Partly because of the fourfold oversampling and partly because of

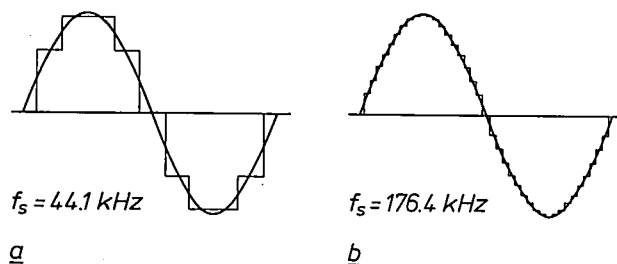


Fig. 1. A sinusoidal signal at 4.41 kHz sampled with a sampling rate f_s of 44.1 kHz (*a*) and with a frequency four times higher (*b*). In (*b*) the 'staircase' curve approximates more closely to the analog waveform, and the high frequencies present in the staircase signal are more easily filtered out.

the feedback of the rounding-off errors in antiphase, rounding off to 14 bits does not result in a higher noise contribution in the audio band. This remains at the magnitude corresponding to a 16 bit quantization (signal-to-noise ratio about 96 dB), so that even though there is a 14 bit digital-to-analog converter it is still possible to think in terms of a 16 bit conversion system.

In comparison with direct 16 bit digital-to-analog conversion, which must be followed by a lowpass filter with a sharp cut-off to give sufficient suppression of signals at frequencies above 20 kHz, our conversion system has a number of advantages. The first is the linear phase characteristic, which can be obtained with a digital filter, but not with an analog filter; the second is a filter characteristic that varies with the clock rate

Ing. D. Goedhart is with the Philips Audio Division, Eindhoven; Ir R. J. van de Plassche and Ir E. F. Stikvoort are with Philips Research Laboratories, Eindhoven.

and is therefore insensitive to variation in the speed of rotation of the disc. Finally, because the quantization steps are smaller, the maximum 'slew rate' that these circuits must be able to process is lower (the slew rate is the rate of variation of output voltage). There is therefore less chance of intermodulation distortion because the permitted slew rate has been exceeded.

44.1 kHz. The frequency spectrum of such a series is illustrated in fig. 3b^[1]. In theory it is infinite; above the baseband 0-20 kHz can be seen integral multiples of the sampling frequency with their left-hand and right-hand sidebands. Between these bands there are transition regions, the first for example being between 20 kHz and 24.1 kHz.

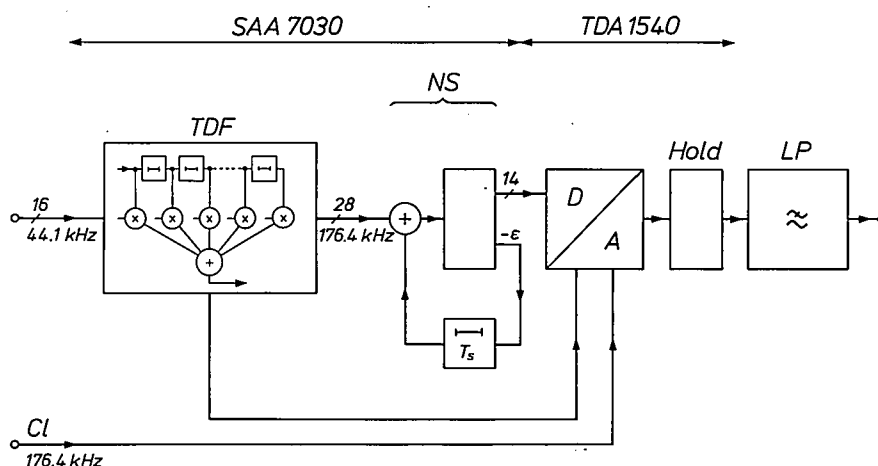


Fig. 2. Block diagram of the digital-to-analog conversion. *TDF* digital transversal filter which brings the sampling rate of 44.1 kHz to 176.4 kHz and attenuates signals in the bands around 44.1 kHz, 88.2 kHz and 132.3 kHz. *NS* noise shaper in which the rounding-off error is delayed by one period T_s after rounding-off to 14 bits and then fed back in the opposite sense. *D/A* 14 bit digital-to-analog converter. *Hold* hold circuit. *Cl* clock signal. *LP* lowpass 3rd-order Bessel filter.

The entire series of operations in the digital-to-analog conversion is shown as a block diagram in fig. 2. The oversampling takes place in the digital filter *TDF* to which the input signal is fed. The filter output signal is then rounded off to 14 bits, and the rounding error is fed back in the opposite sense in the noise shaper *NS*. The digital filter and noise shaper are located in a single integrated circuit in NMOS technology (type SAA 7030). This IC processes both stereo channels. Then follow the digital-to-analog converter *D/A* and a hold circuit, combined in a single IC type (TDA 1540) in bipolar technology; for each stereo channel there is a separate IC. The analog signal finally passes through a lowpass filter.

Suppression of frequencies above the audio band

Direct digital-to-analog conversion of the presented signal provides a series of analog signal samples (fig. 3a). These have the form of pulses that — in theory — are infinitely short, but have a content (duration times amplitude) corresponding to the sampled signal value. The repetition frequency is

This entire spectrum must not be passed on to the player amplifier and loudspeaker. Even though the frequencies above 20 kHz are inaudible, they would overload the player amplifier and set up intermodulation products with the baseband frequencies or possibly with the high-frequency bias current of a tape recorder. Therefore all signals at frequencies above the baseband should be attenuated by at least 50 dB.

To produce such an attenuation, an analog filter after the digital-to-analog converter will inevitably have to contain a large number of elements and require trimming. In addition a linear phase characteristic is required in the passband so that the waveform of pulsed sound effects will not be impaired. In the Philips Compact Disc player these requirements are met in a different way, by means of:

- fourfold oversampling of the signal in the digital phase,
- a digital filter operation,

[1] The principles of the digital processing of audio signals are explained in a very readable account by B. A. Blesser in Digitization of audio: a comprehensive examination of theory, implementation, and current practice, J. Audio Engng Soc. 26, 739-771, 1978.

— a hold function after the digital-to-analog conversion,

— a third-order Bessel filter in the analog-signal path.

A digital transversal filter is used for the filtering after oversampling. To understand the operation of the filter, we can think of it as consisting of 96 elements (*fig. 4a*), while the delay in each element is $(176.4 \times 10^3)^{-1}$ s, i.e. a quarter of the sampling period or $\frac{1}{4} T_s$. Four times in each period the filter takes up new data. At three of these four times the content of this data is zero, since the oversampling is done by the introduction of intermediate samples of value zero. This means that only 24 of the 96 elements are filled at any one time. The contents of each element are multiplied by a coefficient *c*. The filter provides data at a rate of 176.4 kHz; each number is the sum of 24 non-zero multiplications. In this way the filter always cal-

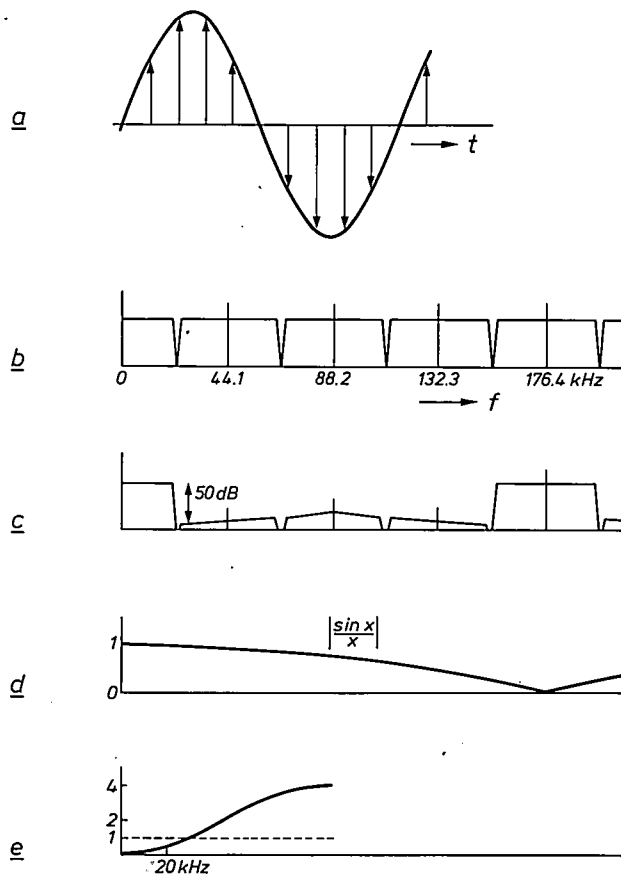


Fig. 3. *a*) A train of periodic pulses that sample an analog signal waveform. *b*) Frequency spectrum of such a pulse train. The pulse repetition frequency is 44.1 kHz, the sampled signal occupies the audio frequency band (0-20 kHz). *c*) Frequency spectrum for oversampling and filtering of the same signal at 176.4 kHz. It is now much easier to filter out the frequencies above the audio band. *d*) A hold circuit after the digital-to-analog converter keeps a signal sample at the same value until the arrival of the next sample. The frequency spectrum in *c* is thus multiplied by the function $|\frac{\sin x}{x}|$ with a first zero at 176.4 kHz. *e*) Noise spectrum after the noise shaper. In the audio range of interest the noise is considerably attenuated compared with the flat noise spectrum (*dashed line*) that would be obtained without noise shaping.

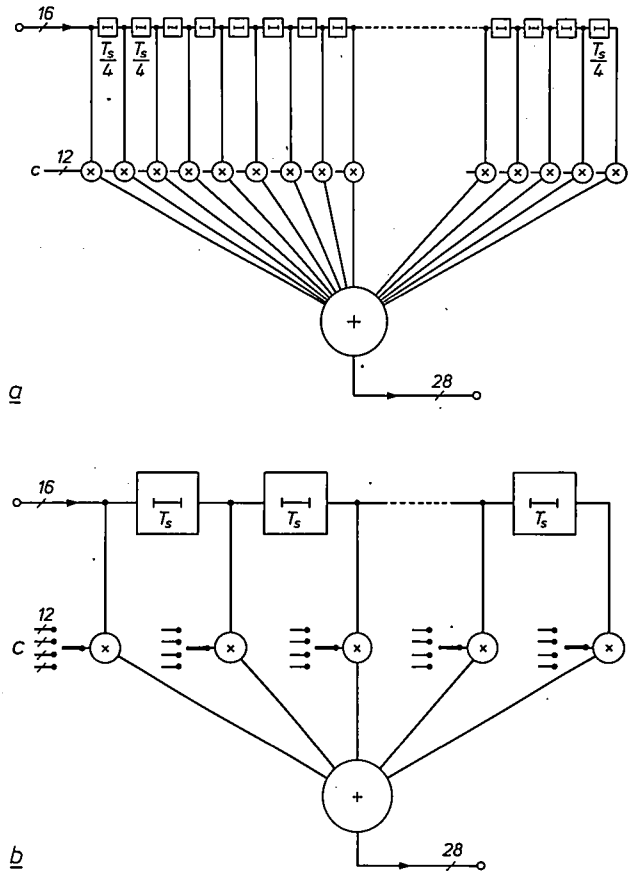


Fig. 4. Digital transversal filter. *a*) Filter consisting of 96 elements. A 16 bit word remains in each element for a quarter of the sampling period T_s . Since a new 16 bit word is only offered once per T_s , three-quarters of the elements are filled by the value zero. During the period T_s there are four multiplications by the 96 coefficients *c*; only 24 multiplications produce a product different from zero. These products are summed; in this way an output is provided four times in each sampling period, i.e. at a frequency of 4×44.1 kHz = 176.4 kHz. This means that there is a fourfold oversampling. *b*) An equivalent circuit that has been used in practice instead of (*a*) because it has 24 delay lines and multipliers instead of 96.

culates three new sample values at the locations of the zero samples.

The practical version of the filter is in fact somewhat different from the version referred to in the above explanation. In practice the filter consists of only 24 delay elements and a 16 bit word remains in each element for a time T_s (*fig. 4b*). During this time T_s the word is multiplied four times by a coefficient *c*, which is different for each multiplication. The products are also summed four times during the time T_s and passed to the output. The frequency at which these summated values appear at the output is therefore $4/T_s = 176.4$ kHz again.

The coefficients are numbers with 12 bits. Each product has a length of $16 + 12 = 28$ bits. The numbers have been chosen in such a way that the summation of 24 products does not introduce extra bits, so that the filter output consists of 28 bits with no rounding off.

The frequency spectrum of the oversampled and filtered signal is shown in fig. 3c. It can be seen that the bands in this spectrum around $1 \times$, $2 \times$ and 3×44.1 kHz are suppressed.

The digital-to-analog converter generates a current whose magnitude is proportional to the last digital

linear phase characteristic in the passband. This filter is simple and requires no highly accurate elements.

The hold function and the Bessel filter introduce some slight attenuation at the top of the passband. The digital filter is designed to correct this with a small overshoot (fig. 5).

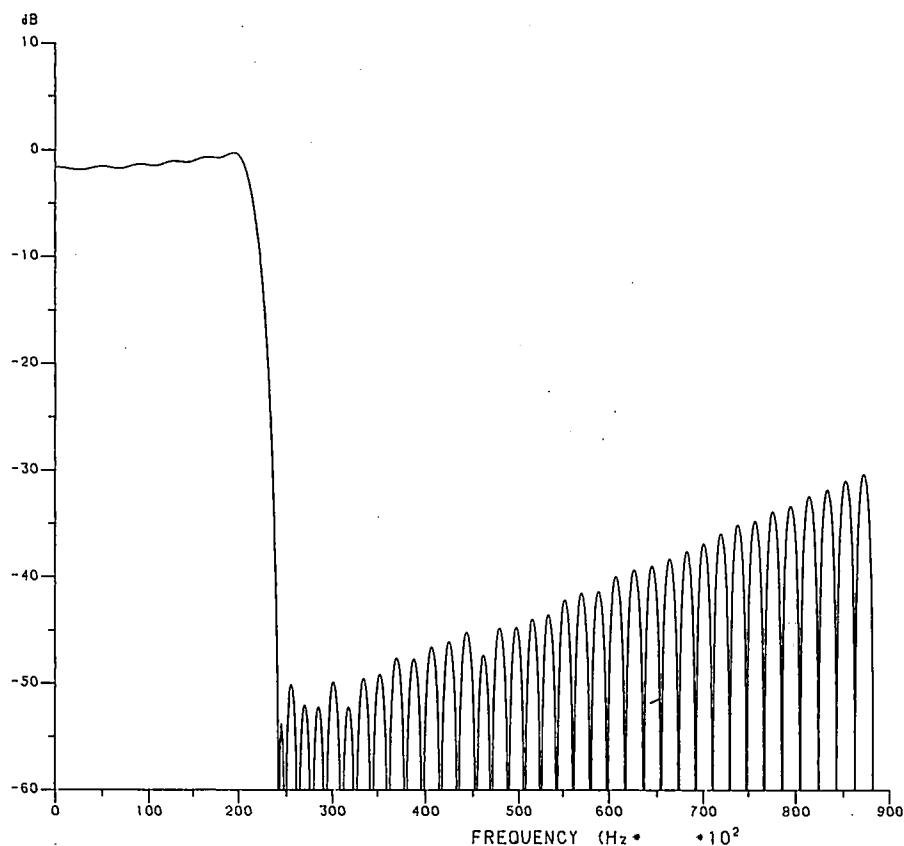


Fig. 5. Computer calculation of the detailed passband characteristic of the digital transversal filter. This has a small overshoot at the highest audio frequencies, which is used to compensate for the slight attenuation produced here by the curve in fig. 3d and the analog Bessel filter. A very sharp lowpass cut-off of 50 dB is obtained. The irregularity in the suppressed band is caused by rounding-off the filter coefficients to 12 bits.

word presented. This current is kept constant in a hold circuit until the next sample value is delivered, producing the staircase curve mentioned above. The signal samples have thus in theory changed from infinitely short pulses to pulses with the duration of a sampling period. This also has consequences for the frequency spectrum; the spectrum in fig. 3c is multiplied by a curve of the form $|(\sin x)/x|$ that has a first zero at 176.4 kHz (see fig. 3d). This gives an attenuation of signals in the 20 kHz sidebands on either side of 176.4 kHz by more than 18 dB. The hold effect causes no phase distortion.

The attenuation is still not sufficient, however. As a supplement, a lowpass Bessel filter of the third order is used, which has its -3 dB point at 30 kHz. The Bessel type of filter has been selected because of its

Suppression of the quantization noise

The presented signal, quantized to 16 bits, will contain some noise on conversion into an analog signal. This reproduces the errors due to the quantization in fixed steps. The root-mean-square value of the noise voltage in the sampled frequency band is $q/\sqrt{12}$, where q represents the magnitude of the quantization step. We see that when the quantization step is doubled, i.e. coding with one bit less, the noise voltage is also doubled, or, in other words, the noise level rises by 6 dB.

The samples that leave the filter at a repetition frequency of 176.4 kHz describe a signal with a bandwidth of 88.2 kHz. The quantization noise added due to the subsequent rounding off to 14 bits is spread over this band. With a signal of sufficient amplitude

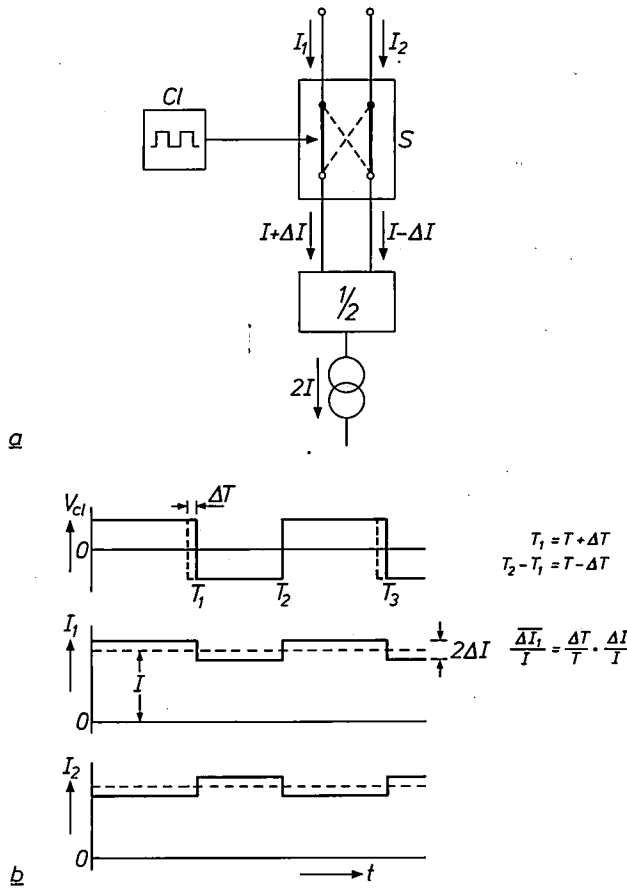


Fig. 6. a) Division of a current $2I$. Cl clock generator. S switches for periodically interchanging the two half-currents. b) The output currents I_1 and I_2 as a function of time t . Their mean value is the same. A difference between the mean output currents can be caused by an asymmetry ΔT of the clock signal V_{cl} . This difference is however an order of magnitude smaller than ΔI .

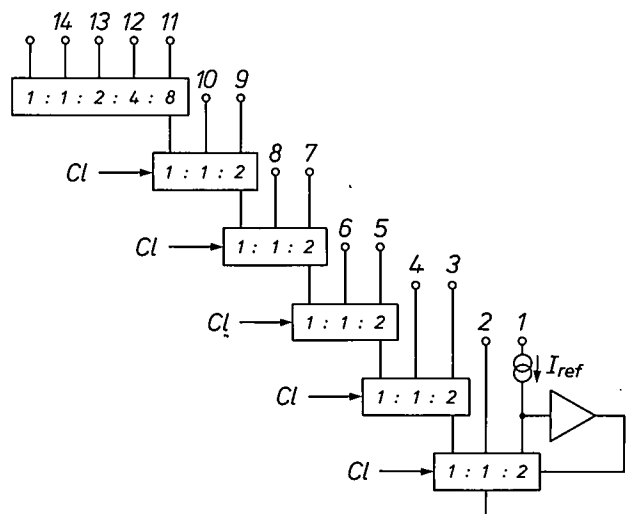


Fig. 7. Cascade of current dividers in the 14-bit digital-to-analog converter TDA 1540. The starting point is the reference current I_{ref} . Currents that are accurately equal to a half and a quarter of the input current are obtained in the divider stages by periodic interchanges; the clock signal Cl controls these interchanges. Only the four least-significant bits 11...14 are obtained by passive division.

and a sufficiently broad frequency spectrum this distribution is uniform, since the quantization errors for successive samples are in principle uncorrelated; the quantization noise is 'white' noise. Only the band from 0 to 20 kHz is relevant; this is only about a fourth part of the sampled band, and the noise power in the band from 0 to 20 kHz is therefore only a fourth part of the total noise power. This means that because of the fourfold oversampling the signal-to-noise ratio in the relevant frequency band is 6 dB better than would be expected with 14 bit quantization. It is thus about 90 dB, which is what would have been obtained with a 15 bit system without oversampling.

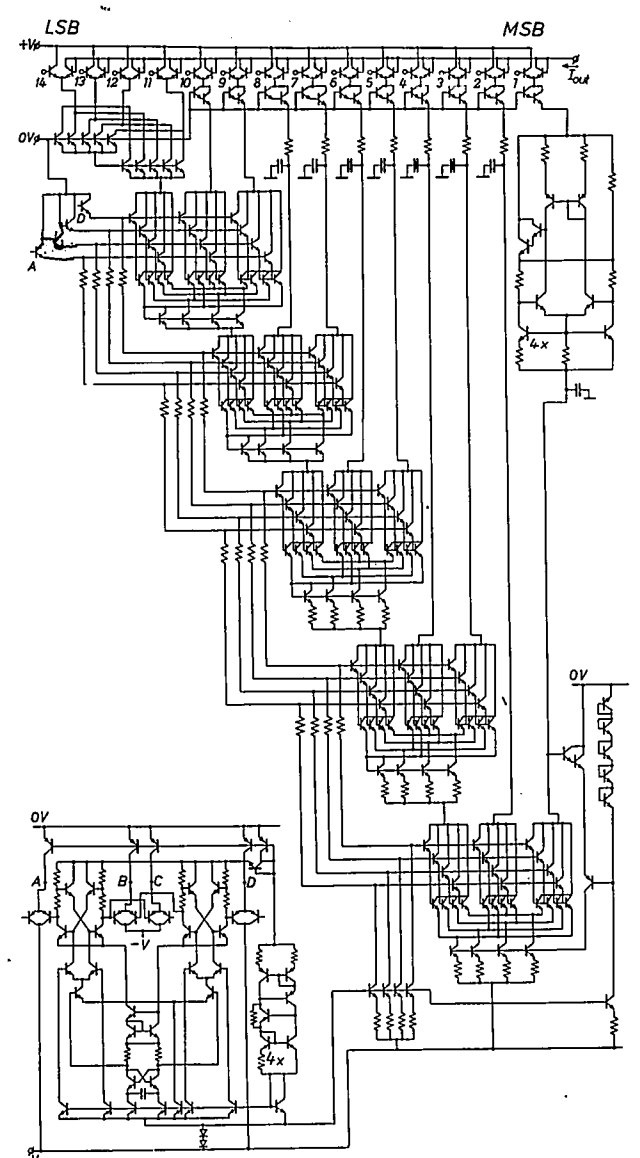


Fig. 8. Complete circuit diagram of the 14-bit digital-to-analog converter. The cascade of current dividers in fig. 7 can be identified here. The capacitors (above), which smooth out the ripple on the divider-output currents, are external. Bottom left: The clock generator.

In rounding off from 28 to 14 bits it is useful to compare successive rounding-off errors. If the analog signal is a direct voltage, successive samples will have the same rounding-off error. The audio signal will not contain any direct current; it will however contain slowly varying signals that will resemble a direct current in a short time interval. If the error produced in the rounding-off from 28 to 14 bits is now changed in sign and added to the next sample to arrive (see fig. 2), the average quantization error for slowly varying signals — i.e. low frequencies — can be reduced. This appears in the shape of the frequency spectrum of the quantization noise (see fig. 3e); at low frequencies the noise level is lower, at high frequencies it becomes higher. With a sampling rate of 176.4 kHz, it follows that a 7 dB gain in signal-to-noise ratio is obtained in the audio band (0-20 kHz). The ratio of the maximum signal to the noise contributed by the entire digital-to-analog conversion system described above is thus brought to about 97 dB, i.e. the value corresponding to a 16 bit quantization.

The digital-to-analog converter

The 14 bit digital-to-analog converter has been dealt with in detail elsewhere^[2]. Here we shall only indicate how it differs from other digital-to-analog converters.

A characteristic feature is the way in which currents are generated that are accurately related by a factor of 2; a digital-to-analog converter requires a set of such currents. The exact ratio is obtained by periodically interchanging the currents that are derived by dividing down by two from a constant reference current (see fig. 6), so that small differences are averaged out. This system is known as 'dynamic element matching'. Accurate division by four can be carried out with a

slightly more complicated circuit, also based on periodic interchange. The full series of current dividers is shown in fig. 7. Here Cl is the clock signal that controls the periodic switching; only for the four least-significant bits are the currents obtained from a passive division by means of differences in emitter area.

Fig. 8 shows the complete switching diagram of the 14 bit digital-to-analog converter. The cascade of divider stages can be seen in the figure. The ripple caused by the periodic switching is smoothed at the seven most significant bits by an RC filter; the seven capacitors (above in fig. 8) are externally connected.

The nonlinearity of the digital-to-analog converter is extremely low: between -20°C and $+70^{\circ}\text{C}$ it is less than 3×10^{-5} , or half the least-significant bit. The TDA 1540 integrated circuit is followed by the low-pass Bessel filter of the third order, and the analog signal appears at the output.

[2] R. J. van de Plassche and D. Goedhart, A monolithic 14-bit D/A converter, IEEE J. SC-14, 552-556, 1979.

Summary. The 16 bit words from the error-correcting circuit are converted into an analog signal by a 16 bit conversion system. This system consists of a digital transversal filter, in which the signal is oversampled 4 times (sampling rate 176.4 kHz) and then filtered in such a way that signals at frequencies above 20 kHz are attenuated by 50 dB after digital-to-analog conversion. The filter is followed by a noise shaper, which rounds off to 14 bits with negative feedback of the rounding-off error of the preceding sample. Next there is a 14 bit digital-to-analog converter, which is followed by a low-pass third-order Bessel filter. The signal-to-noise ratio of the complete system is about 97 dB. Even though the lowpass filter has a sharp cut-off the system is phase linear. The entire system, except for a few operational amplifiers, is contained in three integrated circuits; one for the digital filter (for both of the stereo channels) and two for the two digital-to-analog converters.



The Compact Disc has a diameter of only 12 cm. The information is recorded on the side of the disc shown here, protected by a transparent layer; the other side carries the label. The high density of the information gives a continuous playing time of more than an hour. The disc is packed in a 'de luxe' case, which also contains programme information.

An evacuated tubular solar collector incorporating a heat pipe

H. Bloem, J. C. de Grijs and R. L. C. de Vaan

The Philips VTR141 solar collector consists of a plate with a black coating, an evacuated glass envelope and a heat pipe that transfers the absorbed solar energy into the rest of the system efficiently. Although this seems a simple idea, the design of this collector, which has now been put into production by the Philips Lighting Division, is the result of close cooperation between different departments of Philips. The method of coating the plate was originally developed by the Philips research laboratories in Aachen, for other purposes. The coating is applied to the plate by an electroplating process developed by the Philips Plastics and Metalware Factories. Both the Research Laboratories in Eindhoven and the Electronic Components and Materials Division (Elcoma) have contributed to the adaptation of the heat pipe for this application. The method of coupling the collector to a water-circulation system was designed by the Philips Centre for Manufacturing Technology. The method of fitting the envelope around the heat pipe with its absorber plate and then evacuating it was derived from the glass-to-metal-seal and vacuum technology used in the mass production of tubular low-pressure sodium lamps, with the help of modifications devised by the Philips factory at Turnhout in Belgium.

The product resulting from the combined efforts of Philips specialists in these very different fields is the subject of the article below.

Introduction

The 'solar constant' is 1353 W/m^2 . This is the energy flux density of the sun's radiation in free space at the mean distance of the earth from the sun. Some of the radiation does not reach the earth's surface because of reflection and absorption in the atmosphere, but in favourable conditions the irradiance at sea level can be about 1000 W/m^2 . Ever since the end of the last century attempts have been made, with *solar collectors*, to put this energy flux to use for all kinds of heating purposes. Solar collectors convert solar energy into thermal energy, with an efficiency that can amount to more than 80%. In this article we shall discuss the Philips VTR141 solar collector, whose design is the result of cooperation between various departments in different Product Divisions of Philips^[1].

A body in free space assumes a temperature such that the total thermal radiation from the body is in equilibrium with the incident solar radiation. For a black body near the earth this temperature is about 300 K. If it is desired to utilize the solar energy at a much higher temperature, it is necessary to concentrate the solar radiant flux by reflectors or lenses, to compensate for the much higher thermal radiation loss per unit area of surface. For higher concentrations, the optical system must follow the sun; this requires moving parts, with all the associated complications^[2]. Concentration becomes less necessary at more modest 'load' temperatures, and for use below about 450 K, i.e. about 180 °C, simpler types of col-

H. Bloem is with the Philips Electronic Components and Materials Division (Elcoma), Eindhoven; Ir J. C. de Grijs and Drs R. L. C. de Vaan are with the Philips Lighting Division, Eindhoven.

[1] 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.

[2] For systems using concentration, see for example A. B. Meinel and M. P. Meinel, *Applied solar energy*, Addison-Wesley, Reading, Mass., 1977.

lector — with little or no concentration — are generally adequate. There are problems even with fixed reflectors; they become dirty and are liable to corrode. Also, the benefit from *diffuse* solar radiation increases as concentration is diminished. Even at noon on a clear day in the subtropics, diffuse radiation accounts for as much as 20% of the total solar energy flux at the earth's surface.

The application referred to for fixed collectors — thermal energy at load temperatures up to about 180 °C — is economically of great importance: some 40% of the total energy consumption in the industrialized countries comes into this category. It can be further subdivided into three ranges: load temperatures up to about 100 °C, up to about 130 °C and up to about 180 °C (see *Table I*). The VTR141 is intended for the first range. A collector for the second range, the VTR261, is being developed and one for the third range is planned.

Fig. 1 is a diagram showing how solar collectors are generally used in domestic hot-water systems. The collector *C* is part of a circuit in which water is circulated by a pump. The heat collected is transferred to the water in the storage tank *S*, from which hot water can be drawn off. If the water is not hot enough, it is brought up to temperature by an auxiliary heater *H*. On draw-off, the storage tank is replenished from the cold-water supply (*M*). The pump *P* is stopped by a control device *R* if the water returning from the collector is colder than the water in the tank ($T_C < T_S$). Systems of this type can also be used for other purposes, e.g. for space heating in living or working areas and for various industrial processes.

In sunny parts of the world 'flat-plate collectors' are already widely used for systems in the first category of *Table I*. The construction of a conventional flat-plate collector is illustrated in *fig. 2a*. The solar energy is absorbed by a black-coated plate *A*. This is attached to a tube grid *T*; the heat absorbed is transferred to a liquid, usually water, that flows through the tubes. Heat losses to the environment are reduced by means of a glass cover *G* and back insulation *I*.

In these simple collectors, the major cause of heat loss is radiation. It can be reduced by a 'spectrally selective' coating of the absorber plate, possibly in conjunction with a simple form of concentration. A selective surface has a high absorptance for solar radiation but a small emissivity for thermal radiation. Such surfaces are possible because the spectrum of solar radiation and that of thermal radiation are almost completely separated; there is only a small overlap at a wavelength of about 2.5 μm .

The collector in *fig. 2a* also loses heat by convection and conduction through the air between the absorber

plate and the glass cover. The next stage in the development is therefore *evacuation* of the collector. However, since a flat box is not very suitable for evacuation, *tubular envelopes* are used (*fig. 2b*).

Table I. Applications of fixed solar collectors, arranged in order of load temperature.

Load temperature	Examples of application
1. Up to 100 °C	Hot water: space heating, domestic hot water, industry (e.g. paper, textiles, food processing, etc.) Hot air: drying processes, desalination
2. Up to 130 °C	Low-pressure steam: air conditioning and many industrial processes
3. Up to 180 °C	Medium-pressure steam: conversion of thermal to mechanical energy

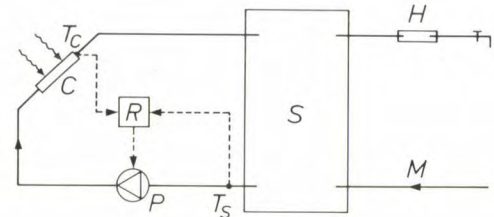


Fig. 1. Solar-collector installation for domestic hot water, schematic. *C* solar collector. *P* pump. *S* hot-water storage tank. *H* auxiliary heater. *M* connection to cold-water supply. The pump *P* is stopped by the control device *R* if the temperature T_C of the water in the collector falls below the temperature T_S of the water in the storage tank.

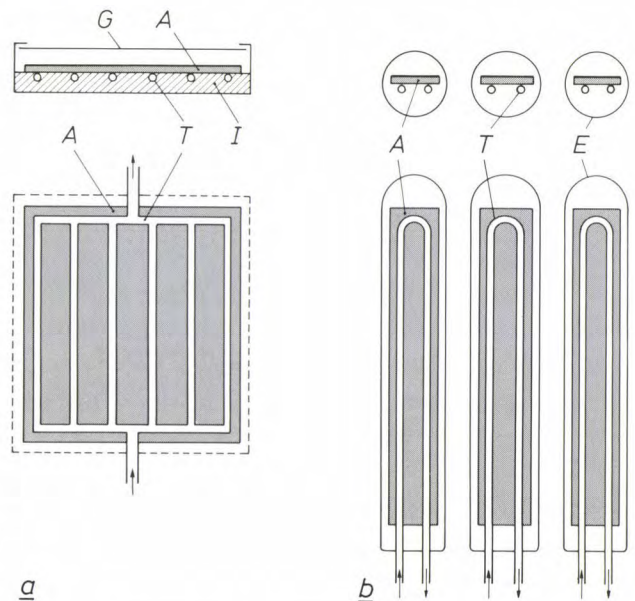


Fig. 2. *a*) Conventional flat-plate collector. *A* black absorber plate. *T* water tubes for heat transfer. *G* glass cover. *I* insulating material. *b*) Evacuated tubular collectors. *A* absorber plate. *T* water tubes. *E* glass envelope.

An essential feature of a collector is good heat transfer from the absorber plate to the circulating water. The water-tubes are often soldered or welded directly to the plate, or integrated with it, to ensure good heat transfer. The heat can also be transferred to the circulating water by a 'heat pipe' [3]. In a heat pipe, latent heat is transferred by means of an evaporation and condensation cycle. The effective thermal conductivity of a heat pipe can be more than 1000 times greater than that of a copper bar of the same dimensions [4]. As we shall see later, heat pipes have some special advantages.

The VTR141 is an evacuated tubular collector with a highly selective absorber surface and a heat pipe. We shall now discuss this collector in more detail.

The VTR141

The construction and operation of the VTR141 are illustrated schematically in *fig. 3a*. The selectively coated plate *A* (shown for clearer presentation at right angles to the normal position) is attached to the heat pipe *H* with good thermal contact. The plate is completely contained by the evacuated glass envelope *E*.

the condenser, and three sizes of clamping block, designed for water pipes of different dimensions.

Even at high irradiance, e.g. 900 W/m², a single VTR141 tube delivers only about 40 W (see p. 187). Practical applications therefore require a large number of tubes. They are grouped together to form 'modules' consisting of a number of collector tubes clamped side by side to a single length of water pipe

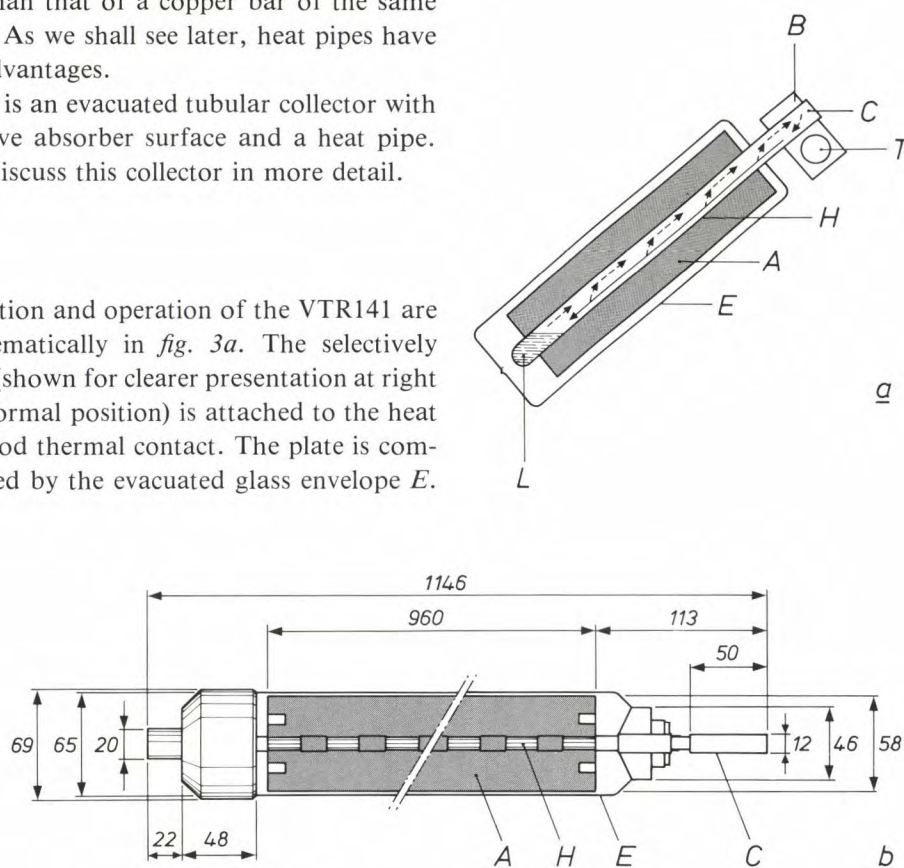


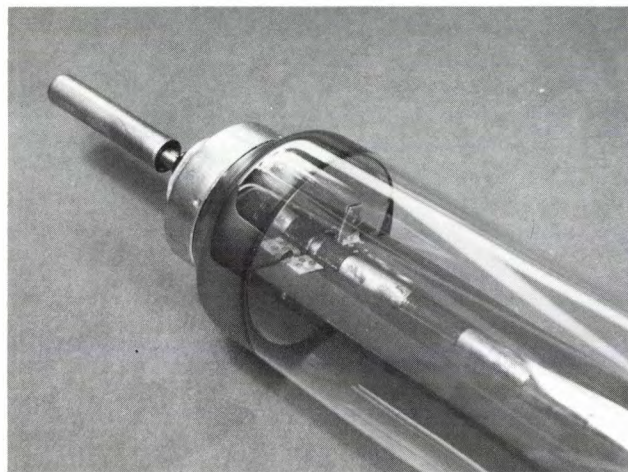
Fig. 3. *a)* Structure and operation of the VTR141, schematic. *b)* Dimensioned drawing (in mm) of the tube. *E* glass envelope. *H* heat pipe. *A* absorber plate (for clarity the plate is shown in *(a)* at right angles to the normal position of operation, in which it is directed towards the sun). *B* aluminium clamping block. *T* water pipe. *C* condenser. *L* working fluid (isobutane in the VTR141). Heating makes the working fluid evaporate (dashed arrows in *(a)*), the vapour condenses in *C* and the liquid flows back (solid arrow). This cycle gives a very high effective conductivity.

At one end the heat pipe protrudes beyond the envelope. The protruding part, the 'condenser' *C*, makes good thermal contact with tube *T* of the water circuit through an aluminium clamping block *B*. As soon as the 'evaporator' — the part of the heat pipe with the absorber plate — becomes hotter than the condenser, the working fluid contained in the heat pipe evaporates; the vapour rises, condenses in *C*, the condensate flows downwards under gravity and the cycle repeats. *Fig. 3b* is a dimensioned drawing of the tube. *Fig. 4* shows the upper part of the collector with

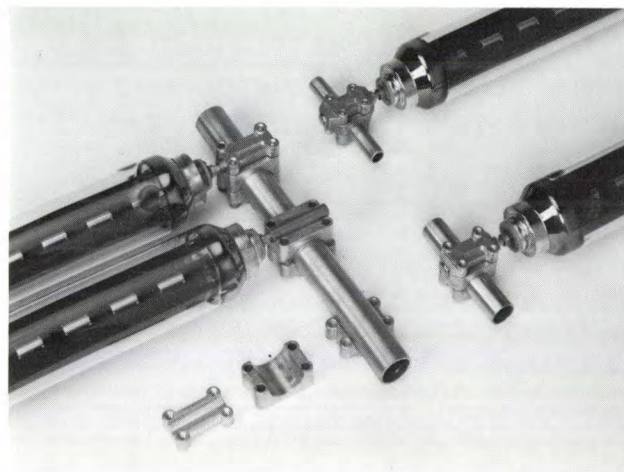
(*fig. 5*). This 'header', with the clamping blocks, is enclosed in a casing with insulating material. The number of collector tubes in a module has to be restricted to 20 or 25; with more tubes there would be difficulties because of the expansion of the header when it becomes very hot — e.g. if the water circulation fails ('stagnation'), preventing the water circuit

[3] See for example U. Ortabasi and F. P. Fehlner, *Solar Energy* 24, 477, 1980.

[4] See for example G. A. A. Asselman and D. B. Green, *Philips tech. Rev.* 33, 104 and 138, 1973, and P. Dunn and D. A. Reay, *Heat pipes*, Pergamon Press, Oxford 1978.



a



b

Fig. 4. a) The end of the collector tube with the condenser. b) Clamping the condensers to the water pipes. The clamping blocks fit on to standard water pipes of outside diameters 15, 22 or 28 mm.

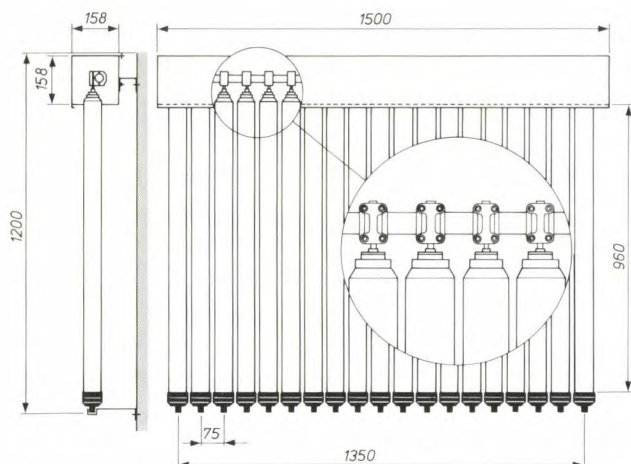


Fig. 5. Dimensioned drawing of a collector module 1.5 m in width with 19 collector tubes. The 'header' with the clamping blocks and upper sections of the collectors is enclosed in insulating material in a casing.

from absorbing heat. The modules can be connected in series, in parallel or in a combination of the two. With the freedom of choice of the header-pipe diameter, the water circuit may be designed to limit the pressure drop and hence the pumping power required.

The clamping blocks

The use of clamping blocks is one of the main features of the VTR141. With this system the risk of leaks in the water circuit is no greater than in an ordinary water circuit. Moreover, since the fluid in the heat pipe is completely isolated from the circulating water, there is no risk of it causing contamination. During installation and maintenance any collector tube can be fitted or replaced without disturbing the water cir-

cuit. Furthermore, this circuit can be very simple in construction so that it can easily be drained if necessary, e.g. to avoid freezing in cold weather.

This is all very different from systems like the ones in fig. 2, where the circulating water flows along or through the absorber plates, and from systems with a heat pipe in which the condenser is directly immersed in the circulating water. In both these systems the water circuit must be broken to insert or replace a collector tube, and each tube requires a water-tight seal. In addition, systems with water circulating through the absorber plates are usually also difficult to drain.

The clamping method of thermal coupling does introduce an additional thermal resistance between

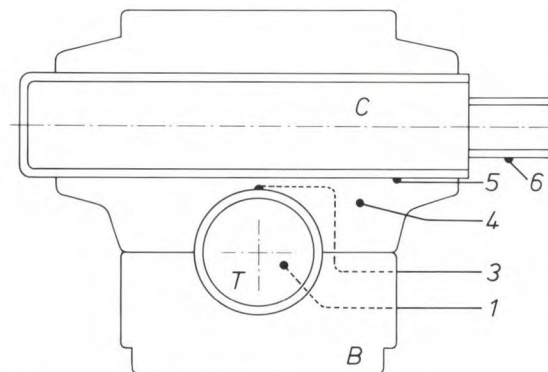


Fig. 6. Measurement of the distribution of the thermal resistance between the water in the pipe *T* and the working fluid in the condenser *C*. The measurement is carried out with six thermocouples: 1 in the water, 2 (not in the figure) in a copper block on the header outside the clamping block (*B*), 3 on the header in the clamping block, 4 in the clamping block, 5 on the condenser in the block, 6 on the condenser neck. The measured resistance distribution is as follows: 1-2: 0.20 °C/W, 2-5: 0.08 °C/W, 5-6: 0.28 °C/W. The total resistance of 0.56 °C/W is thus almost entirely attributable to the liquid-to-wall interfaces; the resistance between the walls via the clamping block is only 1/7 of the total.

the heat-pipe working fluid and the circulating water. However, since the overall thermal resistance arises almost entirely at the fluid-to-wall interfaces, it is very little higher than it would be if the condenser were directly immersed in the water. This has been verified by the measurements quoted in the caption to *fig. 6*. The thermal resistance measured between the walls of the condenser and the header via the clamping block amounts to only a seventh of the total resistance between the two fluids. To achieve this result the blocks must be accurately fitted to both components and must be firmly clamped. To improve the thermal contact between the heat-pipe working fluid and the condenser wall a phosphor-bronze gauze is mounted inside the condenser and soldered to the wall. The total thermal resistance depends in part on the lengths of pipe enclosed by the block. The block as designed gives a temperature difference of about 30 °C between the heat-pipe working fluid and the water for a solar irradiance of 900 W/m² and a water temperature of about 100 °C.

The heat pipe

The use of the heat pipe brings several significant advantages. In the first place, the heat pipe transfers heat to the condenser, but not in the reverse direction, because there is no liquid in the condenser when the absorber is colder than the condenser — at night for example, or when it is cloudy during the day. In addition, the thermal capacity of the heat pipe plus the absorber is low. Any heat absorbed is therefore immediately transferred to the water circuit and no heat is lost, even when the sun is temporarily obscured by cloud.

Furthermore, the system is protected from overheating, and excessive pressure in the water circuit is thus avoided. Evacuated tubular collectors can easily become too hot under high solar irradiance because the thermal losses in this type of collector are very small; the problem can often be a serious one, requiring special precautions. When a heat pipe is used, however, the heat-transfer mechanism cuts off if the condenser reaches the critical temperature of the working fluid: when all the fluid is in vapour form, there can no longer be any evaporation-condensation cycle. (As we shall see later, the cut-off mechanism is in fact somewhat more complicated.) The use of a heat pipe also protects the collector from any damage due to freezing.

A detailed chemical investigation was carried out to assist in the choice of working fluid. Very few liquids met the requirements: internal stability and no reaction with the wall material of the heat pipe for at least 1000 hours at 300 °C. The widely used Freon

refrigerants, for example, are not sufficiently stable. The choice was limited to saturated hydrocarbons with a chain of three carbon atoms: propane, isobutane and neopentane. The critical temperature of propane (96 °C) is just too low. Isobutane, which has a critical temperature of 135 °C, was chosen for the VTR141. With the temperature difference of about 30 °C between the isobutane and the water, as mentioned above, the system water cannot overheat under normal conditions of operation. Neopentane, which has a critical temperature of 161 °C, is used in the VTR261, making it suitable for applications at temperatures up to about 130 °C (see Table I).

The heat pipe itself is made of steel, copper-plated inside and out. The part outside the vacuum envelope must be corrosion-resistant at temperatures up to 300 °C. The condenser is therefore made of a copper-nickel alloy. It is connected to the evaporator by a flexible copper 'neck' that can take up differences in expansion of the frame of a module and the header for varying irradiance and water temperature. With frequent deformation — as would be expected for the neck, especially since the internal pressure reaches 44 bars at 300 °C — the danger of oxidation is increased. The neck is therefore protected from oxidation by a coating of phosphor-nickel.

The thermal contact between the working fluid and the water can be improved by making the condenser longer and increasing the size of the clamping block. For a given overall collector length, however, this reduces the length of the absorber plate. The design has been carefully optimized in this respect as well.

The absorber plate

The plate is a thin sheet of copper-plated steel, coated with a layer of 'black cobalt', a mixture of oxides of cobalt. This material was first used for selective black coatings by the Philips research laboratories in Aachen. Selective surfaces are characterized by two quantities, the 'solar absorptance' α and the 'thermal emissivity' ϵ . The solar absorptance is defined as the fraction of the incident solar power absorbed by the plate. The thermal emissivity is the radiated power density divided by the power density radiated by a 'true' black surface at the same temperature. In practice, the aim is to make α as large as possible and ϵ as small as possible.

Fig. 7 gives the values of α and ϵ for black-cobalt coatings produced by the electroplating process now used, as a function of the electric charge transferred per m²; this quantity is a measure of the thickness of the coating. The thickness for the VTR141 is given a value in the shaded range, resulting in values of 93-95% for α and 5-7% for ϵ . These values are a little

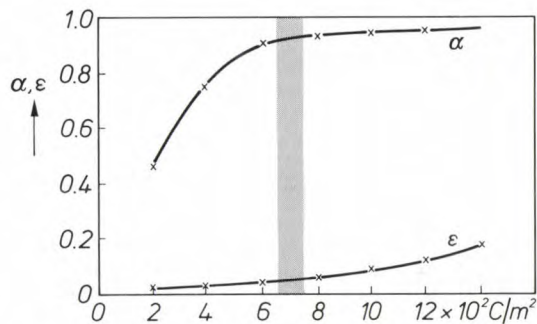


Fig. 7. The solar absorptance α and the thermal emissivity ϵ of copper-plated steel plates provided with black-cobalt coatings of different thicknesses. A measure of the thickness is the charge transferred per m^2 during the electroplating (coulombs/ m^2). The curve for ϵ applies for a temperature of 90°C . In plating the absorber for the VTR141, the charge density is adjusted so that it falls in the shaded area, resulting in values of 93-95% for α and 5-7% for ϵ .

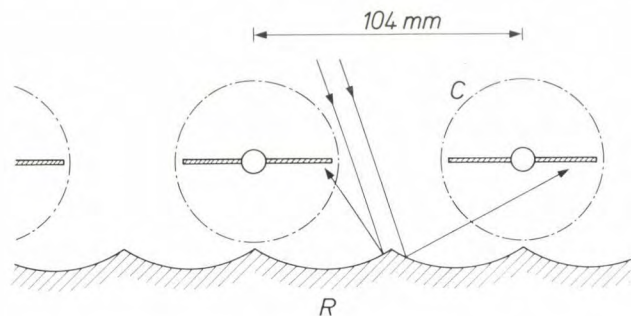


Fig. 8. Ripple reflector (R), enabling optimum use to be made of the radiation incident at the gaps between the absorber plates (schematic). C collector tubes.

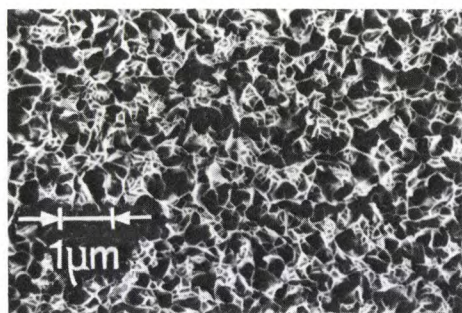


Fig. 9. SEM picture of the black-cobalt layer.

better than those that can be obtained with other commonly used materials, such as 'black chrome'. Our main reason for choosing black cobalt, however, was concerned more with the manufacturing technology; we shall return to this presently. The plate is coated on both sides. This means that the radiation incident on the gaps between the absorber plates can be captured with the aid of a reflector. By mounting the tubes sufficiently far apart (with a centre-to-centre spacing of 104 mm), and making use of the specially designed 'ripple' reflector shown in *fig. 8*, the total

power falling on each absorber plate may be increased by some 40%. Reducing in this way the required number of tubes per unit of aperture area reduces both the thermal losses and the cost.

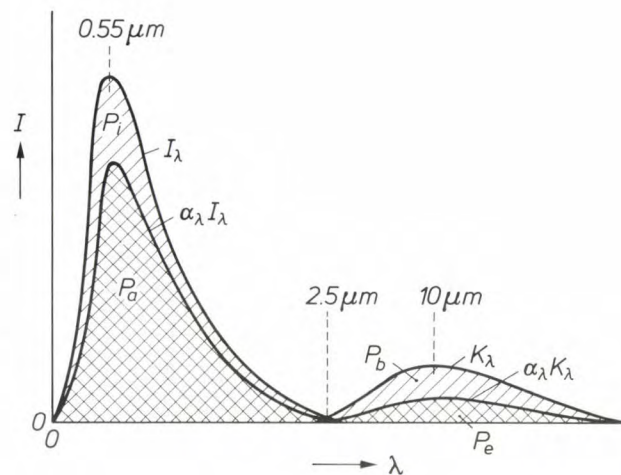


Fig. 10. Illustrating the definition of the solar absorptance α and the thermal emissivity ϵ at the surface of a body. The curves represent 'spectral' powers, i.e. powers per unit of wavelength interval, as a function of wavelength λ . I_λ : the spectral power of the solar radiation incident on the surface; $\alpha_\lambda I_\lambda$: the spectral power absorbed by the surface; α_λ is the 'spectral' absorptance of the surface. K_λ : the spectral power radiated by a 'true' black body at a given temperature T (e.g. 350 K); $\alpha_\lambda K_\lambda$: the spectral power radiated by the surface if it has this temperature (Kirchhoff's law). The areas (P) below the curves are the corresponding total powers (i.e. $P_i = \int I_\lambda d\lambda$, and so on). For clarity the curves are greatly distorted in both the horizontal and vertical directions; in reality the two maxima are at $0.55 \mu\text{m}$ and $10 \mu\text{m}$, and the ratio P_i/P_b is approximately equal to $(5760/350)^4 \approx 7 \times 10^4$. The solar absorptance α is equal to P_a/P_i , the thermal emissivity ϵ to P_e/P_b . The thermal emissivity is a function of temperature. For an 'ideal' selective surface we would have $\alpha_\lambda = 1$ for $\lambda < 2.5 \mu\text{m}$; $\alpha_\lambda = 0$ for $\lambda > 2.5 \mu\text{m}$.

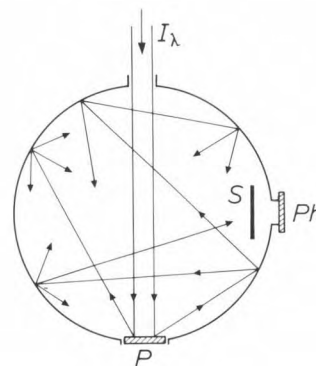


Fig. 11. Measurement of the spectral reflectance ρ_λ of the surface of a sample with an integrating sphere, schematic. This quantity is the fraction of incident radiation of wavelength λ that is reflected, both diffusely and specularly. Owing to the multiple, diffuse reflections at the completely white inner surface of the sphere, a homogeneous radiation distribution is produced. The reading of the photocell Ph is thus representative of the total radiation reflected by the sample (P). The photocell is protected by the plate S from direct radiation from the source or the sample itself. The incident radiation (I_λ) is modulated to enable the reflected radiation to be distinguished from other possible contributions (e.g. thermal emission from P). The instrument is calibrated with surfaces of known reflection coefficient (e.g. MgO , which is 'completely white'). The spectral absorptance α_λ is obtained in this way ($\alpha_\lambda = 1 - \rho_\lambda$).

The choice of the type of coating and the method of manufacture was originally based on two physical principles. In the first place the material itself must have a high absorption for high-energy photons ($\lambda < 2.5 \mu\text{m}$) and a low absorption for low-energy photons ($\lambda > 2.5 \mu\text{m}$); it would thus have to be a semiconductor with a sharp absorption edge at $2.5 \mu\text{m}$ ('wavelength approach'). In the second place the coating should have the structure of a 'mountain landscape' with sharp peaks spaced at about $2.5 \mu\text{m}$; the short waves are entrapped by the peaks, which increases the absorption, while the longer waves see the coating as 'homogeneous' ('wave-front approach'). When a coating having these properties overlays a reflector, the combination absorbs high-energy photons, but it reflects low-energy photons and thus does not emit these either. The black-cobalt-coated plate used in the VTR141 gives an approximation to this behaviour. An SEM picture of the surface is shown in fig. 9. Selective coatings based on other physical principles [5] have also been developed recently. However, making successful coatings is still mainly an empirical process.

For determining α and ϵ the 'spectral absorptance' α_λ — the ratio of the absorbed radiation to the incident radiation at wavelength λ — should be known for the entire wavelength region of interest. Fig. 10 shows plots of the incident solar power (I_λ) and the power (K_λ) radiated by a perfect black body of temperature T (e.g. 350 K), both per unit wavelength interval and per unit area, as functions of wavelength. K_λ is the Planckian distribution for the temperature T , while I_λ is approximately the Planckian distribution at the temperature of the sun (5760 K), modified by the effect of the atmosphere and the glass envelope or cover. The absorbed power is then given by the curve $\alpha_\lambda I_\lambda$, the emitted power (according to Kirchhoff's law) by $\alpha_\lambda K_\lambda$. Then we have (see fig. 10):

$$\alpha = P_a/P_i, \quad \epsilon = P_e/P_b,$$

where P_i is the area under the curve I_λ , and so on.

The function α_λ is determined by measuring the spectral reflectance $\rho_\lambda (= 1 - \alpha_\lambda)$ by means of an integrating sphere (fig. 11), which captures 'all' the reflected radiation (i.e. from both specular and diffuse reflection).

Characteristics of the VTR141

Various kinds of measurements have been made on heat pipes, on individual collector tubes and on modules.

Firstly, the tubes were tested in the Lighting Division laboratory with small solar simulators (fig. 12). These use halogen incandescent lamps with reflectors that are dichroic (i.e. transparent to thermal radiation); the shape of these reflectors was specially developed by Philips for this application. With the exception of a small range in the spectrum ($0.3\text{-}0.4 \mu\text{m}$), the simulator satisfies the Standard 93-77 drawn up for this purpose by the American Society for Heating, Refrigeration and Air-Conditioning Engineers (ASHRAE).

Fig. 13 gives the power (\dot{Q}_t) delivered by a collector tube to the water circuit, as a function of the condenser temperature, for different irradiances (G_T); the conditions are specified in the caption. The measurements were made without background reflection

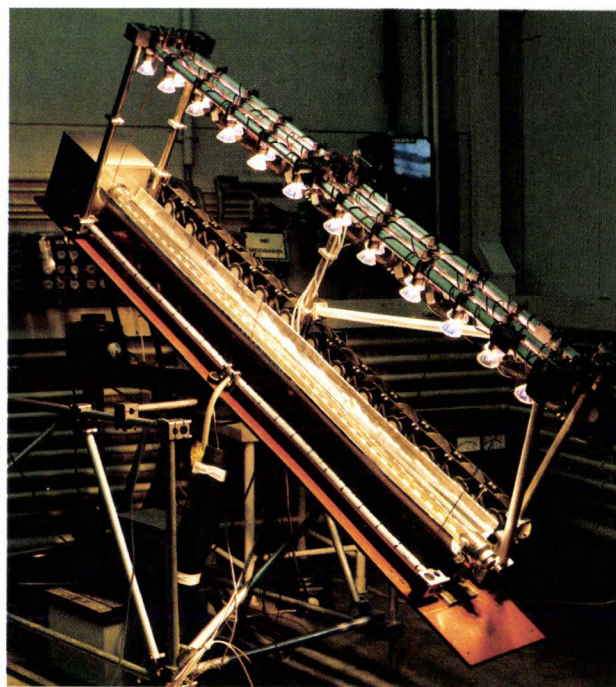


Fig. 12. The solar simulator used in our measurements on individual tubes.

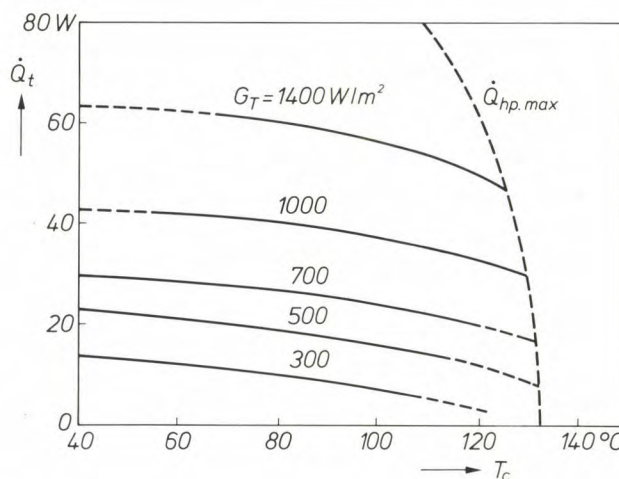


Fig. 13. The power \dot{Q}_t delivered per collector tube (VTR141) as a function of the condenser temperature T_c , for different values of the irradiance G_T . The curves were measured with the solar simulator shown in fig. 12. At $\dot{Q}_{hp,max}$, which is the power at which the heat-transfer mechanism of the heat pipe 'cuts off', the power delivered drops abruptly to a much lower value. In the region of $\dot{Q}_{hp,max}$ instabilities occur, diminishing the reproducibility of the measurements (dashed parts on the right). The measurements were carried out at a (simulated) wind speed of 4 m/s, an ambient air temperature of 25°C , hardly any background reflection (reflectance 0.04) and with the tubes at an inclination of 30° . The area of the absorption plate (one side) is $960 \times 58 \text{ mm}^2$ (see fig. 3b).

(reflectance 0.04). This implies that the curve for $G_T = 1400 \text{ W/m}^2$, for example, can be applied for an irradiance of 1000 W/m^2 if the ripple reflector of fig. 8 is used.

At a given G_T , the value of \dot{Q}_t decreases with increasing T_c , because the losses due to radiation, convection and conduction increase as the components of

[5] See for example O. P. Agnihotri and B. K. Gupta, Solar selective surfaces, Wiley, New York 1981.

the system become hotter. The curves fall sharply at the curve ' $\dot{Q}_{hp,max}$ ', where the heat-pipe mechanism cuts off. The curves are not readily reproducible near $\dot{Q}_{hp,max}$, because of instabilities that we shall return to shortly.

The diagram in fig. 13 is a useful aid in the preliminary design of a solar collector system. If we are making a design for a climate where the solar radiation is 700 W/m^2 for some hours of the day, and the required water temperature is 75°C (the required condenser temperature is about 100°C), then fig. 13 shows that each tube delivers about 25 W . Therefore 40 tubes would be required for each kilowatt of demand during those hours. Further calculations must be made to include the effects of the variation in solar radiation and other climatic conditions, the use of reflectors, the construction of the module, and so on.

The 'maximum heat transfer' $\dot{Q}_{hp,max}$ of the heat pipe is measured as indicated in the diagram of fig. 14a. The evaporator is contained in the temperature-controlled bath B_1 at temperature T_e . The condenser is enclosed in a copper block (temperature T_c), which is connected to a second temperature-controlled bath B_2 at temperature T_b by means of a calibrated thermal resistance W . The heat flux \dot{Q} through W is equal to $(T_c - T_b)/W$. This relation between \dot{Q} and T_c is plotted in fig. 14b for two values of the temperature T_b (solid lines). At a given T_b (e.g. T_{b1}) B_1 is heated up, and T_c increases. As long as the heat pipe is working properly, T_c remains equal to T_e . At a certain moment, however, the heat pipe stops working. T_c then increases more strongly, whereas T_e remains constant or decreases. The heat flux at the reversal point is $\dot{Q}_{hp,max}$ (see fig. 14b).

As long as the heat-pipe mechanism works normally, the temperature in the heat pipe is virtually the same everywhere. The rate of heat transfer is then equal to the effective rate of mass transfer (vapour upwards, liquid downwards), multiplied by the latent heat of evaporation. As the temperature in the heat pipe approaches the critical temperature, the latent heat of evaporation decreases, so that a steadily increasing mass transfer is necessary for a given heat transfer. Now there is a limit to the amount of liquid that can flow back per second, and in fact this limit determines the cut-off points on the straight lines in fig. 14b. When the experiment in fig. 14a is performed at a higher T_b (e.g. T_{b2}), a somewhat higher T_c is reached, but because the latent heat of evaporation is smaller, $\dot{Q}_{hp,max}$ is also smaller. The maximum heat transfer is therefore a decreasing function of T_c , becoming zero at $T_c = T_{cr}$ (dashed curve in fig. 14b). The maximum rate of mass transfer depends on various factors. Near to the critical temperature the return flow is probably

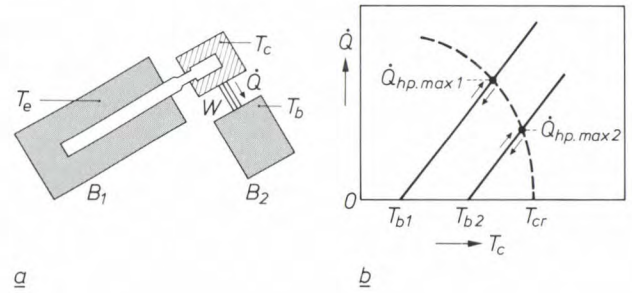


Fig. 14. Measurement of $\dot{Q}_{hp,max}$, the power at which the main heat-transfer mechanism of the heat pipe cuts off. a) The evaporator is contained in a temperature-controlled bath B_1 (temperature T_e). The condenser (temperature T_c) is connected via a known thermal resistance W with a second temperature-controlled bath B_2 (temperature T_b). b) Solid lines: the heat flux \dot{Q} through the thermal resistance W as a function of T_c for two values of T_b ; \dot{Q} is equal to $(T_c - T_b)/W$. If at a given T_b the temperature of bath B_1 is gradually increased, T_c gradually rises until the heat pipe stops working; from that moment on T_c remains constant or decreases. The heat flux at that moment is the desired $\dot{Q}_{hp,max}$.

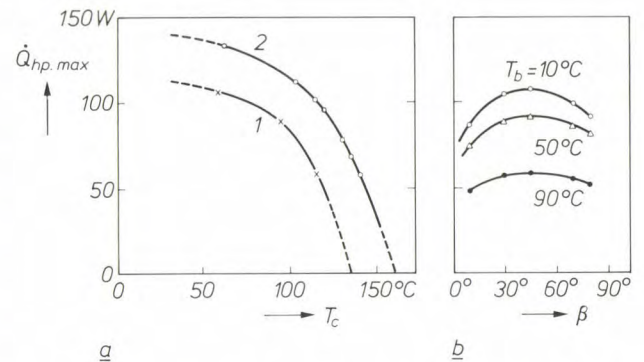


Fig. 15. a) The quantity $\dot{Q}_{hp,max}$ as a function of the condenser temperature T_c for isobutane (1) and neopentane (2), determined by the method of fig. 14. b) The quantity $\dot{Q}_{hp,max}$ as a function of the inclination angle β (angle between the tube and the horizontal plane), for three water temperatures (i.e. T_b in fig. 14a).

hampered by the 'vapour storm' in the opposite direction; the instabilities mentioned earlier are probably associated with this effect. Fig. 15a gives results of measurements for the heat pipes of the VTR141 and the VTR261.

The operation of the tube depends on the force of gravity to make the fluid flow back again, and hence on the inclination of the tube. Fig. 15b shows that this dependence is not very great for a wide range of slopes.

Efficiency

The efficiency η of a collector is defined as the power \dot{Q}_u delivered to the water divided by the incident radiant power. It should be clear, however, whether reference is made to the power incident on the absorber surface, on the aperture area of the col-

lector or on its gross area; the three corresponding definitions of η are all used in various contexts. We have taken 'incident radiant power' to refer to the power incident on the aperture. This seems to us to be the most suitable definition for comparing different systems. The 'aperture' is generally taken to be the projection of the transparent part of the collector on the plane of the absorber plate. For an array of tubular collectors, the effective aperture per tube then amounts to the length of the transparent part multiplied by the centre-to-centre spacing of the tubes. In our results, the VTR141 tubes are assumed to have the minimum centre-to-centre spacing of 75 mm.

We shall now briefly derive an expression for the efficiency η . The power \dot{Q}_u delivered to the water is equal to the power \dot{Q}_a absorbed by the plate less the losses \dot{Q}_l :

$$\dot{Q}_u = \dot{Q}_a - \dot{Q}_l = \bar{\alpha}\tau SG_T - k'\Delta T,$$

where G_T is the irradiance in W/m^2 . The power incident on the aperture S is thus SG_T . The power \dot{Q}_a absorbed by the plate follows from this by multiplying by the factor $\bar{\alpha}\tau$; this is the product of α , the solar absorptance, and τ , the transmittance of the glass, averaged over the aperture. The averaging corrects for the difference in area between the absorber plate and the aperture, for the lower transmission of the glass at the edges as compared with the centre of the tube, and so on. ΔT is the difference between the average temperature of the water in the collector and the temperature of the ambient air. The losses due to convection and conduction are proportional to ΔT . For small values of ΔT this is also true for the radiation losses, because the plate exchanges little thermal radiation energy except with the glass envelope, which is at about ambient temperature. For $\eta (= \dot{Q}_u/SG_T)$ we therefore find:

$$\eta = \eta_0 - k\Delta T/G_T = \bar{\alpha}\tau - k\Delta T/G_T,$$

where k is equal to k'/S . Because of the form of this expression, η is usually plotted as a function of $\Delta T/G_T$.

The solar simulator shown in fig. 12 was used for measuring the efficiency of individual VTR141 tubes. The results are given in fig. 16, together with those for two conventional flat-plate collectors. The curvature of the lines arises mainly because the radiation losses increase more than linearly with ΔT for large temperature differences.

Before taking a closer look at these results, we should note that the really important quantity is the efficiency of the complete system in which the collectors are used. The tubes will usually be contained in modules mounted on a roof, and it will be desirable to

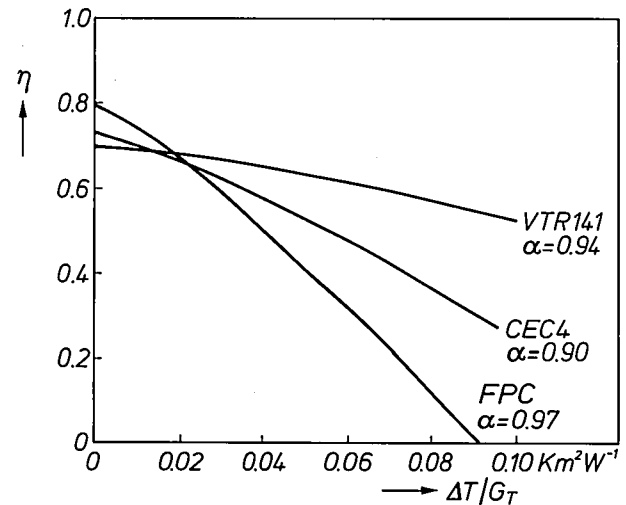


Fig. 16. Efficiency η as a function of $\Delta T/G_T$, for the VTR141 and two flat-plate collectors (FPC, CEC4). ΔT is the temperature difference between the water in the collector and the ambient air, G_T is the irradiance. The simulated wind speed was 4 m/s. The CEC4 is a flat-plate collector chosen as a reference collector by the EEC committee for solar energy. FPC is a simple flat-plate collector of typical performance. The CEC4 and the FPC each had a single glass cover. The FPC, unlike the two others, has a non-selective coating with a high absorptance ($\alpha = 0.97$) and a high emittance ($\epsilon \approx 0.95$); the η -curve therefore starts high but falls rapidly. The curve for the CEC4 also starts somewhat higher than the curve for the VTR141 because the absorption plate completely fills the aperture. At higher values of $\Delta T/G_T$ the VTR141 always has the advantage, because of its very low losses. No reflector was used in the VTR141.

make the best possible use of the available roof area. The optimum arrangement will depend on a variety of factors, such as the temperature of the water, the irradiance, its variation with time, the ratio of diffuse to direct irradiance and so on. The most important variables available for the optimization are the number of tubes, their spacing, the option of using a reflector and the type of reflector. For domestic hot water in our mid-European climate the arrangement in fig. 8 with the ripple reflector is very effective.

Efficiency curves and applications

As fig. 16 shows the VTR141 is slightly less efficient than the other two collectors at very low values of $\Delta T/G_T$. This is partly because the solar absorptance is somewhat lower than for one of the conventional collectors, the 'FPC' (see the values of α in the figure). The absorptance of selective coatings is generally rather lower than that of the 'best' black material. For soot, which is non-selective, α is about 0.98, but for selective coatings it is usually no more than 0.96. The difference, however, arises mainly because the absorber plate occupies less of the aperture in the VTR141 than it does in flat-plate collectors. At larger values of $\Delta T/G_T$ the VTR141 is more efficient than

the two other collectors, because of the lower losses. The crossover point at $\Delta T/G_T \approx 0.02 \text{ Km}^2\text{W}^{-1}$ corresponds at an irradiance of, say, 500 W/m^2 to a temperature difference of only 10°C .

Table II gives some examples of applications for different temperature ranges below 100°C . Fig. 16 shows that the VTR141 is not particularly suitable for

Table II. Applications for different temperature ranges below 100°C .

Load temperature	Application
20-25 $^\circ\text{C}$	swimming pools
25-35 $^\circ\text{C}$	space heating by floor heating or warm air
50-65 $^\circ\text{C}$	domestic water
60-90 $^\circ\text{C}$	space heating by radiators

low load temperatures and high irradiance, e.g. to maintain the temperature of a swimming pool in a sunny climate. The load temperature for domestic hot water is not particularly high either, and for this reason it is often thought that conventional flat-plate collectors would be better for this application than tube collectors. In fact, however, the efficiency of the VTR141 will often be much higher *on average* than that of the other collectors when the variations in irradiance and water temperature are taken into ac-

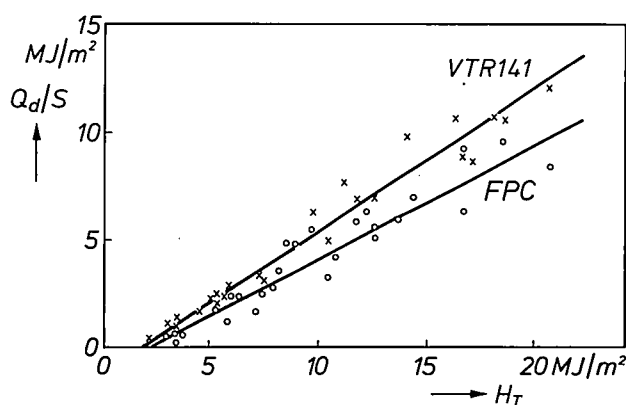


Fig. 17. Daily yield per m^2 of aperture, Q_d/S , as a function of the daily irradiation H_T for a module with VTR141 tubes (crosses) and the FPC (circles) of fig. 16. The quantity Q_d is the heat collected by the water in the system in one day. The 'aperture' S of the module is the product of the length of the transparent part of the tubes, the centre-to-centre spacing of the tubes in the module and the number of tubes; the gaps between the tubes are thus included. The centre-to-centre spacing was 75 mm ; no reflector was used. Each dot and each circle represents the result of one day of measurements (from sunrise to sunset). Only those results are included for which the daily average of the temperature difference between the water in the system and the ambient was $40^\circ\text{C} (\pm 3^\circ\text{C})$.

count. In the morning, for example, the water temperature T_w will be low, of course, but G_T will also be small. During the day, T_w will steadily increase, and in the afternoon G_T decreases again. The quantity $\Delta T/G_T$ thus varies appreciably during the day.

These considerations apart, the VTR141 does have the advantage, as mentioned earlier, that the losses are very small when there is no incident solar radiation (cloudy weather, or at night), because of the low thermal capacity and the 'diode effect' of the heat pipe.

Daily yields

To compare the overall behaviour of the VTR141 with that of a conventional flat-plate collector, the two types of collector were incorporated in a system in our outdoor test installation for two months (September and October 1979). Every 'day' (from sunrise to sunset) measurements were made of the total heat delivered to the water in the system; the total irradiation (in J/m^2) for the day was measured with a solari-meter. Fig. 17 gives the daily yield per m^2 of the 'system aperture' as a function of the daily irradiation; each measured point represents the result for one day, and points are only included if the daily average of the temperature difference ΔT between the system water and the ambient amounted to $40^\circ\text{C} (\pm 3^\circ\text{C})$. The fairly large scatter of points for each collector is due to the daily differences in weather conditions and in the time variation of the radiation; these differences may lead to different daily yields for equal irradiations. It can be seen that a VTR141 is clearly superior to a flat-plate collector. At higher average values of ΔT the advantage of the VTR141 is even greater.

Figs 18 and 19 are photographs of collector modules on a flat roof and on a sloping roof.

Manufacture

In developing the VTR141 we aimed at a product suitable for quantity production, and requiring only small quantities of inexpensive materials. The following points are of interest here.

The heat pipe is simpler in design than heat pipes that may have to be used in different positions. These have to be provided with a 'wick', i.e. capillary material that returns the liquid from the condenser to the evaporator, whatever the position of the tube within a certain range [4]. Since the VTR141 is always used on an incline, with the condenser at the upper end, a wick is unnecessary. The inside of the tubes is however given a surface treatment to improve the surface wetting and hence facilitate the return flow.



Fig. 18. Pilot installation with VTR141 tubes on the roof of the Philips central kitchens in Eindhoven.



Fig. 19. Four modules with VTR141 tubes, set in a tiled roof.

The absorber plate is a thin fin, which is simply crimped — but with good thermal contact — to the heat pipe (see fig. 3*b*). The coating is applied in an inexpensive process in which a metal strip is passed continuously through an electroplating bath. This is made possible by the use of black cobalt as the coating material, since the deposition of black cobalt only requires a very low current density, unlike the deposition of the widely used 'black chrome', for example. In contrast with black chrome, black cobalt is not stable upon exposure to air, but in evacuated collectors this is not a disadvantage.

Finally, the glass-to-metal-seal and vacuum technology developed by Philips during many years of mass production of lamps is used to great advantage for the assembly and evacuation of the glass envelope.

BIBLIOTHEEK NAT. LAB.
N.V. PHILIPS
GLOELAN WABSTEN
POSTBUS 80.000
5600 JA EINDHOVEN

Summary. The Philips VTR141 solar collector is designed for systems with a water circuit in which the temperature of the water is 95 °C or less. Later versions will be produced for load temperatures up to 130 °C and up to 180 °C. A selective black plate, attached to a heat pipe, is contained in an evacuated glass tube. The part of the heat pipe projecting out of the tube, the 'condenser', is clamped by an aluminium block to a water pipe in the system; this method of assembly greatly facilitates installation and maintenance. The present collector tubes are 1.15 m long and 65 mm in diameter. For practical use they are grouped together in modules. The solar heat absorbed is transferred to the condenser by means of an evaporation-condensation cycle of isobutane in the heat pipe. The heat pipe gives a 'diode effect': heat is never transferred from the water circuit to the plate, and a 'cut-off effect': the heat transfer ceases if the water in the circuit becomes too hot. The black-cobalt coating gives excellent selectivity (strong absorption of solar radiation, weak thermal emission) and is applied in a simple, continuous electroplating process. In terms of efficiency and average daily yield the collector is superior to a simple conventional flat-plate collector above an average load temperature of 25 to 30 °C; at higher load temperatures the advantage increases. The tube is corrosion-resistant and freeze-proof; the design is intended for quantity production.

Ink-jet printing

M. Döring

Printing by spraying ink directly on to the paper is not a new idea. Indeed, Lord Kelvin invented his 'siphon recorder' in 1873. This device was capable of recording the telegraph signals sent by cable across the Atlantic. The ink-jet principle is being applied today for printing out computer results. In its latest development, known as the DOD principle (Droplet On Demand), the droplets are impelled directly at the paper on the receipt of control signals. The author has improved the design of the extremely small ejector nozzles in such a way that the ink droplets are always ejected in exactly the same direction. This gives much better printing. At the same time technologies have been developed for making these special nozzles economically. An improved understanding of the dynamic behaviour of the ink in the exit channel has enabled us to double the speed of printing the characters.

Introduction

For some time computer printers have been on the market that print the characters directly on the paper by 'shooting' ink jets or droplets at it. Since they contain few parts moving at high velocity, these printers are quiet in operation and are usually reliable.

Their precursor was an ink-jet oscilloscope, developed in the sixties, that could record high-frequency signals directly on paper. This was done by giving the ink jet an electric charge and deflecting it in an electric field ^[1]. The same principle is used in a printer introduced in 1976 for word processors and computers. The disadvantage of the system employed in these machines is that the ink particles that are not charged and deflected have to be intercepted and returned to the system by means of a pump and filters.

The DOD principle (Droplet On Demand) does not have this disadvantage. In this system a droplet is ejected through a fine aperture and applied directly to paper, without deflection, on receipt of an electrical control signal. The droplet is generated by a pressure wave in the fluid ^[2], produced by applying a voltage pulse to a piezoelectric ceramic ^[3]. There are now various printers that work on this principle, but they

are expensive to manufacture and therefore unsuitable for semi-professional use. The quality of the characters printed with some of them also leaves something to be desired.

Printers operating on the DOD principle have for some years been the subject of investigations at the Philips laboratories in Hamburg. We have investigated a device that generates the pressure wave by means of flat piezoelectric plates, as discussed in this article. The development of a special nozzle has made it possible to eject the droplets with precision in the same direction. This has substantially improved the quality of the characters produced by the print heads. Technologies have also been developed for the economical manufacture of these nozzles, which have an inside diameter of only 50 μm . A print head based on our investigations is used in the P2131 printer developed for the Philips P2000 microcomputer. As a result of improving the dynamic characteristics of the printing system it will be possible to increase the speed of future systems to 6000 droplets per second. This is twice the speed of current DOD printing systems.

In the following we shall first consider the principle of the print head and the mechanism of droplet formation, and then look at the ejector nozzle and

Dipl.-Phys. M. Döring is with Philips GmbH Forschungslaboratorium Hamburg, Hamburg, West Germany.

methods of making it. Next some practical types of print heads will be discussed. The article concludes with a description of the measures that can be taken to achieve a considerable increase in printing speed.

Principle of the print head

The most important part of the print head is the fluid-pressure generator, in which the disc of piezoelectric ceramic *PXE* (see *fig. 1*) is the energizing com-

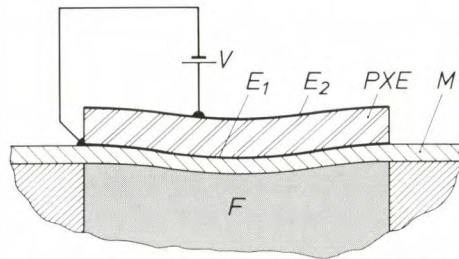


Fig. 1. The fluid-pressure generator. *PXE* plate of piezoelectric ceramic. *E₁* and *E₂* electrodes consisting of evaporated metal films. *M* metal plate. *F* fluid filling the system. *V* direct voltage applied across *E₁* and *E₂*, causing flexure of the bilaminar plate consisting of *PXE* and *M* and setting up a pressure wave in the fluid.

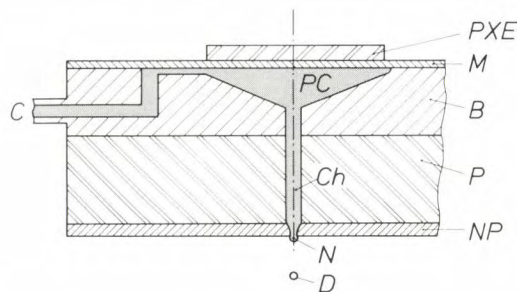


Fig. 2. Part of the print head, built in a sandwich construction. *B* metal body in which the pressure chambers *PC* have been recessed. *P* plastic plate with connecting channels *Ch*. *NP* nozzle plate with nozzles *N*. *D* ejected droplet. *C* common fluid-feed channel. See *fig. 1* for the other symbols.

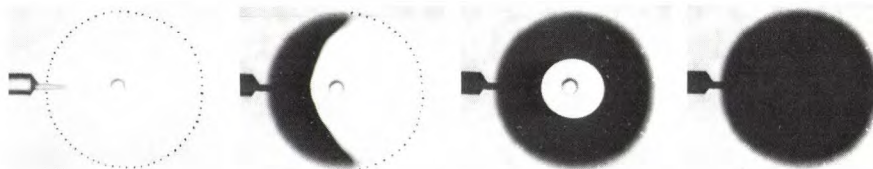


Fig. 3. Filling the pressure chamber with fluid. Since the pressure chamber is shaped like a flat cone and the fluid feed is radial, an air bubble is trapped in the chamber, and is subsequently expelled along the connecting channel *Ch* (see *fig. 2*).

ponent. Attached to the upper and lower faces of the disc are two electrodes *E₁* and *E₂*, consisting of evaporated metal films. The disc is cemented to a metal plate *M*, which is in contact with the fluid *F*. When a direct voltage *V* is applied between the electrodes, the disc becomes thicker or thinner, but a radial contraction or expansion also occurs^[4]. The result is that the

combination of *PXE* and *M*, called the bilaminar plate, flexes as shown in an exaggerated way in the figure. This flexing sets up a pressure wave in the fluid.

Fig. 2 shows how the pressure generator is mounted in the print head, which is a sandwich construction consisting of the metal plate *M*, the body *B* with pressure chambers *PC*, the nozzle plate *NP* with nozzles *N* and the plastic plate *P* with communicating channels *Ch*. The design allows several nozzles to be placed side by side, each with its own pressure chamber. Each pressure chamber *PC* is connected to the common fluid-feed channel *C*.

For sufficient pressure to be generated in the fluid there must be no air bubbles in it. The pressure chamber is therefore shaped like a flat cone and has a radial connection to the feed channel *C*. When the pressure chamber is filled, the capillary action of the gap at the edge of the pressure chamber causes the fluid to flow tangentially into the pressure chamber (see *fig. 3*). The two fluid flows meet at the other side of the pressure chamber, so that a volume of air is enclosed in the centre of the chamber. This air is subsequently expelled through the channel *Ch* and the nozzle.

Droplet ejection

When a short rectangular voltage pulse is applied to the electrodes of the piezoelectric plate a pressure wave is created in the fluid. The pressure wave travels through the fluid into the channel *Ch* (see *fig. 2*), so that the fluid at the nozzle *N* is accelerated and a column of fluid is ejected through the nozzle (see *fig. 4*). When the fluid has returned to its initial state, the ejected fluid column contracts and becomes separated from the fluid in the nozzle. The ejected fluid then forms a droplet, whose velocity depends on the energy contained in the voltage pulse.

[1] R. G. Sweet, High frequency recording with electrostatically deflected ink jets, *Rev. sci. Instr.* **36**, 131-136, 1965.

[2] In the rest of this article the term 'fluid' will be used instead of 'ink'.

[3] E. Stemme and S.-G. Larsson, The piezoelectric capillary injector — a new hydrodynamic method for dot pattern generation, *IEEE Trans.* **ED-20**, 14-19, 1973.

[4] J. van Randerat and R. E. Setterington (eds), *Piezoelectric ceramics*, Philips Application-Book, 1974.

Two kinds of energy play a part in the formation of a droplet. When the fluid leaves the nozzle, it contains a quantity of energy. Some of this energy — the surface energy — is used for generating the droplet. The residual energy is the kinetic energy in the droplet.

the area and volume of a sphere and introducing the density. A plot of these relationships is given in *fig. 5*. A value of 50×10^{-3} N/m has been taken for the surface tension (this value relates to ordinary inks at 20 °C) and the droplet velocity is taken as 2 m/s (a

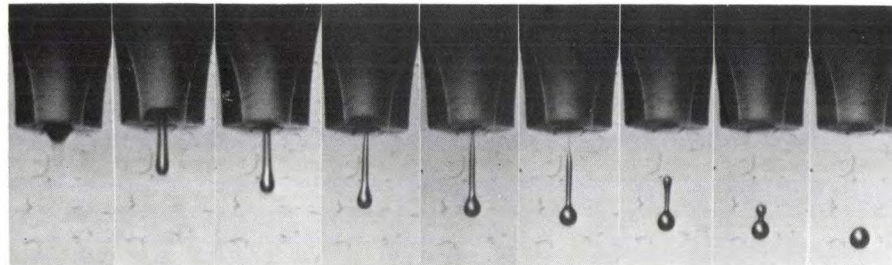


Fig. 4. Ejection of droplets through the nozzle. The photographs were made by stroboscopically illuminating the nozzle at a droplet-ejection rate of about 1000 per second. Although each picture is formed by superimposing about 100 separate images, the definition is sufficient and shows the high stability of the ejection process.

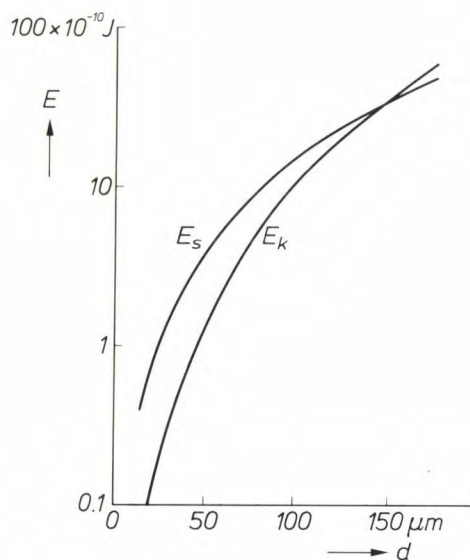


Fig. 5. Kinetic energy E_k and surface energy E_s per droplet, as a function of droplet diameter d . The surface tension is 50×10^{-3} N/m and the droplet velocity 2 m/s. At a droplet diameter of 150 μm the two energies are approximately equal.

practical value). The figure shows that the two energies are identical at a diameter of about 150 μm . However, we want to use droplets with a smaller diameter, determined by the nozzle diameter. This means that each droplet will always have a surface energy greater than the kinetic energy.

It can be seen from *fig. 6* that droplet formation is strongly affected by surface effects. If the emergent fluid wets the area surrounding the nozzle asymmetrically, the droplet is dragged back on the side where the wetting is greatest and is deflected in that direction. In the extreme case the droplet does not break away at

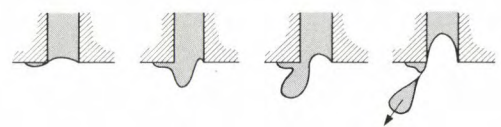


Fig. 6. The effect of wetting the area surrounding the aperture. If the surrounding area is wetted asymmetrically, the droplet is deflected as it leaves the aperture. The drawings are based on photographs made in the same way as those in *fig. 4*.

The surface energy E_s required for forming the droplet surface is

$$E_s = \sigma A,$$

where A is the area of the surface and σ is the surface tension. The kinetic energy E_k of the droplet is

$$E_k = \frac{1}{2} m_d v_d^2,$$

where m_d is the mass of the droplet and v_d its velocity. The energies E_s and E_k can be expressed as a function of the droplet diameter d by using the formulae for

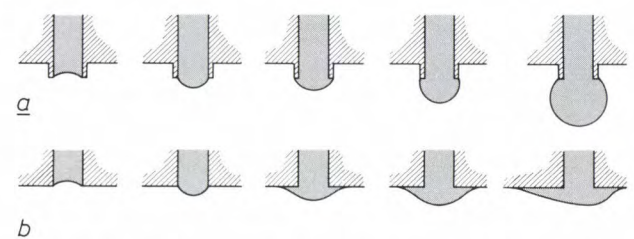


Fig. 7. The modified nozzle with tubular mouth, compared with a nozzle of conventional design. The figure illustrates droplet formation for a low steady fluid pressure. *a)* With the improved design a spherical bulge forms. *b)* With a conventional nozzle the surrounding surfaces are wetted and droplet formation is not symmetrical.

all and the fluid remains behind. It will now be shown that such difficulties with droplet ejection can be prevented by careful attention to the shape of the nozzle.

Design of the nozzle

As we have seen, asymmetrical wetting of the surroundings of the nozzle must be avoided. One sharp edge of 90° is not sufficient to prevent the fluid from spreading over the surface. We therefore designed a nozzle in which two such edges are located closely together, with the aperture of the nozzle in the form of a projecting tube with sharp edges, as illustrated in fig. 7a.

The superior operation of such a nozzle can be seen from a comparison of fig. 7a and fig. 7b. With a steady low fluid pressure, the fluid in our nozzle forms a spherical bulge. For a conventional nozzle, however, the surrounding surface is easily wetted, leading to asymmetrical droplet formation.

Nozzle technologies

If the characters on the paper are to have the desired resolution, the diameter of the apertures of the nozzles must be about 50 µm. To make a nozzle as small as this with the shape shown in fig. 7a will obviously be extremely difficult. We have however developed two technologies for making these nozzles economically.

The stages in the first process are shown in fig. 8a to e. A brass plate has holes drilled in it of diameter greater than the final inside diameter of the nozzle. Next, a layer of nickel is applied by a chemical method; this has the same thickness as the wall of the tube (see fig. 7a). A layer of material of at least the same thickness as the nickel layer is then ground away from the underside of the brass plate. Finally, part of the brass is etched away, producing the desired tubular mouth for the nozzle. Fig. 8f shows a scanning-electron-microscope photograph of a nozzle made in this way.

The second process is illustrated in fig. 9a and b. A spring-steel plate *St* with holes in it larger than the nozzles to be formed is placed underneath a nickel plate. Since the spring-steel plate is only elastically deformed in the process, it can be used several times. A plastic strip *S* is placed underneath the spring-steel plate, and the strip *S* is enclosed by a steel base-plate *BP*. The plastic strip can also be used several times by sliding it out (in a direction perpendicular to the plane of the drawing). A punch tool is located on top of the nickel plate. The punch, which has a tapering diameter, is driven into the nickel plate. Since the spring-steel plate supports the nickel plate and the plastic behaves rather like a fluid, a hole of the desired special

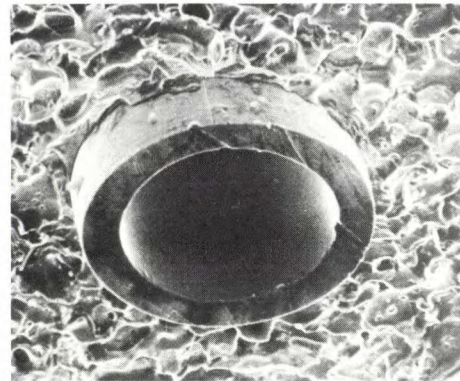
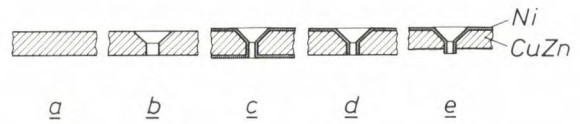


Fig. 8. The chemical process for producing the nozzles. *a)* and *b)*. Holes are drilled in a brass plate. *c)* A layer of nickel is then applied by the 'electroless' method. *d)* The nickel layer, and some of the brass where necessary, is removed by grinding. *e)* Selective etching produces the desired shape of nozzle. *f)* Scanning-electron-microscope (SEM) photograph of a nozzle produced by this process. Magnification about 600×.

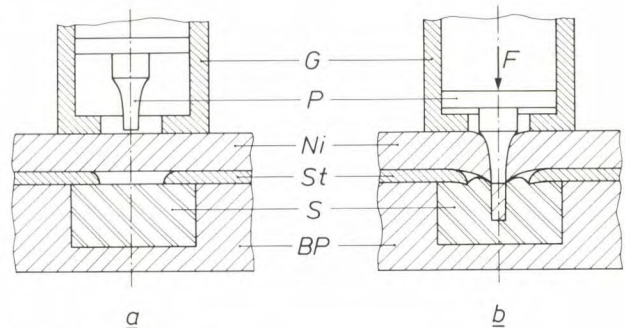


Fig. 9. The mechanical process for producing the nozzles. A spring-steel plate *St* with holes in it is placed beneath a nickel plate *Ni*. A plastic strip *S* is placed beneath *St*, and surrounded by a base-plate *BP*. *a)* A punch tool is located on top of the plate *Ni*; the punch tool consists of a guide *G* and a punch *P*. *b)* The force *F* drives the punch into the plate *Ni*. At the end of the process a part of *Ni* has penetrated into *S*. Because of the supporting action of *St* and the fluid behaviour of *S* a hole without burrs and of the desired shape is produced in *Ni*. *c)* SEM photograph of a nozzle produced by this process. Magnification 600×.

shape and free of burrs is punched in the nickel plate. A nozzle made by this method is shown in fig. 9c. The sharp edge inhibits the wetting of the surrounding surface even more than the tube-like end produced by the previous method. Since no time-consuming drilling is required, the nozzles produced in this way are even cheaper.

Practical design of the print head

In the printing process the print head with its nozzles is moved across the paper so that it prints line by line. The desired resolution of the characters determines how many nozzles are required in the print head [5]. It is not usually possible to put all the nozzles in a single row, for geometrical reasons. Fig. 10 shows how the printing is done when the characters have a height of twelve dots or droplets. The nozzles are positioned in two rows of six, offset from each other by half the spacing in the row (see fig. 10a). During the printing, the appropriate piezoelectric plates in the

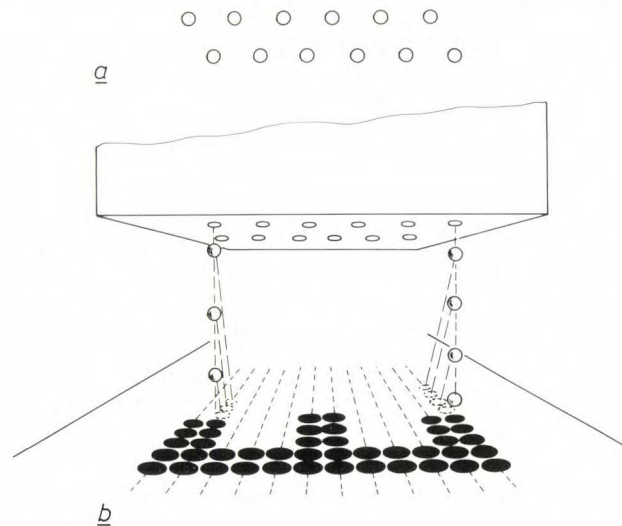


Fig. 10. Ejection of the droplets by a print head with 12 nozzles. *a)* Location of the nozzles in two rows of six. *b)* Formation of characters twelve dots high. Piezoelectric plates in the two rows of nozzles are energized in succession, depending on the particular character required. The droplets strike the paper at an angle equal to that of the resultant of the ejection velocity and the print-head velocity.

pressure chambers of the two rows are energized in succession, depending on the particular character required.

Fig. 11 shows the P2131 printer developed by Philips from the results of our investigations, for use with the P2000 microcomputer. Fig. 12 gives a typical print-out obtained with this machine, showing the three different typefaces. The printing speed is 80 standard characters per second (one line per second). Since a standard character consists of a matrix of 10×12 dots, this corresponds to a rate of approximately 1200 drop-

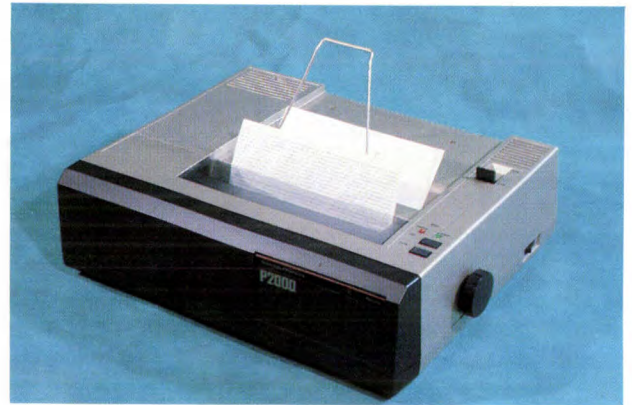


Fig. 11. The P2131 printer developed for the Philips P2000 micro-computer. Its print head has 12 nozzles and prints characters in a matrix of 10×12 dots.

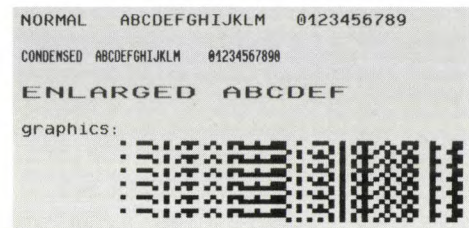


Fig. 12. Examples of the three different typefaces that can be printed with the P2131 machine.

lets per second. We shall see that this rate can be made higher.

For professional purposes a higher resolution is required. We have therefore developed a print head with 24 nozzles in four rows. A photograph of this head can be seen in fig. 13. The plate with the piezoelectric elements has been removed, so that the 24 conical pressure chambers are visible; the connections

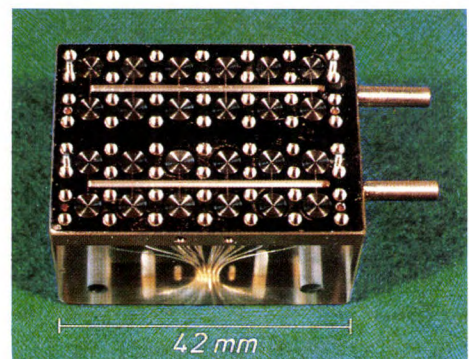


Fig. 13. Photograph of the experimental print head with 24 nozzles that can print characters in a matrix of 24×20 dots. The conical pressure chambers can be seen at the top (the plate with the piezoelectric elements has been removed) and the connections for the ink feed are on the right. The transparent plastic body contains the channels to the nozzles.

for the ink feed can be seen on the right. The channels that connect the pressure chambers with the nozzles are visible in the plastic body.

For printers of this type it is important to have the right combination of ink and paper. The ink in the nozzle must not dry out, but the paper must dry quickly enough to prevent smearing of the characters. A compromise can be found in a type of ink that attracts moisture from the air and therefore does not dry too quickly. Smearing of the characters must then be prevented by using a type of paper that absorbs the ink quickly. Paper with a high ash content is a fast ink absorber, but its coarse fibres give large and blurred dots. Paper with a high sizing content, on the other hand, forms well-defined dots, but does not absorb the ink quickly. The choice of the correct type of paper is therefore a compromise, in which the type of ink is also a factor.

Increasing the printing speed

After a voltage pulse has been applied to the ceramic plate, the pressure in the channel at the nozzle can vary in several ways as a function of time (see fig. 14). The degree of damping determines the time taken for the system to return to rest. If the system is critically damped (curve 1), the time taken is shorter than if it is underdamped (curve 2) or overdamped (curve 3). The next voltage pulse cannot be applied to the ceramic plate until the system is almost completely at rest. For a given resonant frequency of the system the maximum droplet ejection rate n is thus highest when the system is critically damped. Depending on the viscosity of the fluid, critical damping can be achieved by locally narrowing the diameter of the channel between the pressure chamber PC and the feed channel C (see figs 2 and 3). This also has the effect of reducing mutual interference between adjacent pressure generators.

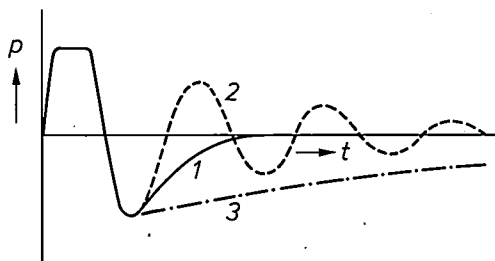


Fig. 14. The variation of the pressure p in the channel at the aperture, after application of a voltage pulse to the ceramic plate, as a function of time t . Curve 1 critically damped, curve 2 underdamped, curve 3 overdamped. For a given resonant frequency of the system the highest droplet-ejection rate is achieved with critical damping.

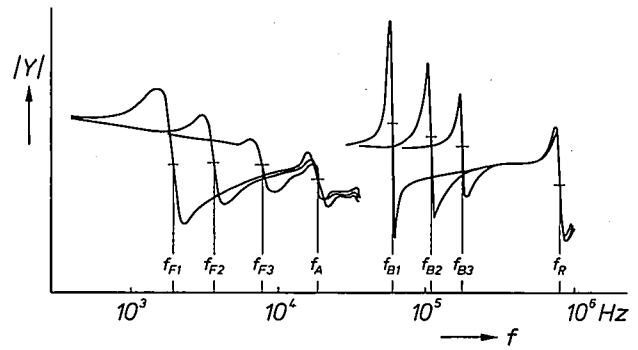


Fig. 15. Measured variation of the modulus of the admittance (the reciprocal of the impedance) of the piezoelectric plate in fig. 1, as a function of the frequency of an applied alternating voltage. On the right are three frequency characteristics of the bilaminar plate alone, for three different plate thicknesses, increasing towards the right. The quantities f_{B1} , f_{B2} and f_{B3} are the lowest natural frequencies of the flexure vibrations. The frequency f_R is the natural frequency of the radial vibrations, and is independent of the plate thickness. The three curves on the left were measured for the complete system filled with fluid, with the same variation of the plate thickness. The quantities f_{F1} , f_{F2} and f_{F3} are the lowest natural frequencies of the flexure vibrations when the system is filled with fluid, since $f_{B1} : f_{B2} : f_{B3} = f_{F1} : f_{F2} : f_{F3}$. The lowest frequency f_A of the acoustic vibrations is independent of the plate thickness. The wavelength associated with f_A is approximately equal to twice the length of the channel Ch in fig. 2.

At the same time, to achieve the maximum ejection rate the lowest resonant frequency of the system should be as high as possible. This depends to a great extent on the resonant frequency of the bilaminar plate in the pressure generator (see fig. 1). Three frequency characteristics of the bilaminar plate alone can be seen at the right-hand side of fig. 15. For a bilaminar plate whose thickness is increased in three steps, f_{B1} , f_{B2} and f_{B3} are the lowest resonant frequencies of the flexure vibrations of the type shown in fig. 1. In these three cases the frequency f_R for radial vibrations remains constant. The left-hand side of fig. 15 shows, for the same variation in plate thickness, the three measured frequency characteristics for the complete system filled with fluid. The lowest resonant frequencies f_{F1} , f_{F2} and f_{F3} originate from the flexure vibrations of the bilaminar plate. The frequency f_A does not depend on the plate thickness and arises from the acoustic vibration with the lowest frequency, the wavelength being roughly equal to twice the length of the channel Ch (see fig. 2). If we increase the flexure frequency f_F by increasing the stiffness of the ceramic plate, f_F approaches the frequency f_A of the acoustic vibrations. The acoustic vibrations determine the dynamic behaviour of the system if the stiffness of the plate is increased further. Eventually a reduction of the length of the channel Ch will give an increase in f_A .

[6] See for example J. Borne, Philips tech. Rev. 29, 205, 1968.

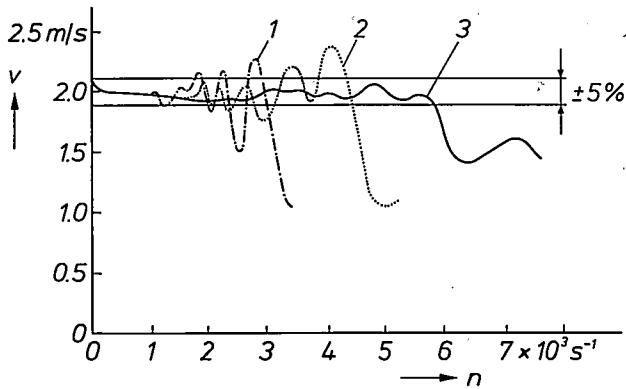


Fig. 16. Velocity v of the droplets generated by a print head as a function of the ejection rate n . To obtain satisfactory printing, v must not vary by more than about 5%. Curve 1 relates to the print head in fig. 13. Curve 2 is obtained when the stiffness of the bilaminar plate is increased by reducing its diameter. Curve 3 applies when the length of the channel to the nozzle is also reduced. A print head constructed to meet these conditions gives twice the maximum droplet-ejection rate and its volume is reduced by a factor of six.

The extent to which the maximum droplet-ejection rate can be increased in this way is illustrated in *fig. 16*, where the velocity v of the ejected droplets is plotted against the ejection rate n . When n approximates to a resonant frequency of the system the ejection rate is increased or decreased. In practice it is found that the

ejection rate should not vary by more than about 5% from the mean value. Curve 1 shows the variation of the ejection velocity of the print head in *fig. 13*. Curve 2 is the characteristic of the same system, but with a smaller diameter for the bilaminar plate, which therefore has increased stiffness. Curve 3 is the characteristic of the system, but now with the length of the channel to the nozzle reduced as well. The graph shows that the ejection rate can now be increased to 6000 droplets per second, which is twice the rate previously considered possible with DOD printers. An incidental advantage of the technology is that the volume of the print head is reduced by a factor of about six.

Summary. A print head of sandwich construction makes it possible to impel ink droplets at the paper in direct response to control signals (the 'droplet-on-demand' (DOD) principle). Special nozzles have been designed that eject the droplets with great directional accuracy, thus improving the quality of the printed characters. Two technologies have been developed for producing these nozzles economically. Some practical print heads of this type are discussed; one of these is used in the P2131 printer developed for the Philips P2000 microcomputer. Improvement of the dynamic characteristics of the printing system allows the droplet-ejection rate to be increased to 6000 droplets per second, which is twice the rate previously considered possible with DOD printing systems.

Three special applications of the Philips high-speed spark-machining equipment

J. L. C. Wijers

For some years the practical potentialities of 'spark erosion' have been the subject of extensive investigations at Philips. The article below gives some instances of this activity and shows that high-speed (micro)spark machining promises to develop into an advanced general-purpose machining technology. The examples given have been chosen to illustrate what can be done with a single relatively simple type of spark-erosion machine. The high-speed spark-machining equipment designed and built at the Research Laboratories — and the result of the combined effort of many people — draws attention once again to the advantages of efficiency and economy that can be achieved with a new general-purpose workshop technology.

Spark erosion

Towards the end of the sixties the first machines were developed at Philips Research Laboratories for machining materials by high-speed (micro)spark erosion. The principle of this workshop technology — the melting or vaporization of workpiece particles in a localized spark discharge and removal of the eroded material by flushing it away with a dielectric fluid — is now generally known; the spark discharge takes place between a moving electrode and the workpiece itself, which acts as the second electrode. The method has been described in detail in an earlier article in this journal ^[1]. Since then all kinds of machining problems of a rather experimental nature, originating from widely divergent areas of research and engineering, have been regularly tackled with our high-speed spark-erosion equipment and usually brought to a satisfactory conclusion. In this way our laboratories have gained a wealth of experience with spark machining, or 'sparking' as we often call it in the machine shop, as a method of machining materials. The method has been particularly useful for 'internal' work, i.e. for making cavities, channels or other recesses inside a block of material.

The 'tool electrode', which produces a 'negative' copy of itself in the material by spark erosion, is made by an *external* operation on the electrode material. Conventional workshop methods such as turning, milling and grinding are generally more suitable for external work than for internal work. Blind holes in particular give many problems when conventional methods are used. Spark machining provides a 'translation' from a relatively easily made external shape to a difficult internal shape.

Spark machining is often the only way of making shapes with accurate rectangular corners, without significant rounding. Another important advantage of spark machining is that there is no direct contact between the tool and the workpiece, so that no mechanical stresses are introduced. It is almost impossible to machine single-crystal aluminium, for example, in any other way without the danger of it becoming polycrystalline.

Fig. 1 gives an idea of what can be achieved with high-speed microspark machining, from data relating to the quality of seven combinations of workpiece and electrode materials now widely employed. 'High-speed' implies the removal of material in quantities of a few cubic millimetres per minute. At these rates a

Ing. J. L. C. Wijers is with Philips Research Laboratories, Eindhoven.

^[1] C. van Osenbruggen, High-precision spark machining, Philips tech. Rev. 30, 195-208, 1969.

surface roughness of less than 1 μm can be achieved [2]. The exact values depend on the particular combination of materials. Table I and the photographs present some typical features of the machine used. In the design of this new high-speed spark machine the experience gained in the last few years has of course been used to the full in making the optimum compromise between machining rate and

inching problem that would have been most difficult with conventional methods. This miniature tool — it is less than 2 mm long — is used for thermocompression bonding the electrical connections of an individual monolithic integrated circuit (IC). High-precision bonding with these bits is of course also intended to facilitate the automation of the assembly of ICs in their electronic 'environment'.

QUALITY FIGURE	ELECTRODE MATERIAL						
	ts	W-Cu	W	W-Cu	tc	W-Cu	Mo
- speed $\cdot \text{mm}^3/\text{min}$	2.05	4.10	0.59	1.15	2.75		0.02
- relative electrode wear %	100	11	4.6	7.2	11.3	10	
- roughness R_a μm	1.4	1.4	1.5	0.6	0.8	1.4	0.8
	ts		tc			C	C
SETTING DATA	hardened tool steel		tungsten carbide			diamond synthetic (sintered)	diamond natural
V (volts)	400	400	200	200	400		800 spec. dielec.
A (amperes)	4.0	4.0	20	20	4.0		
t (μs)	1.6	1.6	0.4	0.4	0.8		

Fig. 1. Seven suitable combinations of workpiece and electrode materials for use with the machine described in Table I. Each box gives the quality figures achieved. The figure in the top left-hand corner of a box gives the erosion rate in mm^3/min , i.e. the rate at which the workpiece material can be sparked away; the figure at the lower right is the surface roughness in μm ; the figure in the middle is the relative electrode wear as a percentage, i.e. the loss of electrode material as a percentage of the loss of workpiece material. Also given are the machine setting data V , A and t : these are the no-load generator voltage, the r.m.s. pulse current, and the pulse duration, respectively. In addition to the combinations of workpiece and electrode materials given here, there are of course others that also yield good results, e.g. the tool steel mentioned here with electrolytic copper, and tantalum with tungsten.

machining accuracy. Working with the machine has strengthened our view that microspark machining can be developed into a general-purpose method of metal machining.

This article will demonstrate some of the accomplishments of our highly skilled 'master-of-all-trades' by discussing two widely different workpieces and one rather unusual machining problem. The section that follows deals with the machining of a 'bonding bit' for joining together extremely small parts (0.01 mm) by thermocompression bonding. Next it is shown how spheres of single-crystal aluminium were made for certain material investigations. The final section deals with the machining of diamond for tools.

Miniature bits for thermocompression bonding

The manufacture of the IC bonding bit to the design shown in fig. 2 is our first example here of high-speed spark machining successfully applied to a mach-

Integrated circuits, especially those for signals at very high frequencies, have little mechanical strength, and therefore it is much too risky to make connections by soldering or thermocompression bonding on the circuit itself. For external connections these ICs have a number of projecting gold strips ('beams', about 0.5 mm long, 50 μm wide and 10 μm thick), and an IC bonding bit can be used for joining the ends of the beams to the contact pins of the adjoining circuits — thus avoiding all danger of damaging the IC itself. The bonding edge (E in fig. 2), which is only 35 μm thick, forms a square into which the IC fits comfortably. When the bonding bit is moved this edge presses the ends of the beams accurately on to the contact pins. A brief application of pressure on the edge (about 0.4 N), with a local temperature increase of 400 $^{\circ}\text{C}$, results in a firm bond between beams and pins.

Making a miniature tool of this type is an unusual problem if only because of the extremely small dimensions. Workpieces measuring about 1 mm are very difficult to set up and machine when conventional methods such as drilling or milling are used. While the bit is being machined, a pressure only slightly greater

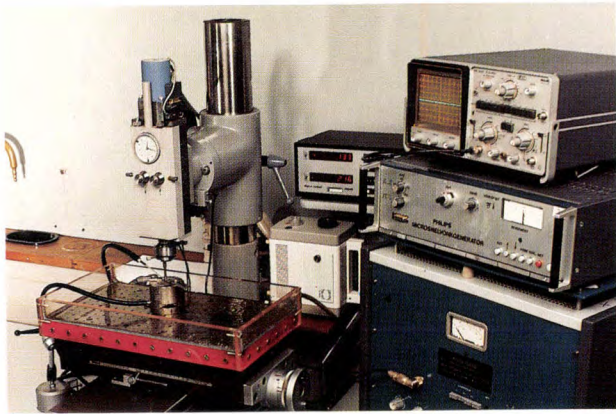


Table I. Some typical data for the Philips high-speed spark-machining equipment (upper photograph). This machine was designed and built in our own laboratories; it differs from the 'die-sinking' type commercially available mainly in having a smaller active electrode surface, a shorter pulse duration, a higher spark frequency and in the use of deionized water as the dielectric fluid and flushing agent. The lower photograph shows a later model of the electrode head of the machine with servo control and spark generator type EDM 81. Quality figures and machine-setting data are shown in fig. 1.

Spark capacity (sinking ^[a] , wire sparking ^[b] , rotary sparking, screw-cutting, copying)	max. 200 mm (x) max. 160 mm (y) max. 240 mm (z)
Clamping area	430 mm × 215 mm
Active area (electrode)	max. 80 mm ² min. about 0.01 mm ²
Generator voltage (no load)	200, 300, 400, 800 ^[c] V
Pulse current (r.m.s.)	4, 20, 40 A
Pulse duration	0.2-3.2 μs
Repetition rate	2-250 kHz
Polarity	workpiece positive
Dielectric fluid	H ₂ O, deionized, conductivity 1 mSm ⁻¹
Electrode servomechanism	electromechanical (stepping motor)
Smallest displacement	½ step = 0.625 μm

^[a] In the 'die-sinking' process the electrode moves vertically (in the z-direction) and 'sinks' into the workpiece.
^[b] In 'wire sparking' the electrode is a thin metal wire stretched over two guide rollers. While it is sparking the wire 'saws' through the workpiece and at the same time moves over the rollers (so that the electrode is continuously replenished).
^[c] This value is only attainable if a special transformer is used.



than it will encounter in thermocompression bonding could permanently damage its bonding edge, which is weak because of its small thickness and relatively large height. What is more, the inside surface of the workpiece requires a great deal of machining (fig. 2) to produce the sharp, rectangular transitions. Spark machining was the obvious method here.

In making the tool we started with small steel balls, as used in bearings. The steel of these balls readily meets a number of the requirements usually necessary for bonding bits. The accurately spherical shape and the smoothness are added advantages, since part of the surface can be used as the upper face of the bit without further machining. The spherical shape and smoothness permit a measure of self-alignment, automatically stabilizing the vertical position of a bit when pressure is applied during the bonding.

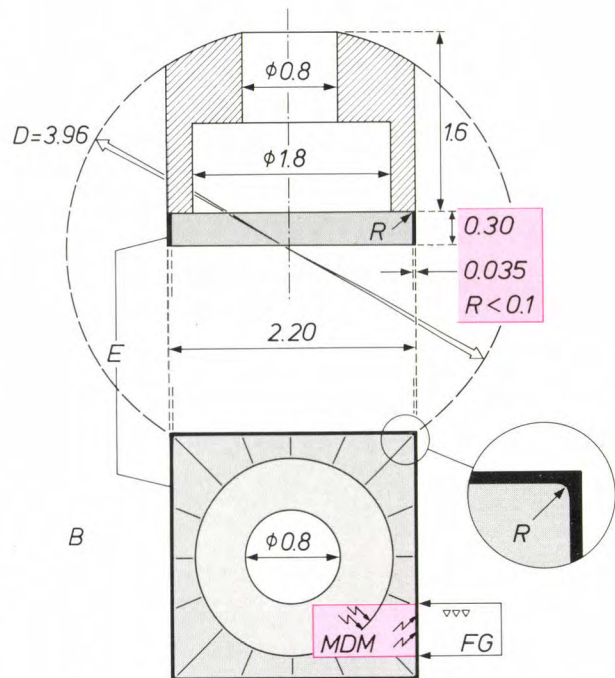


Fig. 2. Dimensional drawing (vertical section, B view from below; measurements in mm) of an IC bonding bit. *E* bonding edge, height 0.30 mm and width 0.035 mm. *R* radius of curvature of the right angles formed by *E*. This miniature tool was made from a steel bearing-ball of diameter *D*. Because of the great deal of internal work, the nearly rectangular corners, and above all the two small dimensions of the bonding edge, microspark machining (sometimes called microdischarge machining, *MDM*) is the only suitable machining technique. The outer face of *E* is ground (*FG*). The total machining time is about 25 minutes.

^[2] The roughness values given in this article are *R_a*-values, which are usually about a quarter of the corresponding *R_{max}*-values (peak-to-valley). A definition of these values will be found in the 'Tool Engineers Handbook' (ASTE), McGraw-Hill, New York 1959, and in the provisional standard ISO/DIS 468 (1980).

Aluminium test spheres

The second manufacturing problem we shall look at here is spark machining of spheres from a rod of single-crystal aluminium. This problem was encountered at the Research Laboratories when it was necessary to find out whether the result of a particular electroplating process depended on the crystallographic orientation of the aluminium used. Small spherical samples were required for the investigation. We started with cylindrical rods of the material made in our laboratories, each rod consisting of a single crystal of aluminium. The sparking hardly affected the single-crystal nature of the material^[3].

Fig. 3 illustrates how we solved the problem of machining the rod to leave a small sphere with a stem. The original rod, fixed in a holder that can rotate about its axis, acts as a rotary anode. The speed of rotation is about 50 revolutions per minute. The axis of rotation has a fixed position in the horizontal plane. The cathode is a tube with a groove cut into one side, and can be displaced vertically. (The width of the groove determines the final diameter of the stem.) The servomechanism of the machine ensures that the sparking edge of the cathode always remains at the same distance from the anode. The rotation of the rod ensures that its entire circumference is exposed to the sparking.

The combined fast rotation of the rod and the slow vertical movement of the cathode eventually leaves a perfect sphere. The entire process takes only an hour. The diameter of the sphere is approximately 0.1 mm smaller than the inside diameter of the cathode, and its roughness is about 1 μm . The wear of the cathode does not adversely affect the dimensional accuracy, provided that the cathode is *long* enough; the active length of the cathode should be greater than is necessary for sparking the complete sphere.

Spark machining of diamond

The idea of machining diamond, the hardest known material, by spark erosion looks at first sight unrealistic. With the exception of a very rare semiconducting variety, diamond is an excellent insulator. However, it is in fact suitable for spark machining, but ingenuity is necessary. A thin coating of adhesive containing powdered silver is applied to the surface of the diamond. The silver makes the coating conduct electrically, and the spark-machining process can be started. When the coating has been eroded away the process does not stop, but continues. The theory of this effect is that after the removal of the top layer the very high temperature (a few thousand kelvins) in the discharge channel continuously converts the newly ex-

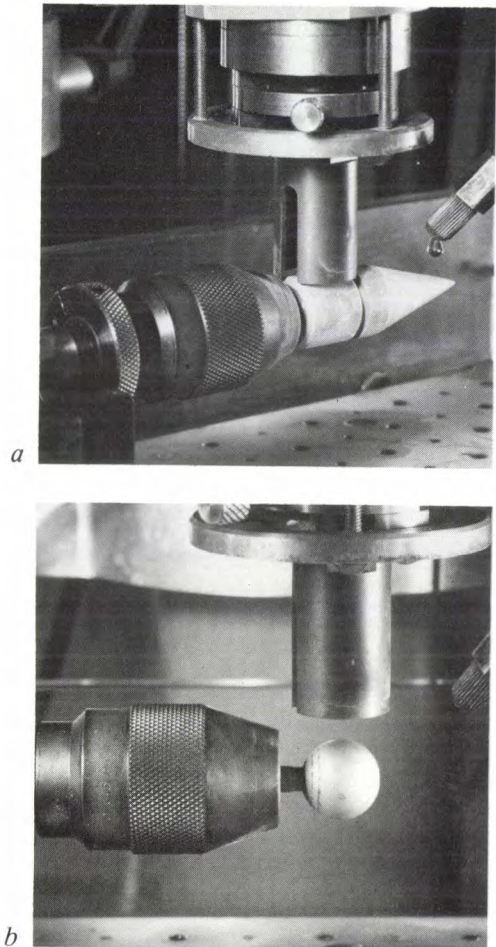


Fig. 3. Making aluminium test spheres (diameter about 17 mm) by rotary sparking. The spheres are formed from a rod of single-crystal aluminium. The electrode has a groove cut in it, so that the test sphere finishes up at the end of a stem. *a*) Intermediate stage. Sparking produces no significant stresses in the material, which therefore remains of a single-crystal nature. *b*) The final situation.

posed material into a graphitic form — which is a good conductor.

Another method for starting the spark-machining process is to heat the diamond material in a carbon-rich atmosphere to a temperature high enough (about 1100 K) to produce a sufficiently conducting graphitic structure on the surface.

If the spark machining of a diamond and the removal of the eroded material from the discharge gap are to proceed in a stable manner, it is important to use the correct dielectric fluid (the 'sparking fluid') in the gap between electrode and workpiece. A concentrated aqueous solution of sodium hydroxide in paraffin (kerosene) with certain additives proves satisfactory. The no-load voltage of the spark generator should be at least 800 V to produce spark pulses, of about 1 μs duration, that have sufficient energy.

The knowledge that diamond can be spark machined — and with great accuracy — has meanwhile acquired considerable technological significance. At Philips tools are now widely used in which an extremely accurately dimensioned natural or synthetic diamond determines the quality of the work. Familiar examples are diamond-tipped tools for machining optical surfaces^[4] and cutters for producing smooth fractures in optical glass fibres^[5]. Diamond is also widely used for wire drawing. Instead of natural diamond the material now chosen for the actual drawing die (fig. 4) is often a composite of a metal alloy containing a polycrystalline form of synthetic diamond.

Making a diamond-tipped precision tool^[4] by conventional methods is an exceptionally laborious and highly specialized job, largely because of all the grinding, and can take at least 50 hours. The high precision that has now been reached in the application of high-speed spark machining is such that the total machining time can be considerably curtailed; at present grinding is necessary only in special cases as a finishing operation. In this way it takes only a few hours to make the complete tool — only about a tenth of the time previously necessary.

A diamond of 1 to 3 carats is used for a precision tool. A chip about 1 mm thick is 'sawn' by spark machining: the electrode is a molybdenum wire of diameter 0.2 mm. The flat chip thus obtained is brazed to the shank of the tool at a temperature of 1300 K. The desired shape of the cutting edge (and of the flank) is obtained by copying a standard shape, again by spark machining with a wire electrode; this takes about 30 minutes. The resultant surface is finished with a polishing operation. For minimum wear of the tool the primary face of the chip, its top surface, should have the same orientation as one of the crystal planes of its cubic crystalline structure, within a tolerance of a few degrees. This orientation is determined beforehand by X-ray diffraction.

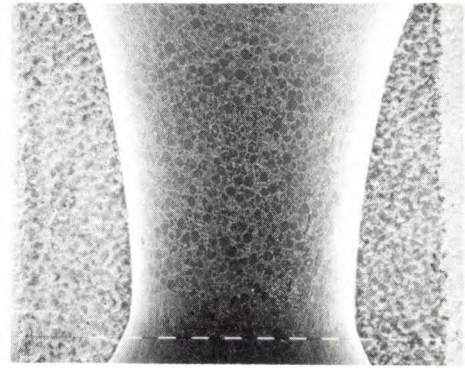


Fig. 4. Photomicrograph of a cross-section of a diamond die. The bell-shaped channel, designed for drawing wire of diameter 1.3 mm, is produced by means of a sparking process. The diamond is a polycrystalline synthetic product. Each dimensional step is 100 μm . The die was sawn through — by wire sparking — to enable roughness measurements to be made and to give a better view. The picture was made with a scanning electron microscope (SEM).

The high-speed spark machine shown in the photograph above Table I can also be useful for 'conditioning' grinding wheels that have lost some of their abrasive power, especially wheels in which the diamond grains that serve as grinding material are contained in a layer of metal, e.g. bronze. The conditioning is sometimes done in two stages. Here the high-speed machine is ideally suited for the second stage, the finer finish. The design of the machine, with removable electrode head, offers the added advantage that the grinding wheel does not have to be taken out of the grinding machine, so that the rotational accuracy of the disc is more easily preserved. In such situations the dielectric fluid is simply sprayed between the electrode and the grinding wheel.

Summary. Three practical cases are described to illustrate the usefulness of high-speed (micro)spark machining as a general-purpose workshop technology. They are the manufacture of a bonding bit for connecting up individual ICs by thermocompression, the manufacture of single-crystal aluminium spheres from a rod of the same material for investigating electroplating behaviour, and the machining of diamond for tools. The high-speed spark-machining equipment built at our laboratories for such machining and manufacturing problems and its principal features are briefly described. The machining rate is a few mm^3 per minute, the attainable surface roughness (R_a) is less than 1 μm . Diamond can be spark-machined after providing it with a conductive layer. With an appropriate spark-generator voltage and dielectric fluid, spark erosion continues after the conductive layer has been eroded away. The machine can also be used for conditioning diamond grinding wheels.

^[3] In the electroplating process no dependence on the crystallographic orientation of the aluminium was found.

^[4] T. G. Gijssbers, COLATH, a numerically controlled lathe for very high precision, Philips tech. Rev. 39, 229-244, 1980. J. M. Oomen of these laboratories has made a substantial contribution to the further development of the method of manufacturing diamond-tipped precision tools.

^[5] A neat illustration of the smooth fracture of a glass fibre accurately perpendicular to its axis by a rotating diamond cutter can be found in Philips tech. Rev. 39, 245, 1980.

Chemical vapour deposition of wear-resistant coatings on tool steel

P. J. M. van der Straten and G. Verspui

The process of covering a substrate with a thin coating from a chemically reactive mixture of gases is often referred to as CVD ('chemical vapour deposition'). The CVD technology is widely used in the electronics industry for the manufacture of integrated circuits. An entirely different application that has recently become the subject of growing interest is the coating of steel tools with a wear-resistant layer. The results of a study started four years ago at the Philips Centre for Manufacturing Technology show that this can appreciably increase the durability and reliability of tools.

Introduction

A material widely used for metal-cutting tools is tool steel, because it is readily workable in the soft state and certain types, after hardening, possess good mechanical properties, such as great hardness and toughness. In many cases, however, tool steel shows an undesirable amount of wear. With the rising costs of labour and materials it is very much in the interest of industry to reduce this excessive wear, thereby lengthening the life and increasing the reliability of tools so that they can be used longer in production processes without giving trouble.

The wear resistance of a number of steels can be improved by means of diffusion processes such as carburizing and nitriding. If the tool steel treated in this way is still not sufficiently wear-resistant, preference is often given to cemented carbides, a group of materials consisting generally of powdered metal carbide and a metal, such as cobalt, as the binding agent. These materials, however, are relatively expensive and as a rule are much more difficult to work than tool steel.

Other materials in solid form cannot usually meet all the requirements. For instance, although borides, carbides and nitrides of the transition metals are hard and resistant to both wear and corrosion, their great brittleness rules them out as a material for tools.

Nevertheless, good use can be made of these materials in tool engineering when they are applied as a thin coating on tool steel. A good method of applying such a coating is to deposit it from a chemically reactive gas mixture. In this chemical vapour-deposition method (CVD for short), the gaseous compounds that contain the elements from which the material to be deposited is composed are allowed to react with each other at high temperature. Under suitable conditions a compact layer is then deposited on a hot substrate, in this case a substrate of tool steel^[1].

In recent years growing interest has been shown in this application of CVD^[2], because tools can be made in this way that combine the desirable properties of tool steel with those of the coating. Under mechanical loads a tool with a hard surface layer produced in this way shows little wear, while the tool-steel substrate takes up the forces acting upon it.

Thin coatings can also be applied by other methods such as evaporation, sputtering and plasma deposition. In these cases the substrate usually has a lower temperature than in the CVD process. One of these methods will generally be preferred for coating substrates that cannot tolerate a high temperature. They can also be used for depositing coatings that are not only wear-resistant but also give good 'lubrication' because of their low friction^[3]. However, if the requirements for the hardness and adhesion of the coating are exacting, the CVD process has distinct advantages.

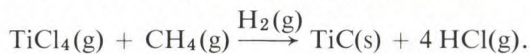
Dr Ir P. J. M. van der Straten and G. Verspui are with the Philips Centre for Manufacturing Technology, Eindhoven.

About four years ago an investigation was started at the Philips Centre for Manufacturing Technology (CFT) into the use of the CVD technology for the improvement of tools. In the course of this work a great deal of experience has been gained with the deposition of titanium-carbide (TiC) and titanium-nitride (TiN) coatings. Questions studied include the dependence of the properties of the coating material on the conditions during deposition, the types of tool steel that are most suitable for a CVD coating, and the best shapes for the tool-steel substrate. Attention has also been paid to potential applications of TiC or TiN coatings on tool steel.

In this article we shall first briefly describe the CVD process and discuss a number of properties of the TiC and TiN coatings. We shall then deal with the effects of the CVD process on the tool-steel substrate. The coating of complex substrate shapes will be dealt with in connection with a number of results obtained with the coating of TiC and TiN on the inner surface of narrow tool-steel tubes. We shall conclude with a discussion of some likely applications.

Description of process and properties of the coatings

Fig. 1 shows a diagram of a CVD installation for depositing TiC and TiN coatings on tool steel. The TiC is formed by the reduction of gaseous titanium chloride (TiCl_4) in the presence of methane (CH_4), using hydrogen (H_2) as reducing agent and carrier gas:



In practice this reaction is usually made to take place at a relatively low pressure (about 90 mbar) and at a

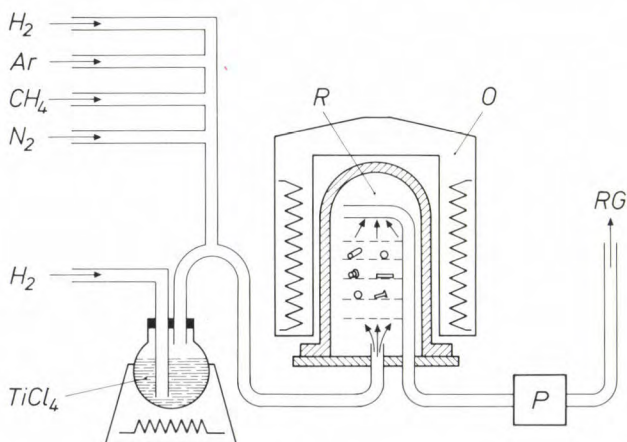


Fig. 1. Schematic arrangement of a CVD installation for applying a TiC or TiN coating. The workpieces to be coated are situated in the reactor *R*, which is kept at the correct temperature in the furnace *O*. The gas flow in the reactor consists of hydrogen (H_2), methane (CH_4) or nitrogen (N_2), and titanium chloride (TiCl_4) obtained by heating liquid TiCl_4 . The residual gases *RG* are extracted by a pump system *P*. For safety, the reactor is filled with argon (*Ar*) before it is opened at the end of the process.

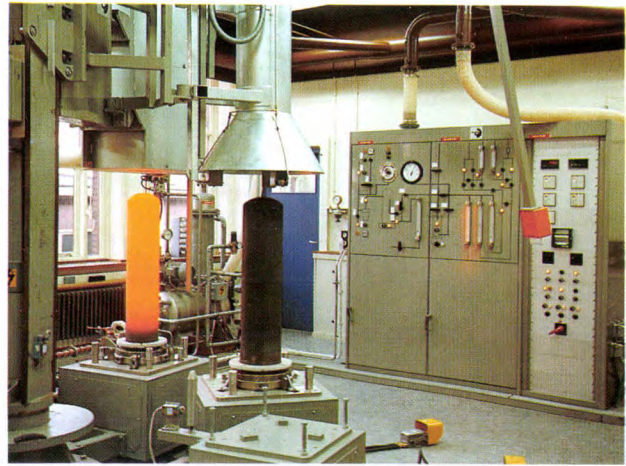
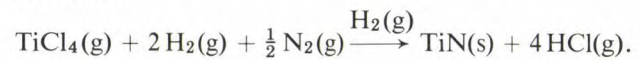


Fig. 2. Photograph of an industrial CVD installation (manufactured by Bernex) at CFT.

temperature of 950 to 1050 °C. The TiC coating has an appearance ranging from bright silver to dull grey, depending on the kind of steel used for the substrate. A yellowish-gold TiN coating is formed in a similar way, but nitrogen (N_2) is then used instead of methane:



This reaction is usually made to take place at almost atmospheric pressure (about 900 mbar) and at a temperature of 850 to 950 °C. The hydrochloric-acid gas (HCl) liberated in the above reactions is transported by the carrier gas to a liquid ring pump filled with caustic-soda solution, which neutralizes it.

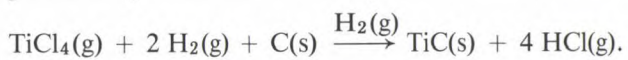
Fig. 2 shows a photograph of an industrial CVD installation, as set up at CFT for coating with TiC and TiN. Normally the coating is allowed to grow at a rate of 1 to 2 μm per hour to a thickness of 2 to 10 μm . There is not usually much point in applying a thicker coating, because the desired increase in resistance to wear and corrosion can be obtained with thin coatings. It is also necessary to take account of the difference in thermal expansion coefficient of tool steel (10 to $13 \times 10^{-6} \text{K}^{-1}$) and that of TiC ($7.5 \times 10^{-6} \text{K}^{-1}$) or TiN ($9.3 \times 10^{-6} \text{K}^{-1}$). This sets up a stress in the coating when the tool is cooled to room temperature. Thicker coatings may then crack or even separate from the substrate. Coatings 2 to 10 μm thick can take up this stress and can even withstand some impact or shock. TiC and TiN coatings of thickness up to about 10 μm adhere well to tool steel.

[1] C. F. Powell, J. H. Oxley and J. M. Blocher Jr., Vapor deposition, Wiley, New York 1966.

[2] H. E. Hintermann, Tribology Int. 13, 267, 1980.

[3] An article on this subject by H. Dimigen and H. Hübsch will appear in a forthcoming issue of Philips Technical Review.

The CVD coating usually gives only a slight increase in the surface roughness: 0.1 to 0.3 μm . When TiC is deposited on tool steel with less than 0.4% carbon, however, the increase in roughness can be much greater. This is connected with the particular part played by the carbon in the steel during the chemical vapour deposition of TiC. A large amount of carbon is extracted from the steel for the formation of TiC, particularly at the start of the process:



Excessive decarbonizing of the surface region of the workpiece is avoided by diffusion of carbon from the bulk of the tool steel [4]. To ensure that the coating will have a fine structure and an acceptable growth rate, the steel should normally contain at least 0.4 per cent of carbon. If it does not, the methane concentration should be increased in the vapour phase. As can be seen in *fig. 3*, the TiC coatings then obtained are appreciably rougher.

The most notable property of the coatings is their great hardness: 3500-4000 HV (Vickers) for TiC and 2000-2500 HV for TiN [5]. The hardness of TiC is only exceeded by that of silicon carbide (SiC), boron carbide (B_4C) and diamond. Because of their great hardness TiC coatings offer good protection from abrasion. Another notable property of TiC and TiN is their high melting point ($> 2500^\circ\text{C}$). In addition, their high chemical stability, the low solubility of metals in the coatings and their low coefficient of friction give TiC and TiN coatings a high resistance to 'adhesive' wear. This means that these coatings have little tendency to 'galling' and 'fretting' and to become eroded. In this respect, TiN coatings are somewhat better than TiC coatings. The coatings are also highly resistant to oxidation by oxygen in the air, even at high temperatures: up to 400°C for TiC and to about 500°C for TiN [5].

The substrate

There are many kinds of tool steel. The choice will depend on the requirements to be met by the tool and on the conditions during the CVD process. Since the workpiece to be coated has to remain at a high temperature for a relatively long time, there is a danger of grain growth in the steel and of stabilization of the austenite phase (the high-temperature phase of steel that should change to the martensite phase on cooling).

After the CVD process the reactor is cooled in a hydrogen atmosphere at a rate of 20 to 50°C per minute. This implies that types of steel that are hardened by quenching in water or oil are nowhere

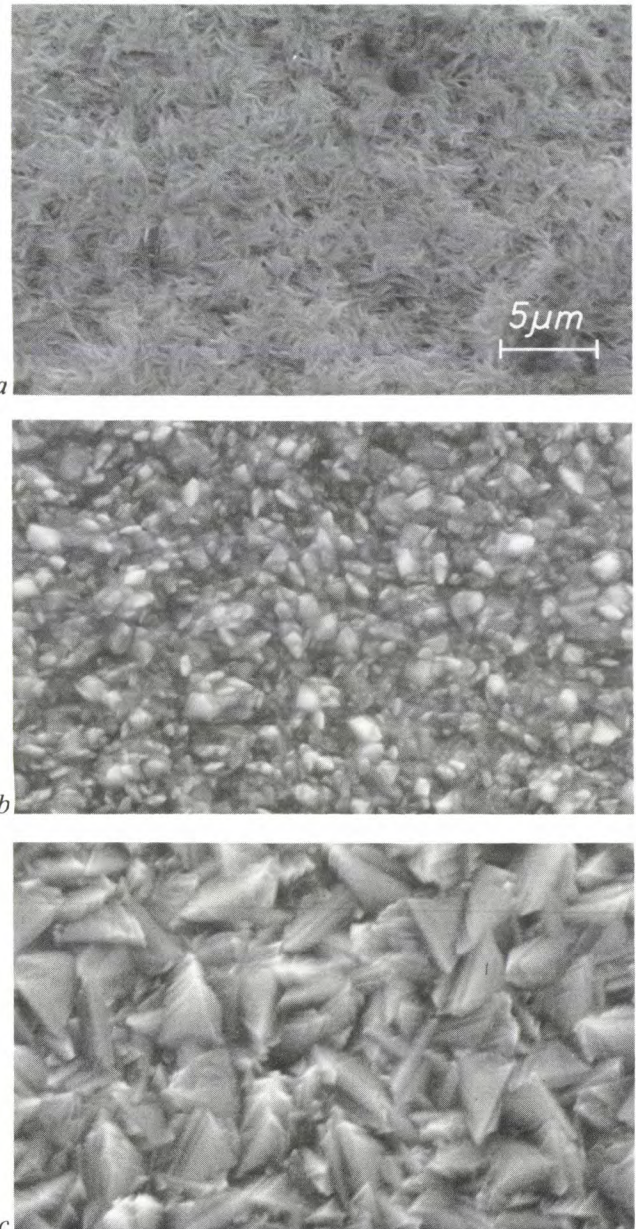


Fig. 3. Photographs of TiC coatings produced by CVD with different CH_4 contents on tool steel containing less than 0.4 per cent by weight of carbon. *a*) CH_4 content of standard process (2.5% by volume): relatively smooth surface. *b*) 5 times as much CH_4 : rough surface. *c*) 10 times as much CH_4 : even rougher surface.

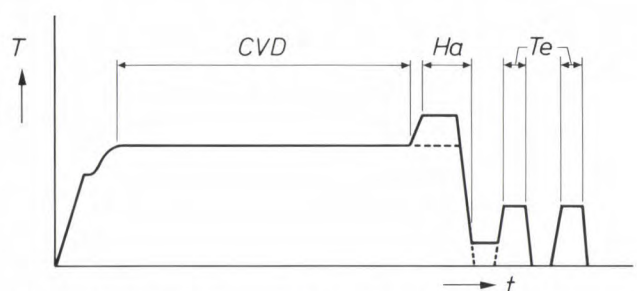


Fig. 4. Schematic diagram of temperature T against time t in coating tool steel whose hardening temperature is above the CVD process temperature. After heating up in steps, the CVD process is carried out, and is immediately followed by the hardening process H_a in the same reactor. Tempering T_e is carried out in two separate stages.

near their maximum hardness. If these types of steel are to be used, and if the substrate is required to have high mechanical strength, then the coated workpiece has to be heated in a separate furnace with all oxygen excluded, and subsequently quenched.

When 'air'-hardened types of steel are used, which do not have to be quenched in a liquid, the deposition of a wear-resistant coating and hardening of the substrate can successfully be combined in a single operation [6]. In this case the steel must have a large hardening depth and the hardening temperature must not differ too much from the CVD process temperature. In addition there must be little tendency to grain growth and little danger of austenite stabilization. Fig. 4 shows a temperature-time diagram for such a combination of CVD process and steel hardening. The CVD process temperature here is somewhat lower than the hardening temperature of the steel. After the deposition the workpiece is heated in the same reactor to the correct temperature for hardening. It is then left to cool, and during the cooling the austenite is converted into martensite. Tempering, which has the effect of converting the martensite into a less brittle structure, is carried out separately from the CVD process.

Certain types of air-hardened steel are found in practice to be most suitable for a CVD coating. When a TiC coating is to be applied, preference is given, for reasons mentioned earlier, to a tool steel containing at least 0.4 per cent of carbon. A CVD coating on air-hardened tool steel gives a very good combination of materials, with the workpiece providing optimum support for the coating. In addition, the dimensional changes due to hardening and tempering are small compared with the changes in types of steel hardened in water or oil. Fig. 5 shows a photograph of a combination of a suitable type of air-hardened steel with a thin TiC coating.

It is important that the workpiece to be coated should be clean and free from rust. It should also be ground or lapped to the exact dimensions required. Since the coating follows the profile of the substrate accurately, the roughness of the substrate surface affects the roughness of the final product.

Coating complex shapes

CVD can be used to deposit a uniform coating on workpieces of complex shape. This is very important if say, machine parts with narrow holes have to be coated with a thin TiC or TiN layer. The uniformity depends on the dimensions of the holes and on the conditions during the CVD process. To some extent these conditions can be varied to obtain the desired uniformity.

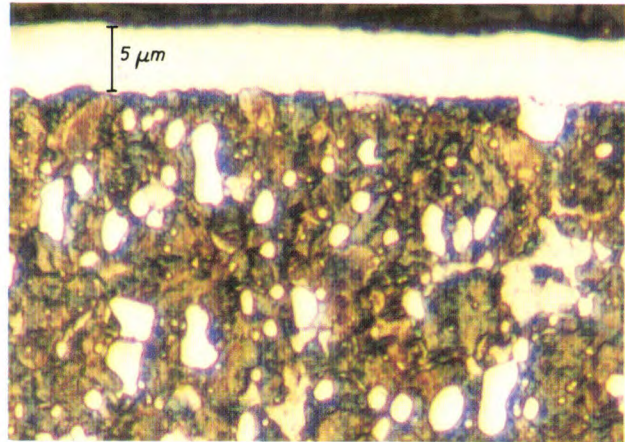


Fig. 5. Photograph of air-hardened tool steel coated with a uniform layer of TiC about 5 μm thick by chemical vapour deposition.

Mass transfer in a CVD process can be regarded as a two-stage process: diffusion of the reactive molecules in the gase phase to the substrate, and the reaction on the substrate surface. Either of these stages can determine the growth rate g . In general the curve of $\ln g$ against T^{-1} (the Arrhenius curve) therefore consists of two sections of differing slope. At low temperature the surface reaction determines g ; the activation energy is then relatively high. At high temperatures the diffusion determines g , and this is reflected in a weaker temperature dependence of the growth rate. CVD can be used to deposit a uniform coating on complex shapes if the surface process is the stage that determines the rate of growth. We shall now deal with the uniformity with the aid of some results of calculations and experiments carried out on CVD coatings in narrow tool-steel tubes [7].

In a narrow tube the CVD deposition rate is a monotonically decreasing function of the axial distance z , measured from the tube opening. This also applies to the thickness d of the deposited coating, which is of course proportional to the local growth rate. C. H. J. van den Brekel, of Philips Research Laboratories, has shown that the axial thickness variation of a CVD coating deposited in a tube open at both ends, of length l and radius r , is given approximately by

$$\frac{d(z/r)}{d(z=0)} = \frac{\cosh \{(z - l/2) \sqrt{2K}/r\}}{\cosh \{l \sqrt{2K}/2r\}}, \quad (1)$$

where K is the effective Sherwood number. This quan-

[4] E. Horvath and A. J. Perry, *Material und Technik* 7, 63, 1979.

[5] H. E. Hintermann and H. Boving, *Die Technik* 33, 387, 1978.

[6] L. B. J. Janssen, *Polytechn. T. Werktuigbouw* 34, 491, 1979 (in Dutch).

[7] C. H. J. van den Brekel, R. M. M. Fonville, P. J. M. van der Straten and G. Verspui, *Proc. 8th Int. Conf. on Chemical vapor deposition, Gouvieux 1981 (Electrochem. Soc. Proc. 81-7)*, p. 142.

tity K depends on r , the order n of the surface reaction, the concentration c_b in the gas, the equilibrium concentration c_e , the diffusion constant D of the rate-determining component in the gas phase, and the mass-transfer coefficient k :

$$K = r n (c_b - c_e)^{n-1} D^{-1} k. \quad (2)$$

The number K determines the uniformity of the coating: at a given l and r the uniformity increases as the value of K decreases. To illustrate this, *fig. 6* gives the calculated thickness variation in an infinitely long

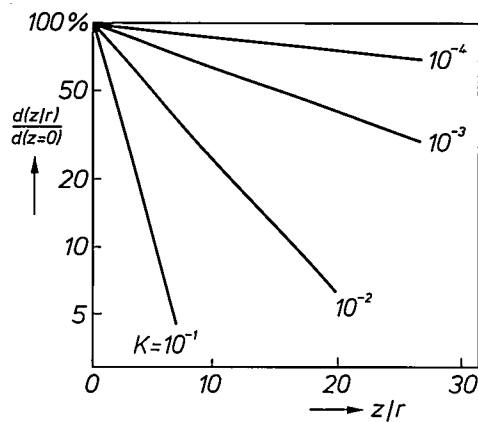


Fig. 6. Calculated axial thickness variation in an infinitely long tube for four values of K , the effective Sherwood number. The ratio $d(z/r)/d(z=0)$, calculated from eq. (3), is plotted as a function of z/r , where z is the axial distance from the tube opening and r is the radius of the tube. The uniformity increases as the value of K decreases.

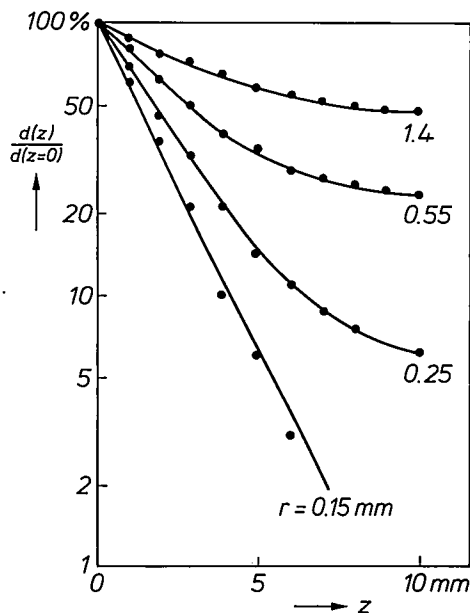


Fig. 7. Measured axial thickness variation for a CVD coating of TiN in four tubes of length 20 mm and different radius r . The deposition took place at a pressure of 900 mbar and a temperature of 850 °C. Provided the tube to be coated is not too narrow (> about 1 mm), a reasonably uniform coating is obtained.

tube for four values of K . In this case the variation is given by

$$\frac{d(z/r)}{d(z=0)} \stackrel{l \rightarrow \infty}{=} \exp\left(-\frac{z}{r} \sqrt{2K}\right). \quad (3)$$

Since the uniformity of TiC and TiN coatings deposited in narrow tubes is virtually independent of c_b , we may put n in eq. (2) equal to 1, so that

$$K = r D^{-1} k. \quad (4)$$

The diffusion constant D depends on the temperature T and the pressure P in the gas phase:

$$D \propto T^m/P, \quad (5)$$

where m is approximately equal to 2. The mass-transfer coefficient k depends on T and the activation energy E of the surface reaction:

$$k \propto \exp(-E/RT), \quad (6)$$

where R is the molar gas constant. Substituting equations (5) and (6) in eq. (4) gives

$$K \propto r P T^{-m} \exp(-E/RT). \quad (7)$$

In most cases the measured thickness variation agrees reasonably well with the calculated variation. As an example of the influence of the tube radius r , *fig. 7* gives some measured thickness variations for TiN coatings deposited in narrow tubes of the same length (20 mm) but with different radii. The coatings were deposited at a pressure of 900 mbar and at a temperature of 850 °C. As equations (1) and (2) indicate, the uniformity of the coating in a tube of radius 1.4 mm is much better than in a tube of radius only 0.15 mm.

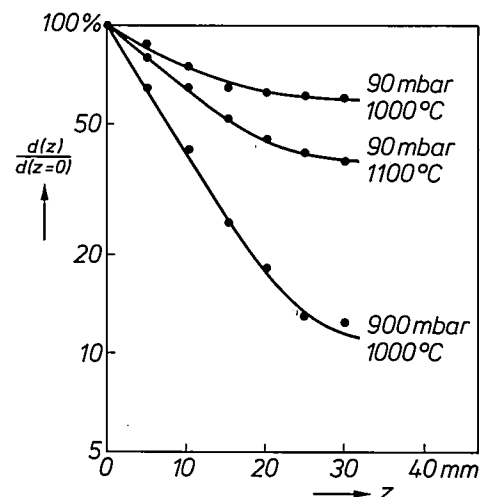


Fig. 8. Measured axial thickness variation for a TiC coating produced by CVD in a tube of length of 70 mm and a radius of 0.5 mm, for three combinations of pressure and temperature during the CVD process. Reducing the pressure or the temperature gives greater uniformity.

The measured pressure and temperature dependences are also in good agreement with theory. As equations (1) and (7) indicate, it was found that the uniformity of the deposited coating increases as the total pressure and temperature decrease during the CVD process. Fig. 8 gives the measured thickness variation of a TiC coating deposited in a narrow tube ($l = 70$ mm, $r = 0.5$ mm) for three P, T combinations. Owing to a pressure decrease from 900 to 90 mbar and a temperature decrease from 1100 to 1000 °C the coating is much more uniform. A further decrease of pressure and temperature may produce a further improvement in the uniformity of the coating, but has the disadvantage that the coating process can take too long because of the marked decrease of the growth rate.

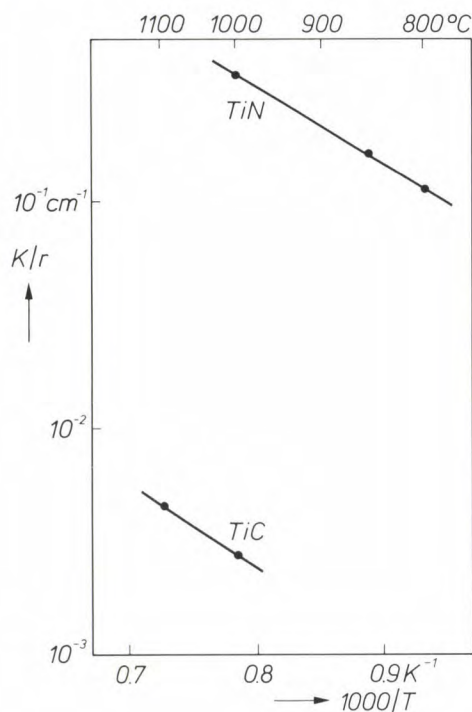


Fig. 9. Ratio K/r (K = the effective Sherwood number) plotted on a logarithmic scale against $1000/T$ (T is the absolute temperature) for TiC and TiN coatings produced by CVD in a tube of radius r . The values of K were determined from eq. (1) and the measured thickness variations. The pressure was 90 mbar during the deposition of TiC and 900 mbar during the deposition of TiN. The great uniformity (small K values) of TiC coatings is mainly due to the low pressure during the CVD process and to the high activation energy of the surface reaction (eq. 7).

In general the uniformity of a TiC coating is greater than that of a TiN coating on the same workpiece. This is illustrated in fig. 9, which gives a logarithmic plot of the experimental K/r values of TiC and TiN against $1000/T$. The K -values were determined from the measured thickness variations after coating with TiC at 90 mbar and with TiN at 900 mbar. For the

TiC coating the K/r -value at 1000 °C is about a factor of 100 smaller than for the TiN coating. If the pressure during the CVD process is the same, say 900 mbar, then the difference is still about a factor of 10. This difference is largely the result of the higher activation energy of the surface reaction for TiC.

The low pressure (90 mbar) during the CVD process and the high activation energy of the surface reaction permit a uniform TiC coating on all kinds of complicated workpieces, containing features such as narrow gaps or even blind holes. For instance, in a blind hole with a radius of 0.5 mm and a length of 40 mm the TiC coating thickness at the bottom of the hole is about 70% of that on the outside of the workpiece. The uniformity of the coating is a great advantage compared with other coating techniques, which cannot produce a uniform coating on complicated workpieces, or can only do so with great difficulty.

Applications

The deposition of a hard, wear-resistant TiC or TiN coating on a tool-steel workpiece by means of CVD can greatly improve the durability and reliability of tools [4][8]. The great resistance of these coatings to abrasion and adhesive wear, and their low coefficient of friction, greatly reduce the problems of 'galling', 'fretting', corrosion and erosion. In the past few years, with the assistance of the Tribology Project Group of CFT, the Metallurgical Laboratory of the Philips Plastics and Metalware Factories (PMF) and the Materials Laboratory of the Philips Engineering Works ('*Machinefabrieken*') we have investigated various potential applications, some of which we shall now mention.

In the processing of plastics with abrasive fillers such as crushed quartz and magnesium oxide there is much less wear of the extrusion dies when they are coated with TiC. An uncoated die of a particular tool steel was no longer usable after only 1000 extrusions, but the same die coated with 5 μm of TiC still gave excellent service after 8000 extrusions.



Fig. 10. Photograph ($1.5 \times$ full size) of an uncoated riveting tool (left) and a riveting tool with a coating of 5 μm of TiC formed by CVD (right).

[8] H. Benninghoff and H. Zickler, *Metalloberfläche* 32, 118, 1978. E. Horvath and H. Zimmermann, *Fertigung*, Nov. 1979.

An even greater extension of tool life is found with riveting tools used for fitting pins in fluorescent-lamp bases. *Fig. 10* shows two such tools, one coated with $5\ \mu\text{m}$ of TiC, the other uncoated. An uncoated tool

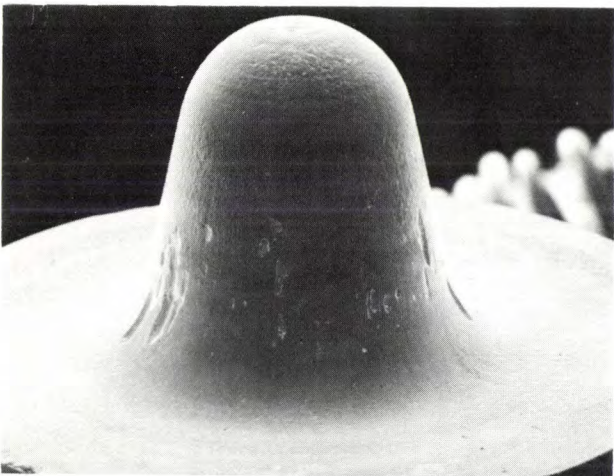
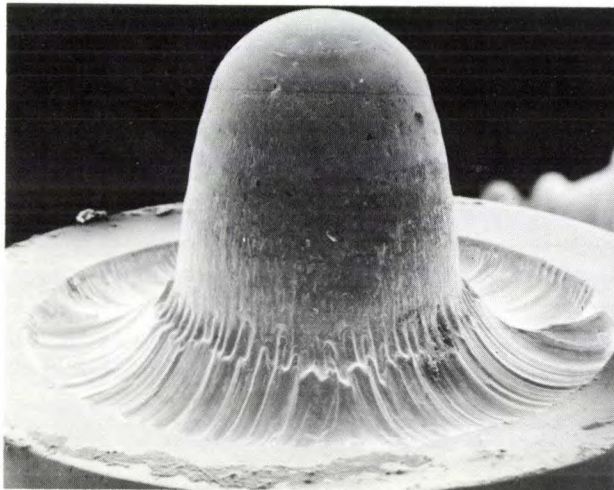


Fig. 11. Photographs of two used riveting tools. *Above:* An uncoated tool after 200 hours use in production. The wear has made the tool unserviceable. *Below:* A tool coated with $5\ \mu\text{m}$ of TiC, formed by CVD, after 3900 hours. There is hardly any wear even though the tool has been used for much longer.

has to be taken out of service after 200 hours use in production (6000 operations per hour) because of the wear shown in *fig. 11*. A coated tool, however, has hardly any wear after 3900 hours.

Sheet-steel scrapers, used for scraping away paint from a wheel, last four times as long when they are coated with $4\ \mu\text{m}$ of TiC.

In addition to these three examples there are many other promising applications for hard, wear-resistant TiC and TiN coatings on tools such as cutting and shearing tools, guide pins, deep-drawing tools, dies and guillotines. We are planning to extend the technology to coating materials that have other interesting properties in addition to those of TiC and TiN^{[9][10]}. These materials include aluminium oxide (Al_2O_3), which unlike TiC and TiN is an electrical insulator and is very hard at high temperatures, and tungsten carbide (W_2C), which can be deposited at relatively low temperatures ($500\text{--}800\ ^\circ\text{C}$). Other materials, such as SiO_2 , Si_3N_4 , SiC, boron and tungsten are also being studied.

- ^[9] R. V. Leverenz, *Manufacturing Engng* **79**, 38, 1977.
R. Funk, H. Schachner, C. Triquet, M. Kornmann and B. Lux, *J. Electrochem. Soc.* **123**, 285, 1976.
H. O. Pierson, E. Randich and D. M. Mattox, *J. less-common Met.* **67**, 381, 1979.
- ^[10] W. Schintlmeister, O. Pacher, W. Wallgram and J. Kanz, *Proc. 3rd Eur. Conf. on Chemical vapour deposition, Neuchâtel 1980*, p. 181, and *Metall* **34**, 905, 1980.

Summary. The life and reliability of a workpiece made of tool steel can be substantially improved by coating it with a thin ($\leq 10\ \mu\text{m}$) wear-resistant layer of TiC or TiN by means of chemical vapour deposition (CVD). The improvement depends on the CVD process conditions, the kind of tool steel to be coated and the shape of the workpiece. Good results have been obtained for TiC and TiN coatings on several types of air-hardened steel. A uniform coating of TiC can be formed on workpieces of complex shape. The combination of the desirable properties of a TiC or TiN coating with those of the tool steel provides great scope for many successful applications.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- B. Andlauer** (Fraunhofer-Institut für Angewandte Festkörperphysik, Freiburg) & **W. Tolksdorf**: Associated incorporation of boron, lead, and oxygen vacancies in garnets. *J. appl. Phys.* **50**, 7986-7995, 1979 (No. 12). *H*
- D. E. Aspnes** (Bell Laboratories, Murray Hill, N.J.) & **J. B. Theeten**: Dielectric function of Si-SiO₂ and Si-Si₃N₄ mixtures. *J. appl. Phys.* **50**, 4928-4935, 1979 (No. 7). *L*
- M. Auphan & J. Matthys**: Reflection of a plane impulsive acoustic pressure wave by a rigid sphere. *J. Sound and Vibr.* **66**, 227-237, 1979 (No. 2). *L, B*
- P. Baudet**: Comportement en bruit et réalisation pratique des transistors à effet de champ petits signaux en arséniure de gallium pour applications hyperfréquences. *Acta Electronica* **23**, 111-118, 1980 (No. 2). *L*
- P. Baudet**: Les transistors à effet de champ de puissance en arséniure de gallium: conception et technologie. *Acta Electronica* **23**, 119-125, 1980 (No. 2). *L*
- J. R. van Beek & P. J. Rommers**: Behaviour of the secondary lithium electrode on alloying substrates in propylene carbonate based electrolytes. *Power Sources* **7**, Proc. 11th Int. Symp., Brighton 1978, pp. 595-622; 1979. *E*
- C. Belouet, E. Fabre, S. Makram-Ebeid, Ngo-Tich Phuoc & C. Texier**: Early assessment of the photovoltaic potentialities of RAD polysilicon sheets. Photovoltaic Solar Energy Conf., Berlin 1979, pp. 114-122. *L*
- F. Berz**: Step recovery of *p-i-n* diodes. *Solid-State Electronics* **22**, 927-932, 1979 (No. 11). *R*
- M. Binet & P. Baudet**: Caractérisation hyperfréquence des transistors à effet de champ: mesures du facteur de bruit, du gain, et de la puissance sur banc hyperfréquence. *Acta Electronica* **23**, 127-136, 1980 (No. 2). *L*
- L. A. Boatner** (Oak Ridge National Laboratory), **E. Krätzig & R. Orłowski**: KTN as a holographic storage material. *Ferroelectrics* **27**, 247-250, 1980 (No. 1-4). *H*
- D. Boccon-Gibod**: Modèle analytique et schéma équivalent du transistor à effet de champ en arséniure de gallium. *Acta Electronica* **23**, 99-109, 1980 (No. 2). *L*
- O. Boser**: The behavior of Inconel 625 in a silver environment. *Mat. Sci. Engng.* **41**, 59-64, 1979 (No. 1). *N*
- R. Brehm, K. van Dun, J. C. G. Teunissen & J. Haisma**: Transparent single-point turning of optical glass: A phenomenological presentation. *Precision Engng.* **1**, 207-213, 1979 (No. 4). *E*
- A. Broese van Groenou & P. E. C. Franken**: Second phase in electroceramics: its detection and usefulness. *Proc. Brit. Ceramic Soc.* **28**, 243-266, 1979. *E*
- E. Bruninx**: X-ray fluorescence analysis with a Seemann spectrometer, improved pulse analysis and crystal dispersion. *Anal. chim. Acta* **113**, 97-106, 1980 (No. 1). *E*
- K. Bulthuis, M. G. Carasso, J. P. J. Heemskerk, P. J. Kivits, W. J. Kleuters & P. Zalm**: Ten billion bits on a disk. *IEEE Spectrum* **16**, No. 8, 26-33, Aug. 1979. *E*
- K. H. J. Buschow**: H₂-absorption in intermetallische verbindungen. *Chem. Weekbl. Mag.* 1980, p. m 39 (Jan.). *E*
- K. H. J. Buschow & N. M. Beekmans**: Formation, decomposition, and electrical transport properties of amorphous Hf-Ni and Hf-Co alloys. *J. appl. Phys.* **50**, 6348-6352, 1979 (No. 10). *E*
- K. H. J. Buschow & A. M. van Diepen**: Moment formation upon hydrogen absorption in Hf₂Fe. *Solid State Comm.* **31**, 469-471, 1979 (No. 7). *E*

- K. H. J. Buschow & W. W. van den Hoogenhof:** Note on the formation and the magnetic properties of amorphous Eu-Ag alloys.
J. Magn. magn. Mat. **12**, 123-126, 1979 (No. 2). *E*
- K. H. J. Buschow, P. R. Locher & M. Leger** (CNRS, Meudon): Nuclear magnetic resonance and pressure effects in $\text{UNi}_{5-5x}\text{Cu}_{5x}$.
J. Physics F **9**, 2483-2490, 1979 (No. 12). *E*
- P. J. Cameron** (Merton College, Oxford), **J.-M. Goethals & J. J. Seidel** (Eindhoven University of Technology): Strongly regular graphs having strongly regular subconstituents.
J. Algebra **55**, 257-280, 1978 (No. 2). *B*
- B. L. Cardozo & K. G. van der Veen** (Institute for Perception Research, Eindhoven): Estimation of annoyance due to low level sound.
Appl. Acoustics **12**, 389-396, 1979 (No. 5).
- A. I. Carlson & B. M. Singer:** Variable gain X-ray image intensifier tube.
IEEE Trans. ED-26, 1711-1717, 1979 (No. 11). *N*
- M. Cathelin:** Technologies de circuits intégrés sur arsénure de gallium.
Acta Electronica **23**, 193-204, 1980 (No. 3). *L*
- V. Chalmeton & E. Rammos:** Study of X-ray resist sensitivity with synchrotron radiation.
Proc. Microcircuit Engineering '79 — Microstructure Fabrication, Aachen 1979, pp. 98-107. *L*
- T. A. C. M. Claasen & W. F. G. Mecklenbräuker:** Extension of a method for sharpening the response of digital filters based on the amplitude change function.
IEEE Trans. ASSP-27, 559-560, 1979 (No. 5). *E*
- T. A. C. M. Claasen & W. F. G. Mecklenbräuker:** The Wigner distribution — a tool for time-frequency signal analysis, Part III: Relations with other time-frequency signal transformations.
Philips J. Res. **35**, 372-389, 1980 (No. 6). *E*
- W. A. P. Claassen & J. Bloem:** Kiemvorming van silicium op SiO_2 en Si_3N_4 .
Ned. T. Natuurk. A **45**, 119-122, 1979 (No. 3). *E*
- W. J. Dallas:** Pre-folded Fourier-filter reconstruction of coded-aperture images.
Optica Acta **26**, 1359-1365, 1979 (No. 11). *H*
- J. P. M. Damen, J. M. Robertson & M. A. H. Huyberts:** Thin film growth of $(\text{Mn}_x\text{Zn}_{1-x})\text{Fe}_2\text{O}_4$ by liquid phase epitaxy on Bridgman grown Zn_2TiO_4 substrates.
J. Crystal Growth **47**, 486-492, 1979 (No. 4). *E*
- T. T. Dao** (Signetics Corp., Sunnyvale, Calif.), **M. Davio & C. Gossart:** Complex number arithmetic with balanced ternary logic.
Proc. 9th Int. Symp. on Multiple-valued logic, Bath 1979, pp. 281-289. *B*
- M. Delfino, G. M. Loiacono, W. A. Smith, A. Shaulov*, Y. H. Tsuo* & M. I. Bell*** (* Yeshiva University, New York): Thermal and dielectric properties of LiKSO_4 and LiCsSO_4 .
J. solid State Chem. **31**, 131-134, 1980 (No. 1). *N*
- M. Delfino, J. A. Nicolosi & J. Ladell:** Crystal data for α -iminodiacetic acid.
J. appl. Cryst. **12**, 621-622, 1979 (No. 6). *N*
- P. Delsarte, Y. Genin & Y. Kamp:** Schur parametrization of positive definite block-Toeplitz systems.
SIAM J. appl. Math. **36**, 34-46, 1979 (No. 1). *B*
- P. Delsarte, Y. Genin & Y. Kamp:** The Nevanlinna-Pick problem for matrix-valued functions.
SIAM J. appl. Math. **36**, 47-61, 1979 (No. 1). *B*
- P. Delsarte, Y. Genin & Y. Kamp:** Characterization theorems for bounded and positive two-variable functions.
Int. J. Circuit Theory and Appl. **7**, 257-266, 1979 (No. 2). *B*
- P. Delsarte, Y. Genin & Y. Kamp:** A survey of two-dimensional Toeplitz systems.
Int. Symp. on Mathematical theory of networks and systems, Delft 1979, pp. 17-30. *B*
- P. Delsarte, Y. Genin & Y. Kamp:** Generalized Schur parametrization.
Int. Symp. on Mathematical theory of networks and systems, Delft 1979, pp. 290-296. *B*
- P. A. Devijver:** Statistical methods of pattern recognition.
Actes Journ. Int. Reconnaissance des Formes, GAST, Lyon 1979, pp. 10.01-10.29. *B*
- J. G. Dil & C. A. Wesdorp** (Philips Electro-acoustics Division, Eindhoven): Control of pit geometry on video disks.
Appl. Optics **18**, 3198-3202, 1979 (No. 18). *E*
- C. Z. van Doorn:** De optische eigenschappen van een gedraaid nematisch display.
Ned. T. Natuurk. A **45**, 151-153, 1979 (No. 4). *E*
- J. P. Dougherty & R. J. Seymour:** Automated, simultaneous measurement of electrical properties of pyroelectric materials.
Rev. sci. Instr. **51**, 229-233, 1980 (No. 2). *N*
- H. Duifhuis, L. F. Willems** (both with Institute for Perception Research, Eindhoven) & **R. J. Sluyter:** An outline of pitch analysis in speech: a hearing theory approach.
Hearing mechanisms and speech, ed. O. Creutzfeldt, H. Scheich & Chr. Schreiner (Exper. Brain Res. Suppl. II), pp. 254-259; Springer, Berlin 1979. *E*
- G. E. J. Eggermont, Y. Tamminga & W. K. Hofker** (Philips Res. Labs, Amsterdam Dept): Elemental and dose dependent threshold for Nd-YAG laser induced recrystallization of silicon.
AIP Conf. Proc. **50**, 321-324, 1979.
- J. van Esdonk:** Bonding techniques for ultra-high-vacuum systems.
Preprints 3rd Int. Brazing and Soldering Conf., London 1979, paper 30, 11 pp. *E*
- E. Fischer:** Combustion of zirconium in oxygen at high pressures.
Philips J. Res. **35**, 390-440, 1980 (No. 6). *A*

- W. E. Fischer:** The interface between CAD/CAM-software and CODASYL-data base management systems. Eurographics 79, Proc. Int. Conf. and Exhib., Bologna 1979, pp. 178-189. *H*
- E. Fischmann:** Retracing Röntgen's discovery. Diagnostic Imaging **48**, 294-303, 1979 (No. 5). *E*
- B. Fitzhenry:** Optical effects of adsorption of dyes on pigment used in electrophoretic image displays. Appl. Optics **18**, 3332-3337, 1979 (No. 19). *N*
- P. E. C. Franken & W. T. Stacy:** Microstructuur van gesinterd MnZn-ferroferriet. Ned. T. Natuurk. A **45**, 111-114, 1979 (No. 3). *E*
- R. C. French:** The effect of fading and shadowing on channel reuse in mobile radio. IEEE Trans. VT-**28**, 171-181, 1979 (No. 3). *R*
- P. Gerthsen*, K. H. Härdtl & N. A. Schmidt* (* Universität Karlsruhe):** Correlation of mechanical and electrical losses in ferroelectric ceramics. J. appl. Phys. **51**, 1131-1134, 1980 (No. 2). *A*
- J.-M. Goethals:** Schémas d'association et designs sphériques. Cahiers CERO **20**, 287-292, 1978 (No. 3/4). *B*
- R. G. Gossink, H. A. M. de Grefte & H. W. Werner:** Analysis of sodium depth profiles in glasses using secondary ion mass spectrometry (SIMS). Silicates ind. **44**, 35-41, 1979 (No. 2). *E*
- S. Gourrier, A. Mircea, J. B. Theeten & M. Bacal (Ecole Polytechnique, Palaiseau):** Use of a multipole plasma for the oxidation of GaAs. 4th Int. Symp. on Plasma chemistry, Zürich 1979, pp. 181-186. *L*
- H. de Graaf, W. J. Huiskamp, R. C. Thiel, H. Th. LeFever (all with Rijksuniversiteit Leiden) & K. H. J. Buschow:** Mössbauer effect and magnetic properties in EuMg_5 and $\text{Eu}_2\text{Mg}_{17}$. Physica **98B**, 60-64, 1979 (No. 1/2). *E*
- P. Günter (ETH, Zürich), P. M. Asbeck & S. K. Kurtz:** Second-harmonic generation with $\text{Ga}_{1-x}\text{Al}_x\text{As}$ lasers and KNbO_3 crystals. Appl. Phys. Letters **35**, 461-463, 1979 (No. 6). *N*
- G. J. van Gorp, G. E. J. Eggermont*, Y. Tamminga*, W. T. Stacy & J. R. M. Gijsbers (* Philips Res. Labs, Amsterdam Dept):** Cellular structure and silicide formation in laser-irradiated metal-silicon systems. Appl. Phys. Letters **35**, 273-275, 1979 (No. 3). *E*
- G. J. van Gorp, J. L. C. Daams, A. van Oostrom, L. J. M. Augustus & Y. Tamminga:** Aluminum-silicide reactions. I. Diffusion, compound formation, and microstructure. J. appl. Phys. **50**, 6915-6922, 1979 (No. 11, Part I). *E*
- G. J. van Gorp & W. M. Reukers:** Aluminum-silicide reactions. II. Schottky-barrier height. J. appl. Phys. **50**, 6923-6926, 1979 (No. 11, Part I). *E*
- P. Hansen & H. Heitmann:** Influence of nuclear tracks on the magnetic properties of a $(\text{Gd,Bi})_3(\text{Fe,Ga})_5\text{O}_{12}$ garnet film. Phys. Rev. Letters **43**, 1444-1447, 1979 (No. 19). *H*
- P. Hansen & M. Urner-Wille:** Magnetic and magneto-optic properties of amorphous GdFeBi-films. J. appl. Phys. **50**, 7471-7476, 1979 (No. 11, Part II). *H*
- J. Hasker:** Recombination structures and their effects in depleted low-pressure gas discharges. J. appl. Phys. **50**, 5007-5011, 1979 (No. 7). *E*
- E. E. Havinga & P. van Pelt:** Electrochromism of substituted polyalkenes in polymer matrices; influence of chain length on charge transfer. Ber. Bunsen-Ges. Phys. Chemie **83**, 816-821, 1979 (No. 8). *E*
- W. K. Hofker, G. E. J. Eggermont, Y. Tamminga & D. P. Oosthoek (Philips Res. Labs, Amsterdam Dept):** The removal of implantation damage in silicon by laser irradiation and its thermal stability. AIP Conf. Proc. **50**, 425-428, 1979.
- K. Holford:** Vehicle control by portable traffic lights. 2nd Int. Conf. on Automotive electronics, London 1979 (IEE Conf. Publ. No. 181), pp. 180-184. *R*
- L. Honds & H. Meyer:** Unipolar-Linearmotor mit offenem Statorkreis. Elektrotechn. Z. **101**, 285-289, 1980 (No. 5). *A*
- L. Honds & H. Meyer:** Nichtlinearität der Kraft-Weg-Kurven von Tauchspullinearmotoren. Feinwerktechnik & Messtechnik **88**, 162-166, 1980 (No. 4). *A*
- A. J. Huart (Philips Elcoma Division, Eindhoven):** De Europese IC-industrie. Micro-elektronica, 'uitdaging voor allen', lectures KIVI-Symp., Delft 1979, pp. 41-51; KIVI/Misset, Doetinchem 1979.
- J. B. Hughes, E. S. Eilley & G. J. Glynn:** Modulo- N counter technique for the u.h.f. band. Electronics Letters **15**, 640-641, 1979 (No. 20). *R*
- R. H. Johnston:** Speed measurement by radar for anti-lock braking. 2nd Int. Conf. on Automotive electronics, London 1979 (IEE Conf. Publ. No. 181), pp. 185-188. *R*
- H. D. Jonker & W. van Erk:** Segregation of Ca and Ge in LPE growth of magnetic YSmCaFeGe garnet films. J. Crystal Growth **48**, 131-140, 1980 (No. 1). *E*
- J. J. Kelly:** The influence of fluoride ions on the passive dissolution of titanium. Electrochim. Acta **24**, 1273-1282, 1979 (No. 12). *E*
- J. Kittler (Oxford University) & P. A. Devijver:** Statistical properties of error estimators in performance assessment of recognition systems. Proc. 1979 IEEE Computer Soc. Conf. on Pattern recognition and image processing, Chicago, pp. 44-51. *B*

- A. G. Knapp:** The effect of electron bombardment on the secondary electron emission from Na_3AlF_6 .
J. appl. Phys. **50**, 5961-5965, 1979 (No. 9). *R*
- H. Köstlin & H. Schaper:** Electrochemiluminescence by dc of rubrene displaying highly organized electrohydrodynamic convection.
Physics Letters **76A**, 455-458, 1980 (No. 5,6). *A*
- G. Kowalski:** Future development of computerized tomography.
Biomed. Technik **24**, Ergänzungsband, 223-224, 1979. *H*
- E. Krätzig & R. Orlowski:** Light induced charge transport in doped LiNbO_3 and LiTaO_3 .
Ferroelectrics **27**, 241-244, 1980 (No. 1-4). *H*
- H. Kropp** (T.H. Darmstadt), **E. Dormann** (Universität Bayreuth) & **K. H. J. Buschow:** Electric field gradient in cubic intermetallic europium compounds with unstable europium valence.
Solid State Comm. **32**, 507-510, 1979 (No. 7). *E*
- H. Kropp, W. Zipf** (both with T.H. Darmstadt), **E. Dormann** (Universität Bayreuth) & **K. H. J. Buschow:** Indirect exchange in intermetallic europium compounds.
J. Magn. magn. Mat. **13**, 224-230, 1979 (No. 1/2). *E*
- M. H. Kuhn:** Access control by means of automatic speaker verification.
J. Physics E **13**, 85-86, 1980 (No. 1). *H*
- M. H. Kuhn:** Speaker recognition accounting for different voice conditions by unsupervised classification (cluster analysis).
IEEE Trans. **SMC-10**, 54-57, 1980 (No. 1). *H*
- L. J. M. Kuijpers, G. A. A. Asselman & G. v.d. Berk-mortel:** Reed valve simulation: a comparison of numerical and experimental results.
Proc. XV Int. Congress of Refrigeration, Venezia 1979, paper B2-14, 6 pp. *E*
- D. Küppers:** Recent developments in plasma activated chemical vapor deposition.
Proc. 7th Int. Conf. on Chemical vapor deposition, Los Angeles 1979 (Electrochem. Soc. Proc. **79-3**), pp. 159-175. *A*
- D. Küppers & H. Lydtin:** Preparation of optical waveguides with the aid of plasma-activated chemical vapour deposition at low pressures.
Topics in Current Chemistry **89**, 107-131, 1980. *A*
- A. van Lamsweerde & M. Sintzoff:** Formal derivation of strongly correct concurrent programs.
Acta Informatica **12**, 1-31, 1979 (No. 1). *B*
- P. K. Larsen, J. H. Neave & B. A. Joyce:** Angular resolved photoemission from surface states on reconstructed {100} GaAs surfaces.
J. Physics C **12**, L 869-874, 1979 (No. 22). *E, R*
- C. Le Can:** Perspective of gallium arsenide integrated circuits.
Acta Electronica **23**, 191-192, 1980 (No. 3). *L*
- M. Lemke, W. Hoppe, W. Tolksdorf & F. Welz:** Magnetically tunable millimetre-wave filter with single-crystal barium ferrite.
IEE J. Microw. Opt. Acoust. **3**, 253-254, 1979 (No. 6). *H*
- J. Lohstroh & J. D. P. van den Crommenacker:** First-order modeling and temperature behaviour of standard ISL-gates.
5th Eur. Solid State Circuits Conf. — ESSCIRC 79, Southampton 1979 (IEE Conf. Publ. No. 178), pp. 91-93. *E*
- G. M. Loiacono, M. Delfino, W. A. Smith, M. I. Bell*, A. Shaulov* & Y. H. Tsuo*** (* Yeshiva University, New York): Dielectric, pyroelectric, and thermal properties of LiNH_4SO_4 and LiND_4SO_4 .
Ferroelectrics **23**, 89-95, 1980 (No. 1/2). *N*
- S. Makram-Ebeid:** A computer model for polycrystalline Si n^+/p solar cells.
Photovoltaic Solar Energy Conf., Berlin 1979, pp. 792-799. *L*
- D. Meignant:** Etude de la fiabilité des transistors à effet de champ à barrière Schottky en GaAs pour hyperfréquences.
Acta Electronica **23**, 151-164, 1980 (No. 2). *L*
- R. Memming:** Solar energy conversion by photoelectrochemical processes.
Electrochim. Acta **25**, 77-88, 1980 (No. 1). *H*
- J. Michel & B. G. Martin:** A new diffusion process for silicon solar cells.
Photovoltaic Solar Energy Conf., Berlin 1979, pp. 181-188. *L*
- A. R. Miedema:** The formation enthalpy of monovacancies in metals and intermetallic compounds.
Z. Metallk. **70**, 345-353, 1979 (No. 6). *E*
- B. J. Mulder:** Protective glassy layers passivating copper at 500 °C, Part II.
J. Electrochem. Soc. **126**, 1425-1426, 1979 (No. 8). *E*
- B. J. Mulder:** A simple system for inserting samples into an ultrahigh-vacuum system through a vacuum lock.
J. Physics E **12**, 908, 1979 (No. 10). *E*
- B. J. Mulder:** Unbacked aluminium windows for helium discharge lamps directly mountable as part of the copper seal between UHV flanges.
J. Physics E **12**, 1036-1039, 1979 (No. 11). *E*
- P. H. Oosting:** Signal transmission in the nervous system.
Rep. Prog. Phys. **42**, 1479-1532, 1979 (No. 9). *E*
- A. van Oostrom:** Some aspects of Auger microanalysis.
Surface Sci. **89**, 615-634, 1979 (No. 1-3). *E*
- C. van Opdorp & H. Veenlyet:** On the relation between threshold current and interface recombination velocities in double-heterojunction lasers.
IEEE J. **QE-15**, 817-821, 1979 (No. 8). *E*

- R. Orlowski & E. Krätzig:** Holographic investigation of charge transport in electro-optic crystals. *Ferroelectrics* **26**, 831-834, 1980 (No. 1-4). *H*
- J. A. Pals & J. Dobben:** Measurements of microwave-enhanced superconductivity in aluminum strips. *Phys. Rev. B* **20**, 935-944, 1979 (No. 3). *E*
- M. Parisot, M. Binet & A. Rabier:** Caractérisation automatique en hyperfréquences du transistor à effet de champ. *Acta Electronica* **23**, 137-149, 1980 (No. 2). *L*
- P. Piret:** Graphe d'automorphisme des codes convolutionnels. *Cahiers CERO* **20**, 419-426, 1978 (No. 3/4). *B*
- P. Piret:** Addendum to 'Generalized permutations in convolutional codes'. *Information and Control* **40**, 332-334, 1979 (No. 3). *B*
- J. Polman, A. K. de Jonge & A. Castelijns:** A capacity-controlled free piston electrodynamic compressor. *Proc. XV Int. Congress of Refrigeration, Venezia 1979*, paper B2-12, 6 pp. *E*
- A. Rabier & M. Parisot:** C.A.D. of an octave band power FET amplifier. *SPACECAD '79, Computer-aided design of electronics for space applications, Proc. Int. Symp. Bologna 1979*, pp. 405-406. *L*
- J. M. Robertson, M. W. van Tol, J. P. H. Heynen, W. H. Smits & T. de Boer:** Thin single crystalline phosphor layers grown by liquid phase epitaxy. *Philips J. Res.* **35**, 354-371, 1980 (No. 6). *E*
- M. Rocchi:** Outil CAO pour circuits intégrés numériques sur arséniure de gallium. *Acta Electronica* **23**, 223-242, 1980 (No. 3). *L*
- M. Rocchi & M. Gavant:** Circuits intégrés numériques sur arséniure de gallium pour applications au-delà du gigahertz. *Acta Electronica* **23**, 243-267, 1980 (No. 3). *L*
- E. Roza, J. G. M. van Thuyt & W. van Doorn** (Philips Telecommunication Industries, Huizen): A novel concept for a hybrid 140 Mbit/s system. *IEEE Trans. COM-27*, 1584-1593, 1979 (No. 10). *E*
- A. M. J. G. van Run:** Computation of striated impurity distributions in melt-grown crystals, taking account of periodic remelt. *J. Crystal Growth* **47**, 680-692, 1979 (No. 5/6). *E*
- F. L. J. Sangster:** Master-slave charge-transfer device. *IEEE J. SC-14*, 624-626, 1979 (No. 3). *E*
- P. Schagen:** X-ray image intensifiers: design and future possibilities. *Phil. Trans. Roy. Soc. London A* **292**, 265-272, 1979 (No. 1390). *R*
- H. Schauer & E. Arnold:** Simple technique for charge centroid measurement in MNOS capacitors. *J. appl. Phys.* **50**, 6956-6961, 1979 (No. 11, Part I). *N*
- A. Sereny & V. Castelli** (Xerox Corp., Scarsdale, N.Y.): Experimental investigation of slider gas bearings with ultra-thin films. *Trans. ASME, J. Lubr. Technol.* **101**, 510-515, 1979 (No. 4). *N*
- H. R. Sethi & D. H. Paul:** MPlot3 — a system-independent plotting package. *Software — Pract. and Exper.* **9**, 891-905, 1979 (No. 11). *R*
- J. G. Siekman:** Analysis of laser drilling and cutting results in Al₂O₃ and ferrites. *AIP Conf. Proc.* **50**, 225-231, 1979. *E*
- M. Sintzoff:** Principles for distributing programs. *Lecture Notes in Computer Science* **70**, 337-347, 1979. *B*
- J. W. Slotboom & A. H. M. Goorman:** An efficient quasi three-dimensional bipolar transistor analysis program. *Numerical analysis of semiconductor devices, Proc. NASECODE I Conf., Dublin 1979*, pp. 280-289. *E*
- G. A. C. M. Spierings:** Optical absorption of transition metals in alkali lime germanosilicate glasses. *J. Mat. Sci.* **14**, 2519-2521, 1979 (No. 10). *E*
- G. A. C. M. Spierings:** Extraction of water from alkali germanosilicate glasses for optical fibres. *J. Mat. Sci.* **14**, 2919-2923, 1979 (No. 12). *E*
- H.-P. Stormberg:** Line broadening and radiative transport in high-pressure mercury discharges with NaI and TII as additives. *J. appl. Phys.* **51**, 1963-1969, 1980 (No. 4). *A*
- J. L. Teszner:** Introduction (*to issue on Gallium arsenide field effect transistor*). *Acta Electronica* **23**, 97-98, 1980 (No. 2). (*In English and in French.*) *L*
- A. Thayse:** Integer expansions of discrete functions and their use in optimization problems. *Proc. 9th Int. Symp. on Multiple-valued logic, Bath 1979*, pp. 82-87. *B*
- A. Thayse:** Encoding of parallel program schemata by vector addition systems. *Int. J. Computer and Inform. Sci.* **8**, 209-218, 1979 (No. 3). *B*
- J. B. Theeten, S. Gourrier & M. Steers:** Caractérisation optique des couches minces et des interfaces. *3ème Coll. Int. sur la Pulvérisation cathodique et ses applications, Nice 1979 (Suppl. Le Vide No. 196)*, pp. 517-528. *L*
- G. E. Thomas:** Bombardment-induced light emission. *Surface Sci.* **90**, 381-416, 1979 (No. 2). *E*
- N. C. de Troye:** Wat is micro-élektronica? *Micro-elektronica, 'uitdaging voor allen', lectures KIVI-Symp., Delft 1979*, pp. 9-20; *KIVI/Misset, Doetinchem 1979.* *E*

- H. J. M. Veendrick:** Design aspects and reliability of a synchronizer made in MOS technology. 5th Eur. Solid State Circuits Conf. — ESSCIRC 79, Southampton 1979 (IEE Conf. Publ. No. 178), pp. 8-10. *E*
- H. Veenliet, C. van Opdorp, R. P. Tijburg & J.-P. André:** Growth and characterization of MO/VPE double-heterojunction lasers. IEEE J. QE-15, 762-766, 1979 (No. 8). *E, L*
- J. D. B. Veldkamp:** Effective properties of polyphase materials. J. Physics D 12, 1375-1384, 1979 (No. 8). *E*
- C. H. F. Velzel & R. P. Brouwer:** Output power and coherence length of stripe-geometry double-heterostructure semiconductor lasers in incoherent feedback. IEEE J. QE-15, 782-786, 1979 (No. 8). *E*
- H. Verweij:** Raman study on glasses and crystalline compounds in the system K_3AsO_4 - $KAsO_3$. Appl. Spectrosc. 33, 509-515, 1979 (No. 5). *E*
- H. Verweij, H. van den Boom & R. E. Breemer:** Raman study of the reactions in a $30K_2CO_3 \cdot 70SiO_2 \cdot 1As_2O_3$ glass-forming batch. Silicates ind. 44, 71-77, 1979 (No. 3). *E*
- J. C. van Vessem** (Philips Elcoma Division, Eindhoven): Microelectronics, the revolution in consumer equipment. 5th Eur. Solid State Circuits Conf. — ESSCIRC 79, Southampton 1979 (IEE Conf. Publ. No. 178), pp. 20-24.
- L. Vriens:** Processen in en modellen van lagedrukmetaaldamp-edelgasontladingen. Ned. T. Natuurk. A 45, 164-168, 1979 (No. 4). *E*
- C. Werkhoven, J. H. T. Hengst & C. van Opdorp:** Annihilation of frozen-in point defects in GaP by thermal and recombination-induced processes. Appl. Phys. Letters 35, 136-138, 1979 (No. 2). *E*
- H. W. Werner:** Microanalysis of materials used in the electronics industry. Mikrochim. Acta, Suppl. 8, 25-50, 1979. *E*
- G. A. Wesselink:** The fluorescent light-window lamp. Philips J. Res. 35, 337-353, 1980 (No. 6). *E*
- G. F. Weston:** Measurement of ultra-high vacuum: Part 1. Total pressure measurements, Part 2. Partial pressure measurements. Vacuum 29, 277-291, 1979 (No. 8/9), & 30, 49-67, 1980 (No. 2). *R*
- A. W. Witmer & E. W. J. M. van Meijl:** The application of long wavelength radiation in X-ray analysis. Spectrochim. Acta 34B, 415-422, 1979 (No. 11/12). *E*
- J. P. Woerdman:** Dopplervrije tweefotonenspectroscopie. Ned. T. Natuurk. A 45, 30-32, 1979 (No. 1). *E*
- L. E. Zegers:** Moderne Signalverwerkingmethoden, ihre Realisierung durch Größtintegration und ihr Einfluß auf die Entwicklung der Telekommunikationsnetze. ntz-Archiv 1, 165-176, 1979 (No. 7). *E*
- H. Zijlstra:** Coping with Brown's paradox: the pinning and nucleation of magnetic domain walls at antiphase boundaries. IEEE Trans. MAG-15, 1246-1250, 1979 (No. 5). *E*
- Published in Conf. Proc. Optical Communication Conf., Amsterdam 1979 (5th Eur. Conf. on Opt. Comm. & 2nd Int. Conf. on Integrated Opt. and Opt. Fiber Comm.):*
- D. Küppers, J. Koenings & H. Wilson:** Influence of substrate temperature on the deposition properties for the plasma activated chemical vapour deposition process (paper 3.5, 4 pp.). *A*
- E. T. J. M. Smeets & J. Politiek:** Very low noise silicon avalanche photodiodes made by channeling of aluminium in $\langle 110 \rangle$ silicon (paper 4.1, 4 pp.). *E*
- G. D. Khoe:** New coupling techniques for single-mode optical fibre transmission systems (paper 6.1, 4 pp.). *E*
- G. D. Khoe, H. W. W. Smulders & A. J. J. Franken:** Demountable optical fibre connectors specially suited for realization using injection-moulded thermoplastics (paper 9.2, 4 pp.). *E*
- P. Geittner:** Dispersion measurements on PCVD optical fibres using a mode locked synchronously pumped dye laser system (paper 14.2, 4 pp.). *A*
- P. J. de Waard:** A novel single mode laser having periodic variations in the stripe width ('super DFB') (paper 18.4, 3 pp.). *E*
- U. Killat & G. Rabe:** Binary phase gratings for use in fiber optic communication systems (paper 21.2, 4 pp.). *H*
- K. Mouthaan & J. R. Schlechte** (Philips' Telecommunicatie Industrie, Huizen): A 140 Mbit/s optical transmission system with 8 km repeater spacing and line section length of 96 km (paper 22.4, 4 pp.).
- Published in Proc. Int. Conf. on Land mobile radio, Lancaster 1979 (IERE Conf. Proc. No. 44):*
- R. C. French:** Common channel multi-transmitter data systems (pp. 31-47). *R*
- M. J. Underhill:** Wide range frequency synthesizers with improved dynamic performance (pp. 171-182). *R*
- R. I. H. Scott & M. J. Underhill:** FM modulation of frequency synthesizers (pp. 183-191). *R*
- P. A. Lewis & M. J. Underhill:** Two stage tuning and matching of h.f. mobile antennas (pp. 213-222). *R*
- R. Wells:** The application of single sideband modulation in the 450 MHz and 960 MHz land mobile radio bands (pp. 291-298). *R*
- C. K. Davis & R. F. Mitchell:** Traffic on a shared radio channel (Suppl., pp. 341-349). *R*
- C. K. Davis & R. F. Mitchell:** Blind queueing — a signalling protocol for trunked, mobile radio systems (Suppl., pp. 351-357). *R*

Recent United States Patents

Abstracts from patents that describe inventions from the following research laboratories, which form part of or cooperate with the Philips group of companies:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brevannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

4 318 112

Optical recording disc

P. J. Kivits

M. R. J. de Bont

A. W. de Poorter

An optical recording disc which comprises a disc-shaped substrate plate, for example, a plate manufactured from synthetic resin, for example, polymethylmethacrylate, which has on at least one side an ablative recording layer, for example, an amorphorous layer of a tellurium mixture, in which the substrate on the side of the recording layer has extrinsic instabilities, in particular surface discontinuities. For example, the surface of the substrate is roughened or provided with scratches. The substrate may also be provided with a separate coating layer comprising the discontinuities, for example, an inhomogeneous vapor-deposited layer, a coarse-granular layer of a dye, or a light-cured lacquer layer which has been roughened or provided with scratches.

E

4 319 559

Solar collector for heating a gaseous heat transport medium

W. Hermann

H. Hörster

A solar collector for heating a gaseous heat transport medium comprises an open-top housing provided with a transparent cover layer formed of a plurality of adjacently arranged, substantially mutually contacting, sealed and evacuated transparent cover tubes. The inner surface of each cover tube is provided with a selective absorption layer over the half facing the interior of the housing. The surface of the housing bottom facing the cover tubes has a vaulted profile corresponding to the shape of the cover tubes. An inlet and an outlet are respectively provided in the opposite sidewalls of the housing for flow of the gaseous heat transport medium transversely of the cover tubes in direct thermal contact herewith.

A

4 321 454

Method of and welding torch for arc welding

G. A. M. Willems

G. W. Tichelaar

Plasma-MIG welding in which a thermally ionizable gas stream is flowed through a nozzle non-consumable electrode having a central orifice and a surrounding annular opening toward a workpiece and is thereby split into a central gas column enveloped by an annular gas shield. A consumable electrode is fed through the central gas column toward the workpiece, with the establishment of a MIG-arc therebetween. A plasma arc is then spontaneously established by means of the MIG-arc between the nozzle non-consumable electrode and the workpiece. The central plasma gas column is accelerated by constriction of the annular gas shield down-stream of the nozzle non-consumable electrode.

E

4 319 280

Apparatus for converting cinematographic pictures into video signals

J. Roos

K. Compaan

A. C. van Kasteren

An apparatus is described for converting cinematographic images into video signals. The film gate in which a frame to be scanned is located is illuminated uniformly. The film is scanned by means of a rotatable scanning mirror, which with the aid of an imaging system images one line of the film frame onto one row (three rows for color film) of radiation sensitive detectors. The detectors of a row are read sequentially. The illumination beam has a large aperture angle in a first plane transverse to the direction of film transport and a small aperture angle in a plane transverse to the first plane. The film movement is detected with the aid of a sprocket hole detection system, in which an elongate narrow light spot, whose longitudinal direction is transverse to the direction of film transport, is projected onto the strip of film with the sprocket holes. The position of the scanning mirror is detected with the aid of two gratings of special shape.

E

4 321 552

Amplifier comprising a first and a second amplifier element

N. V. Franssen

G. A. van Maanen

H. G. J. M. Kockelmans

An amplifier comprising a first and a second amplifier element, whose outputs are connected to a load, a difference circuit which comprises a comparator, in which the output signal of the first amplifier element is compared with the input signal thereof to produce

E



a correction signal which is applied to second amplifier element, a second substantially identical difference circuit which comprises a comparator in which the output signal of the second amplifier element is compared with the input signal thereof to produce a correction signal which is applied to the first amplifier element; a desired signal to be amplified being applied to both amplifier elements. The two amplifier elements, which may be preceded by pre-amplifiers, take the form of power amplifiers. Because the distortion signals appear at the load as "common mode" signals, the overall distortion becomes substantially zero in the case of equal attenuation in the two difference circuits. The power amplifiers may be provided with a voltage output as well as a current output.

4 322 624

X-ray tube having a magnetically supported rotary anode

G. A. A. F. Cornelissen
H. J. Gerrits
E. M. H. Kamerbeek

E

The invention relates to an X-ray tube. The tube has a rotary anode which is rotatably journaled by means of a magnetic bearing having a stator comprising a magnet yoke which intersects the outer walls of the tube and which is magnetizable by magnets arranged outside the tube. The magnet yoke comprises radially extending pole faces which enclose gaps in conjunction with radially extending pole faces of a rotor connected to the rotary anode. The rotor consists of a soft-magnetic disk, which provides the magnet yoke which closes the magnet yoke of the stator.

4 322 808

Coding and decoding artifact-free images of objects

H. Weiss

H

A method of coding and decoding of objects by means of a large number of point-like radiation sources which are subdivided into two groups. Two recordings are made of an image object. A first recorded image is made on a first recording medium with a first radiation source distribution. A second recorded image is made with a second radiation source distribution which is determined by the autocorrelation function of the first distribution. In the second radiation source distribution, the source which is determined by a function value in the origin of the autocorrelation function is omitted. In a first decoding step, the first recorded image is decoded with a point image function of the first radiation source distribution. In the second decoding step, the second recording is subtracted from the decoded first recorded image, so that an artifact-free image of the object is produced.

4 322 811

Clamping circuit for an adaptive filter

J. O. Voorman

E

In an adaptive filter having a delay circuit wherein amplitude control circuits are comprised in taps of that delay circuit, which control circuits are controlled by means of an error signal which indicates the difference between an output signal of the filter and an output signal which has been brought into a desired waveform by a threshold circuit of the filter, a clamping control circuit is used which utilizes the integrated error signal as the control signal for matching the level of the output signal to a threshold level of the threshold circuit so that unwanted direct currents do not affect the amplitude control circuits.

4 322 821

Memory cell for a static memory and static memory comprising such a cell

J. Lohstroh
C. M. Hart

E

A memory cell for integration into a static memory includes two transistors with cross-coupled base and collector regions. The collector regions are connected to p-n junction diode load elements having at least one region of polycrystalline silicon material. The collector regions of the transistors are connected to the regions of the diodes which are of the same conductivity type as the collector regions.

4 323 618

Single crystal of calcium-gallium germanium garnet and substrate manufactured from such a single crystal and having an epitaxially grown bubble domain film

J. P. M. Damen
J. A. Pistorius

E

A novel non-magnetic monocrystalline garnet substrate material in the form of calcium-gallium-germanium garnet. Single crystals of calcium-gallium-germanium garnet can be grown at much lower temperatures by means of the Czochralski method than single crystals of the conventional rare earth-gallium garnets. These single crystals are very suitable to epitaxially grow bubble domain films thereon, in particular films on the basis of $\text{Lu}_3\text{Fe}_5\text{O}_{12}$.

4 323 817

Color display tube

J. A. Vennix
C. Stapel

E

Reference studs on the outside of a cone of a color display tube, against which a system of deflection coils is to be placed, are manufactured from non-ferromagnetic material. The side of each stud which is secured against the cone consists of stainless steel, titanium or a titanium alloy to prevent interference with deflection fields generated by the deflection coils. The studs are secured to the cone by a glue which permits removal if necessary.

4 323 856

Injection laser

P. J. de Waard

E

An injection laser includes two substantially parallel mirror side faces and a substantially stripe-shaped contact member having a comparatively large width. The stripe-shaped contact member has a shape and disposition with respect to the mirrors such that laser action can be obtained only in a comparatively narrow stripe of the active layer. For example, the contact member may be arranged to extend obliquely with respect to the mirrors instead of at right angles thereto.

4 325 054

Folding circuit for an analog-to-digital converter

R. J. van de Plassche

E

In a folding circuit of an analog-to-digital converter a chain of emitters of transistors which are interconnected by threshold elements and fed by direct current sources are used to reduce the distortion. The circuit is controlled by a current source which produces the input signal.

4 325 084

Semiconductor device and method of manufacturing same, as well as a pick-up device and a display device having such a semiconductor device

G. G. P. van Gorkom

E

The invention relates to a semiconductor cathode based on avalanche breakdown in the p-n junction. The released electrons obtain extra accelerating energy by means of an electrode provided on the device. The achieved efficiency increase makes the manufacture of such cathodes in planar silicon technology sensible. Such cathodes are applied, for example, in cathode ray tubes, flat displays, pick-up tubes and electron lithography.

4 325 135

Optical record carrier and apparatus for reading it

J. G. Dil
P. J. Heemskerck

E

A record carrier is disclosed having an optically readable information structure comprising trackwise arranged information areas, as well as an apparatus for reading said record carrier. By using two

mutually perpendicularly polarized beam components for reading and by giving the information areas of adjacent track portions different geometries such that areas of one geometry can be read in an optimum manner by only one of the beam components and is virtually not observed by the other beam component, the track distance can be reduced without increasing the crosstalk, so that the information content of the record carrier is increased.

4 326 211

N+PP-PP-P+ Avalanche

E. T. J. M. Smeets

E

A radiation-sensitive semiconductor device includes a radiation-detecting avalanche diode which has a semiconductor layer structure made up of four layers of the same type conductivity. The fourth semiconductor layer is located above the third layer and has a higher doping concentration than that of the third layer. This fourth layer substantially improves the noise properties of the device, by a factor of about two. The radiation-sensitive semiconductor device is manufactured by a method in which the first and third layers of the semiconductor layer structure are provided by epitaxial growth, while the second and fourth layers of the structure are provided by ion implantation. The structure and method of the invention are particularly useful in the manufacture of avalanche photodiodes with an improved noise factor.

4 326 282

Apparatus for reproducing digitally coded information recorded on an optically readable disc-shaped record carrier

J. J. Verboom

E

M. G. Carasso

Disclosed is an apparatus for reproducing digitally coded information recorded on a disc-shaped record carrier in the form of optically detectable, unit information areas arranged in accordance with a concentric or spiral track pattern. The apparatus comprises a filter which is tuned to a frequency $f = V/L$, V being the nominal tangential speed of the record carrier and L the nominal center-to-center distance of said unit areas. This frequency component has a frequency equal to twice the bit frequency of the coded information and is situated at a zero point of the power spectrum of said information, so that the extracted signal may be employed for the generation of a clock signal.

4 327 299

Stepping motor

B. H. A. Goddijn

E

A stepping motor having a first and a coaxially disposed second annular stator section, which stator sections are axially spaced by a coaxially disposed permanent-magnetic ring. Each stator section terminates on the inner side in two coaxially disposed annular systems of teeth. The systems of teeth of each stator section have been mutually shifted nominally by half a tooth pitch relative to the rotor tothing and the systems of teeth situated on the outer side relative to the permanent-magnetic ring as well as the systems of teeth situated on the inner side relative to the permanent-magnetic ring exhibit a mutual shift relative to the rotor tothing of nominal a quarter of a tooth pitch. In this stepping motor stepping-angle errors can be nearly eliminated for a specific energization by selecting the teeth geometries in such a way that the amplitude of the permeances of the air gap between the inner systems of stator teeth and the cooperating rotor tothing as a function of the rotor position is smaller than the amplitude of the permeance of the air gap between the outer systems of stator teeth and the cooperating rotor tothing as a function of the rotor position.

4 327 354

Learning device for digital signal pattern recognition

E. H. J. Persoon

E

A device for recognizing overall patterns of digital signals which are arranged according to at least one coordinate comprises: first storage means for the storage of a pattern to be presented for com-

parison; a window device for selecting a local pattern of a number of digital signals from the overall pattern; a stepping mechanism for coordinate-wise adjustment of the window device over a step; a gradient determining device for determining, for predetermined coordinate values within the window, the absolute value of the coordinate-wide gradient of the values of the digital signals prevailing at said coordinate values; a moment generator for receiving these absolute values, for determining the moment therefrom with respect to a focusing point of the window, and for generating an output signal when a predetermined minimum value of said moment is reached in order to indicate the reaching of a local maximum in said coordinate-wise gradient; second storage means for the storage of a number of reference patterns; and a difference determining device for comparing, under control of the output signal of the moment generator, a local pattern selected by the window device with the reference patterns and for alternatively forming a non-correspondence/correspondence signal. If the difference is too large, the "new" pattern is stored as a further reference, so that the learning phase and the use phase occur simultaneously. The comparison does not take into account a difference near substantially coincident acute amplitude gradients.

4 327 963

Coupling element with a lens for an optical transmission system

G. D. Khoe

E

J. P. M. Gieles

G. Kuyt

A coupling element for an optical transmission system, in which the coupling element comprises a convex lens having a refractive index which is dependent on the radius r of a shell in the lens, and a holder. Use is preferably made of lenses having a refractive index $N_{(r)} = N_1 \cdot (2 - (R/R_0)^2)^{1/2}$, where N_1 is the refractive index of the core of an optical fiber (approximately 1.55) and R_0 is the radius of the lens. The coupling of monomode fibers via lenses of this kind can be effected with large dimensional tolerances.

4 332 078

Method of manufacturing a semiconductor device

H. L. Peek

E

M. G. Collet

In manufacturing a semiconductor device, a semiconductor body is first provided with a first insulating layer having a homogeneous dielectric thickness. A first conductor pattern of polycrystalline silicon is then provided on the first insulating layer. A second insulating layer is formed by oxidation of the first conductor pattern in such manner that the dielectric thickness of the first insulating layer remains approximately constant. Insulating paths are then formed in spaces below edges of the second insulating layer by successive deposition and etching steps. During the deposition step, a temporary layer is deposited to a thickness exceeding half the height of the spaces. During the etching step, the temporary layer is removed from the second insulating layer. Finally, a second conductor pattern is provided on and beside the second insulating layer.

4 333 022

Semiconductor device for digitizing an electric analog signal

M. V. Whelan

E

C. J. P. F. Le Can

An analog/digital converter implemented on an MIS structure having a gate electrode in the form of a resistance layer across which the potential gradient is applied. The analog signal is supplied to the gate electrode and converted into a shift of an inversion layer below the gate electrode. The digital signal is obtained by means of determining the number of surface regions which are present below the gate electrode and which can be electrically connected to the inversion layer.

4 333 158

Automatic gain control circuit for an adaptive filter

J. O. Voorman E

An adaptive filter having a delay circuit, taps of which comprise amplitude control circuits controlled by means of an error signal which indicates the difference between the output signal of the filter and a reference signal, includes a further control circuit which controls the amplitude of the reference signal such that an unwanted stable state for the filter is prevented from occurring.

4 334 128

Echo canceller for homochronous data transmission systems

W. A. M. Snijders E

Echo canceller for use in a homochronous data transmission system, comprising two-wire and four-wire connections and arranged for suppressing echo signals occurring in the four-wire connection. This echo canceller comprises an adjustable signal processing arrangement to which the data signal to be transmitted is applied and which produces a synthetic echo signal. A signal formed by a received data signal and an echo signal is present in the receive path of the four-wire connection. To generate a residual signal, the synthetic echo signal is subtracted from the signals in the receive path. For the adjustment of the signal processing arrangement there is added to this residual signal, outside the receive path of the four-wire connections, an auxiliary signal which is not correlated to this residual signal. The sum signal thus obtained is sampled with a suitably chosen sampling rate which is harmonically related to the symbol rate. The signal samples obtained are applied to a limiter circuit which converts each signal sample into a positive or a negative pulse, which is indicative of the polarity of the signal sample. The pulses thus obtained are applied as the control signal to an adjusting device for adjusting the signal processing arrangement.

4 334 139

Apparatus for writing patterns in a layer on a substrate by means of a beam of electrically charged particles

S. Wittekoek
T. A. Fahner E

An apparatus is described for writing patterns in a layer on a substrate by using a beam of electrically-charged particles. The apparatus is provided with an optical height-measuring system for determining a deviation between the desired and the actual position of the surface to be inscribed relative to the charged particle lens system. The deviation can be measured accurately and continuously without the use of additional markers on the substrate.

4 335 330

Low-pressure mercury vapor discharge lamp

R. C. Peters
L. E. Vrenken
W. D. Couwenberg
F. A. S. Ligthart E

Low-pressure mercury vapor discharge lamps provided with a luminescent layer and having a relatively high power consumption. Lamps are loaded by at least 500 W per m² surface area of the luminescent layer. In some embodiments the lamps have a nominal length of 60 to 150 cm and consume a nominal power of 0.25 to 0.50 W per cm length of the discharge tube.

4 339 688

Cathode-ray tube having a magnetic quadrupole shadow mask

J. Verweel E

A cathode ray tube for displaying colored pictures of the post-focusing type. The color selection means are formed by a plate which is provided with a large number of apertures, which plate is

provided on one side with a first plate of parallel strips of hard magnetic material provided between the rows of apertures. The color selection means are magnetized in a direction perpendicular to the plane of the plate so that a magnetic quadrupole is obtained in each aperture.

4 340 869

Amplifier for use in a line circuit

J. F. P. van Mil E

An amplifier for a subscriber's line circuit comprising an extra transistor whose base-emitter junction is connected across a resistor which is arranged in series with the main current path of the output transistor of the amplifier. Furthermore, a supply source is connected to the output terminal of the amplifier via the main current path of the additional transistor and an impedance for supplying at least a portion of the line current of a subscriber's line connected to the output terminal, in order to reduce the dissipation in the semiconductor elements of the amplifier.

4 341 010

Fabrication of electroluminescent semiconductor device utilizing selective etching and epitaxial deposition

R. P. Tijburg
T. van Dongen E

An electroluminescent semiconductor device such as a semiconductor laser has epitaxial monocrystalline layers, including an active layer, grown on a substrate. The epitaxial layers are etched in the presence of an etching mask to form nonplanar mirror surfaces which in the longitudinal direction bound active regions. To form flat and parallel mirrors an epitaxial monocrystalline protective layer is grown from the gaseous phase on the mirror surfaces after etching. The etching can be carried out in two stages using different etchants. With the first etchant the etched layers taken on a swallow-tail profile and then with the second etchant they take on a concave profile.

4 342 916

Method of and apparatus for tomographic examination of structures by X-ray or gamma ray scanning

M. Jatteau
V. Chalmeton
J. Pauvert L

A method of and an apparatus for tomographic examination by X-ray or gamma ray scanning of structures, intended for the display of images of slices of the structure examined by determination of the electron density in each volume element of these slices. Using an X-ray or gamma ray source, a stationary principal detector and an auxiliary detector which is aligned with the source, for each volume element of each slice two measurements are performed of the diffusion by the Compton effect towards the principal detector and two measurements are performed of the transmission between the source and the auxiliary detector, starting with the slice nearest to the principal detector and proceeding to the slice which is situated furthest from this detector. The measuring results are progressively used for calculating the exact electron density throughout the structure examined.

4 342 994

Display device having a liquid crystal

J. H. J. Lortetje E

In a display device having a liquid crystal the gate circuits for the row selection switches can be considerably simplified by coupling a terminal of the supply source for the selection switches to the ground conductor by means of an auxiliary supply source producing a periodical pulsating direct current voltage. To simplify the drive of column excitation switches a further auxiliary supply source can be used in a similar manner for a display device which drives liquid crystal display elements of the rms-type in time-division multiplex. All the required voltages can be obtained from one central supply source by means of at least one current source circuit.

Software

For many years it was self-evident that the functions of a piece of electronic equipment were determined by its physical characteristics, especially its electrical characteristics. However, this situation changed fundamentally when the digital computer appeared on the scene: this machine can be used for many different applications, provided it has been properly programmed for the particular application in hand. The hardware has become subordinate to the software: it is used to execute programs that are conceived entirely independently. A completely new discipline has emerged as a result of all this — the development and study of computer programs and programming techniques (sometimes called 'computer science'). From time to time Philips Technical Review has reported on Philips research activities in this field. In recent years, for example, we have published articles on experimental programs for electronic information storage, automatic picture processing, simulation of computer hardware and for extracting information from a data base by asking it questions in ordinary English.

Meanwhile, the more fundamental research studies on software have also attained a considerable significance for the electronics industry. Yet they remain somewhat underexposed, no doubt because of the rather intangible nature of this new discipline. Nevertheless, developments now under way in this field could well be of great importance to us all. A direct result of the steadily decreasing cost of the processors is that all kinds of new architectures for computer hardware are being investigated, so that the 'boundary conditions' of programming are in a state of flux. It seems as if in the future the development of new hardware and the development of new software techniques will have to go hand in hand.

With this background, it was clearly time to devote an issue of Philips Technical Review to the theme of 'software'. In addition to the two introductory articles on programs and programming, some systems are presented in which software is used as an aid or tool, while the final four articles give an impression of the problems that are now being tackled in software research. Interspersed with the articles there are a number of photographs of Philips products in which software is applied.

A great deal has been contributed to the production of this issue by Ir R. J. H. Scha of Philips Research Laboratories. He has been responsible for the editorial treatment of almost all of the articles: in close cooperation with the staff editors he has turned obscure and difficult material into articles that are readable and comprehensible for the non-specialist.



The display screens of the 'Signaal' Automatic Radar Processing System in the air-traffic control centre at Schiphol show the result of the successive operations by the various computers on the radar echoes. The simple radar 'blip' has been replaced by a clearly defined symbol, accompanied by the call-sign of the aircraft and the present and anticipated height. Air routes and approach areas are marked on the display. The flight path of the aircraft is extrapolated and the predicted distance in three dimensions from other aircraft is compared with the permitted distance. In fact, long before each flight started, a computer analysis was made of the schedules and flight plans. All of these operations require considerable amounts of software, and Schiphol therefore has its own specialized software centre.

Software

R. J. H. Scha

Introduction

Many different tasks can be performed today by one particular kind of machine: the programmable digital computer, or simply the *computer*. We see such machines performing impressively in many different areas. They calculate the orbits of rockets to the moon, keep our bank accounts updated and control the processes in chemical factories. They are also beginning to improve in areas that we still regard as typically human: they answer questions put to them in ordinary English and in a few years time they will no doubt win the world chess championship.

Since the invention of the electronic computer (at the end of World War II) computer applications have steadily increased in all sectors of business life. Meanwhile a new phase has emerged: computers have become so cheap and compact that they are now being sold and used not just as independent machines but also as integrated components of other machines or systems in which they perform a wide variety of functions. A notable example of this is the role played by the computer in modern telephone exchanges. But even in smaller consumer equipment, such as radio and television sets, we shall soon find more and more computer-type components.

Good use of the capabilities of computers is thus of vital importance to an electronics company. The computer and its applications have therefore long been a major field of study at Philips Research Laboratories. The research at Philips relates both to the computer as a piece of electronic equipment (the *hardware*) and to the programs (the *software*) that the computer requires to perform a particular task in a given situation.

This issue of Philips Technical Review deals with applications of computer software developed at Philips^{[1]-[3]}, and with new methods of computer programming now being studied at the various Philips laboratories^{[4]-[6]}. Although software is the theme of this issue, we shall not entirely neglect the hardware on which the programs run. Indeed, a strong interaction between ideas on hardware and ideas on

Ir R. J. H. Scha is with Philips Research Laboratories, Eindhoven. The author gratefully acknowledges important contributions to this article from Drs C. S. Scholten, a Scientific Adviser with Philips Research Laboratories.

programming is to be found in some interesting new developments.

This article introduces a number of basic concepts that are necessary for a proper understanding of the rest of the issue; it also outlines the background to problems and developments that are dealt with in greater depth in the articles that follow.

The computer

A computer is a simple machine, only capable of carrying out a limited repertoire of straightforward instructions. The operations of a computer consist of little more than the application of simple logical and arithmetical manipulations, reading and writing at particular locations in the computer memory, providing output via its output media and reading input data via its input media. Most computers in current use are *sequential* computers (other types will be mentioned later). These computers possess a central processing unit (CPU) which carries out the operations, and a 'memory' (storage) in which the program and the data are stored. The memory usually consists of a relatively small and fast main memory, and a relatively large and slow background memory. Display terminals, line printers, magnetic-tape units, etc. ('peripherals', which exchange data between the machine and the outside world) are also connected to the CPU.

Such a relatively simple machine can only perform its varied and complicated tasks because of the programs written for it. Different programs make a computer capable of performing different tasks — often quite complicated tasks, because computer programs can be extremely complex and powerful.

[1] J. H. A. Melis, O-BUS: a system for flexible public transport by means of on-call buses, this issue, p. 231.

[2] D. J. Burnett, INDA, a software tool for the production engineer, this issue, p. 237.

[3] W. E. Fischer, A data-base management system for CAD and CAM, this issue, p. 245.

[4] J. L. W. Kessels and A. J. Martin, Parallel programs, this issue, p. 254.

[5] B. L. A. Waumans, Software aspects of the PHIDIAS system, this issue, p. 262.

[6] A. J. Martin, Distributed computations on arrays of processors, this issue, p. 270.

Programs in machine language

In essence, a computer can only do two things: it can read a program into its memory and execute a program that has been read in. A program is a 'recipe' that defines a complicated task in terms of simple 'elementary' operations specified beforehand. A 'bare' computer, on which no programs are as yet running, must be programmed in 'machine language': the program must use, as its 'elementary operations' the instructions that the machine can in fact carry out.

It is possible to write such programs for a surprisingly wide variety of problems. The situation is rather like that in plane geometry, where one can also define very complicated constructions in terms of simple ruler-and-compass operations.

When a computer is actually carrying out a program it has become, in a manner of speaking, a different machine — for example a machine that reads in data, performs certain operations on it and displays the results on a screen. The program executed is a detailed specification of what this 'new machine' does. The program consists of instructions such as 'Load the contents of memory location No. 1234 into register No. 3' and 'Add the contents of register No. 3 to the contents of register No. 2'. These instructions are contained in a list, which is worked through step by step by the CPU. The instructions are executed one by one, in the order in which they occur in the list. The order may be altered, however, by special instructions. The list may contain 'branch' or 'jump' instructions stating for example: 'Now continue at point X in the program'. If at the end of a program a jump instruction is added, stating 'Continue at the beginning of the program', one can construct a program that continues forever. A further degree of complexity is achieved by means of 'conditional jump instructions', which transfer control to another program instruction if some specified condition is satisfied, e.g. 'Jump to X if condition Y is fulfilled, otherwise proceed directly to next instruction'. With instructions like this the execution of the program can be made to depend on results computed earlier (and stored in the memory). For instance, a sequence of instructions can be made to repeat itself until a particular result is achieved, or the behaviour of the program can be made to depend on data entered via an input device.

Programming languages

The first programs that were written when the programmable computer had just been invented were computer programs of the type described above. The 'language' in which they were written was 'machine code'. A program in machine code consists of a series

of binary numbers that identify the instructions and the data. This code is very cumbersome to work with, and therefore *assembly languages* were soon developed, which use combinations of letters, rather than binary numbers, to identify their instructions. This code is easier to memorize as it evokes associations with the function of the instruction (e.g. LOAD, STORE, JUMP). A program in assembly language cannot immediately be executed; it must first be processed by an *assembly program* or 'assembler'. A program written in assembly language is input to the computer as data and then translated into a program in machine code, which the computer can then execute.

A later development was the invention of programming languages of a higher level than machine languages and assembly languages. The instructions and structures in high-level languages do not correspond to the elementary operations and the sequential structure of the machine, but to concepts that are convenient to the programmer who wishes to describe a complex calculation, for example. Programs written in such a high-level language use a notation that is easier to understand and are simpler than programs in machine or assembly language.

An example of the way in which a high-level language makes programming easier can be seen in the execution of a command like 'PRINT(x)', which can be given in most programming languages. (Here 'x' is, say, the name of a numerical variable whose value we want to see printed.) This simple command triggers off a whole series of instructions, not only sending all the decimals of the number to be printed to the printer but also ensuring that non-significant zeros are suppressed, numbers are rounded off where necessary, a decimal point is inserted at the right place and that spaces are reserved after the number to separate it from the next one — while line feed is operated or another printed form is chosen as required, etc.

Programming languages have continued to develop through the years. In many respects, the new languages are less affected by the details of the operation of the computer. Instead, they are geared more towards the possibility of a simple and elegant description of the task that one wishes the computer to execute. If languages developed in consecutive periods are compared, e.g. FORTRAN, ALGOL 60 and PASCAL, this progress towards a 'higher level of abstraction' can be seen [7].

Programs at different levels

A program written in a high-level language cannot be executed directly by the hardware of a computer. The computer must first be provided with a *compiler*

for that 'source language'. A compiler is a program that accepts programs in a source language as input and translates them into programs in machine language, which can then be executed by the computer.

The program that supervises all the activities of this kind is known as the 'operating system' of the computer. The operating system accepts commands such as 'compile program A', 'execute program B', and so on. It also gives error messages if the compilation or execution of a program fails for one reason or another.

Programs thus exist at different 'levels'. Compilers and operating systems are called 'system software' — they are programs that 'belong' to the computer and make it a usable machine. Every programmer employs this system software to implement his own specific programs on the computer. We call these programs the 'application software' or 'user programs'. Such application programs, in their turn, often process input that is generated by an 'end-user' at a display terminal. The end-users need not be seen as programmers, however, because the data they produce as input for the computer is extremely simple.

The 'paper machine'

A computer that can be programmed in a high-level language is a remarkable device. It enables us to design all kinds of information-processing machines without having to think about the hardware. All we have to do is to produce a mathematical description of the behaviour of the desired machine, in the form of a computer program, and feed that program into the computer. A computer used in this manner is often referred to as a 'general-purpose computer'. This does not mean that the machine can be used for all purposes, but that its repertoire of primitive actions has been composed in such a way that, by providing the computer with the right software, we can convert it into a wide range of special-purpose machines.

The complexity of the machines that can be built with software in this way seems to be virtually limitless. The reliability of these machines should be equally phenomenal. Computers are now designed in such a way that the effect of unavoidable occasional malfunctions in the hardware can be almost completely neutralized. When the information is stored 'redundantly' (this is standard practice in all computer memories), error-correcting techniques can be used, with the result that the computer does everything it is instructed to do, with absolute obedience.

Nevertheless, computers do sometimes behave unexpectedly. They suddenly fail or give the wrong results. The explanation of this paradoxical situation

seems to be that in designing computer software one faces a kind of problem not usually encountered with 'hard' machines: the program designer finds himself incapable of specifying correctly what he wants the machine to do.

Problems

The programs written are often so complicated that the programmer does not have a good overview of his own programs. There are various reasons for this. In the first place many programs are extremely long. (A program of a million instructions, for example, is quite usual.) It is obviously not easy to coordinate such large numbers of operations, interrelated in so many ways. (The problem becomes even more acute when a single program is written by a team of people — the rule rather than the exception with long programs.) In the second place there is plenty of opportunity for the programmer to make errors of commission or omission. A calculation that has to be specified in all its details may turn out to be more complicated than was originally thought. The programmer must be prepared for any exception or borderline case that may occur in the execution of the calculation. If a computer has been told how to calculate a/b , for example, it will follow that instruction to the letter, even if b is equal to zero. Then, of course, it will produce no result or an incorrect one unless it has also been told to check to see whether b is equal to zero and, if so, to respond with 'instruction cannot be executed'.

Another problem in the design of software is a consequence of something that is in fact a great advantage of software over 'real' machines: a program is a 'paper machine', which can easily be modified and extended. It is therefore tempting to design a computer program by a process of trial and error, checking whether it works and making the necessary corrections until the whole thing seems to be functioning satisfactorily. In the design of a 'hard' machine, which cannot be modified simply by a stroke of the pen, one is less likely to adopt such a procedure.

This trial-and-error method for designing software has its problems. One of the problems is that program testing assumes an unduly important role in the design process. Unfortunately, it is usually impossible to determine by testing whether a program is working absolutely correctly. Since the number of situations that can occur in the execution of a complex program runs into astronomical figures, it is not possible to try out all these situations.

Then there is the added difficulty that applying corrections to corrections to corrections makes a pro-

^[7] F. E. J. Kruseman Aretz, *Abstraction*, this issue, p. 225.

gram increasingly dense and more difficult to understand. In practice the trial-and-error method never succeeds in eliminating all the 'bugs' in a program.

A more systematic method of design is therefore extremely important in the development of software.

The design process

Scientific research in software is not simply concerned with devising programs for new kinds of problems. As we have seen, it is also necessary to find ways of managing the complexity of the program — by developing techniques and aids that will enable programmers to write reliable programs.

Let us now look at the process we should follow to program a particular computer for a particular task. The first phase is the specification of the package of *requirements* that the program has to satisfy, i.e. the specification of the external behaviour of the computer that must execute the program.

This specification does not as yet say anything about the structure of the program. The basic structure of the program is developed in the second phase, called the *design* phase. A recommended procedure for the design phase is the method of 'stepwise refinement': the description of the program is first stated in terms of very general operations, and more detail is introduced later in stages.

The implementation of the program is the last step in such a process of refinement. (This step is rather special and may therefore be treated as a separate phase in the whole procedure.) During the implementation phase the algorithm that has been designed is described in an appropriate programming language, for which the computer has a compiler.

The program produced in this way can now be *tested*: it is entered into the computer and translated by the compiler into a machine-language program that the computer can execute, and tests are made to see if it is working correctly. (It has already been pointed out that too much importance should not be attached to testing: errors can indeed be discovered by testing, but is rarely possible to guarantee that the program is correct in this way.)

Finally, there is the *maintenance* of the program: once it has been put into operation, it may be necessary to correct errors discovered later, or to make small modifications at the request of the user.

The greatest problem in software, as we saw, is the problem of managing the complexity: how can we be certain that the very complicated programs we write are correct? The answer is that we should not make a program more complicated than is really necessary. The method of stepwise refinement in the design of

software is particularly useful since it produces programs that have a clear-cut structure. Sometimes it may even be possible for the design of a program to be accompanied by a mathematical proof of its correctness. If, on the other hand, the programmer goes into the details of implementation at too early a stage in the design process, the result will be an obscure and untidy program.

New developments

Now that processors are becoming smaller, faster and cheaper, it is becoming more attractive to think of computers whose structure differs from that of the sequential machines we have been dealing with so far. Nowadays we see more and more configurations in which a large number of processors cooperate in a common task: parallel and distributed machines. This change in hardware structure also has far-reaching consequences for the software.

In the established programming languages the programmer always writes sequential programs. Such programs specify separate steps in a one-by-one sequence. Languages are now being developed, however, for describing other types of programs. Such programs describe 'parallel' or 'distributed' processes — processes that can be performed simultaneously, largely independently of each other, but synchronized where necessary.

These new description methods are of great value when new kinds of hardware configurations are used, with independent processors operating in parallel. But if a parallel or distributed description method is more suited to the structure of the problems to be solved, it may also be useful to implement a language of this kind on a conventional computer.

Parallel programs have their own problems, however. For example, they can run into a *deadlock* — a situation in which the execution of the program is blocked. Such problems call for special techniques in the design of parallel programs. For example, methods are being developed for deriving correctly operating programs from programs that may run into deadlock^[8].

Hardware and software: two worlds?

A noteworthy effect of the availability of distributed hardware is that the traditional dichotomy between the world of the hardware designers and that of the programmers is to some extent disappearing. It is not difficult to see why such a dichotomy came about; at the time it did look as if the two activities were completely different. The programmer wrote 'sequential' programs (that is to say, he wrote machine instruc-

tions in the appropriate order), and took for granted that in the execution of the program a machine instruction would be completed before the next one was acted upon.

The hardware designer, who ensured that this was the case, was primarily concerned with an entirely different set of problems: the problems of synchronization. A logic circuit is usually built up from thousands of smaller circuits, i.e. from thousands of independent entities, each continuously adjusting its output to suit its input. The speeds at which this happens are not predetermined; only their minimum values are fixed in advance. To make such systems manageable the 'clock pulse' was introduced. (Some of the first machines worked without clock pulses, and there has been a revival of interest in systems of this kind recently.) However, the difficulties reappeared when hardware designers started to interconnect different circuits, each with its own clock pulse.

The programmer's raw material, the software, was therefore rather more manageable than that of the hardware designer. This led to two developments. First, the designs that the programmer could develop were orders of magnitude more complex and larger than those that could be tackled by the hardware designer. This led to a deliberate policy of keeping the machine as simple as possible. 'Anything can be programmed' was the watchword of the day. Secondly, in the software area it was much easier to start developing theories providing the programmer with a solid basis for his activities.

However understandable the dichotomy between the 'hardware' and the 'software' cultures may be, the problems encountered in the design of hardware and of software have a great deal in common. In both

cases complex entities are designed whose behaviour has to be determined on the basis of their individual components. (One possible fundamental difference is that the hardware designer, unlike the programmer, works with signals that are essentially analog signals. He has deliberately ignored this aspect, however, right from the beginning; this was the only way of applying concepts such as Boolean algebra, for example.)

In view of this historical background, it is an interesting development that the gap between the 'hardware' and 'software' cultures has narrowed considerably in recent years. VLSI (Very-Large-Scale Integration) has made it possible to have thousands of microcircuits all cooperating with a common objective. The question then arises: how can the programmer program these devices so that they will all cooperate in the correct way? To do this, he must leave the familiar ground of sequential methods and start to use parallel methods — formerly the exclusive domain of the hardware designer. The breaking down of these barriers has brought into view a problem area that always existed but was difficult to deal with: weighing up the relative merits of software and hardware for attaining a given objective. To ensure the proper balance one must think in software terms when designing hardware: one must think the problem through as far as possible in terms of a functional description of the task of the machine, and only *then* try to subdivide the machine into hardware components in the best possible way. The hardware that emerges at the end of this sort of design process may be called 'petrified software'.

[8] M. Sintzoff, Transformation methods for improving parallel programs, this issue, p. 278.



The PRX/D telephone exchange is completely digital — even the speech signals are converted into digital signals and then processed in time-multiplex form by the exchange. The program that must be written for the central processing unit, which controls the path of the signals, is extremely complicated and is different for each exchange because of differences in the telephone networks they serve. Adapting the program for a particular exchange can require one man-year.

Abstraction

F. E. J. Kruseman Aretz

Introduction

In the opening article^[1] of this issue it has been said that the most important problem in developing reliable programs lies in the control of their complexity — the programmer must be able to obtain a good general view of the complete process in spite of the multiplicity of operations contained in a program. An important technique for achieving this involves the use of 'abstraction': describing data and operations at a level that is 'high' enough for attention to be concentrated on the real problems, and not distracted by less relevant details.

The use of abstraction is nothing out of the ordinary. It is inherent to all mental activity; we abstract continuously without being aware of it. If someone says 'Don't forget to get your car licence renewed in time', we have a simple expression that conjures up a whole multitude of activities. While the actual details will not be particularly vivid in many people's minds, the overall effect, the acquisition of a new car licence, valid for another year or six months, will be clear to everyone.

A more detailed breakdown of 'Get your car licence renewed' might be: 'Find the form sent by the Licensing Authority, or get another one from the Post Office, and fill it in. Collect together the registration document, the certificate of insurance, the fee payable, the test certificate (if the car is more than 3 years old) and send or take them all to an appropriate Post Office or Vehicle Licensing Office'. All these details are themselves abstractions; they tend to describe the effect of a number of activities rather than indicate exactly how those activities should be carried out.

Prof. Dr F. E. J. Kruseman Aretz is with Philips Research Laboratories, Eindhoven.

[1] R. J. H. Scha, Software, this issue, p. 219.

This everyday example at once illustrates three properties of the concept of abstraction that are also important in a programming context.

— Abstraction is used when we describe something (a complex object or a complex activity) with the omission of details that are unimportant at the level of discussion. The intention is to make the truly relevant aspects more intelligible and more easily manageable.

— Often abstraction will be employed to describe the *function* of a machine or procedure, without actually specifying any of the details of its construction or implementation.

— It will often be useful to distinguish between different degrees of abstraction. A detail that is irrelevant at one level of discussion may be of fundamental importance at another level.

The most important function of abstraction, which we have already mentioned in passing, is that of making knowledge and data manageable by a kind of divide-and-rule technique. By proceeding one step at a time, concentrating on one part of the problem and considering just a few main features of the remainder, a complicated problem can be mastered gradually.

As has already been said, all human thought and work makes use of abstraction. We can use a television set without any knowledge of electronics and an electronic engineer can design circuits without knowing anything about quantum mechanics. The fact that we have no hesitation in calling two such entirely different objects as a coffee table and a dining table a 'table' neatly demonstrates that all concepts are in fact abstractions.

With the concept of 'table' there are borderline cases where it is not immediately obvious whether something should actually be called a table or not. Such uncertainties occur even more with concept

pairs such as alive/dead, animal/vegetable, healthy/sick or normal/abnormal. In the exact sciences (especially in mathematics), however, abstractions are sharply defined and abstract is not synonymous with 'vague'.

What does all this have to do with programming? This question can be answered in any number of ways. As it is impossible to give a complete picture of the part played by abstractions in the design of programs, we shall look at some important examples.

Numbers

A computer operates on strings of bits. A bit string can be used to represent many different kinds of entities, including integers. This is done in a number of ways that differ mainly in the way in which negative numbers are represented. The range of numbers that can be represented in a particular computer will always have an upper and a lower limit. In programming, however, we think and reason in terms of the entire set of integers without being aware of the machine representation and the restrictions introduced by it: we 'abstract'. It is therefore pointless for programming languages to introduce 'binary numbers' in addition to 'integers'. We want to have as little as possible to do with the representation of numbers in the machine.

Assembly languages

Machine instructions are also coded internally with strings of bits, and the first programmers had to know the machine code for each instruction. Symbolic notations were soon introduced for instructions and these developed into assembly languages. In this way it was possible to abstract from the actual coding in the machine. This abstraction excluded certain once-popular possibilities, such as the use of the same bit string in the memory to represent both an instruction and a piece of data operated on by instructions. But programming is difficult enough without tricks like this, and fortunately we no longer have to be so drastically economical with storage capacity.

Programming languages

The introduction of programming languages made it possible to write programs without being aware of the instruction set of the machine. A translation program (a compiler) translates the abstract constructions of the programming language into instructions for the machine on which the program is to be executed. As a result of this a program can be executed on any

machine that has a compiler for the language used in the program, so that programs become 'machine-independent'. The machine origin was still clearly recognizable in the oldest programming languages: the 'flow-of-control' structures had been taken from the conditional branch instruction and the procedure had a parameter mechanism based on parameter transfer with the aid of registers — if such a mechanism was provided at all. We are gradually beginning to forget the machine in the design of a programming language and now the emphasis is on the use and not the implementation. Only now are we beginning to find machine architectures designed to suit the programming languages, instead of the other way round.

This 'getting away from the machine' has been a difficult process and is still not quite complete. Compilers continue to generate extensive tables that show how the program is represented in the machine (e.g. where every variable and every label have been located in the memory). These tables are used for tracing errors: if something goes wrong during the execution of the program, everything stored in the memory of the machine can be printed and this 'dump' can then be interpreted by using these tables. Usually the errors are not in the computer, but in the program. It would therefore be useful if the programmer could consult a higher-level description of the state of the machine at a particular moment rather than a 'dump'.

Variables

Another example of a development that leads away from the physical machine is the way one thinks about the memory function. In assembly languages the memory of the machine is considered to be a collection of identical memory locations, which may be referred to by using a number, the 'address' of the location. (To make the programs easier to read, names may be assigned to locations.)

However, if we write a program in a high-level language, we are free to introduce symbolic entities that we use for storing information: these are *variables* in which *values* are stored. When introducing a new variable the *type* of this variable is also specified in the program. The type determines the kind of values that may be stored in that particular variable. It is possible to define variables that are intended to contain an integer, and also variables that can retain entire structures, such as the description of a transistor or of a complex organic molecule. The number of variables that a program uses may vary while the program is being executed, depending on the information storage requirement. It is the task of the compiler and of the program responsible for the management of the

memory to represent the variables in a linearly ordered memory of identical locations: in the high-level languages we have abstracted from this hardware structure.

Procedures

A computer program describes a complicated operation for performing a desired task in terms of a set of elementary operations specified previously. A comparable situation can be found in plane geometry, where complicated structures are built up from a small number of simple operations with ruler and compasses.

To perform the desired task, the computer program specifies a complex interplay of elementary operations that cannot be taken in at a glance. For this reason the algorithm is not built up directly from the elementary operations actually available; instead sub-tasks that might serve as intermediate steps are described at different levels of abstraction. This is another phenomenon that occurs in plane geometry. Here we find non-elementary activities such as 'construct the perpendicular from a given point P to a given line l '.

The specification of such an activity, which describes the effect obtained by its execution, has a dual function. For the person who uses the activity as a 'module' for performing more complex operations, the specification describes the net effect without indicating the way in which that effect is obtained. For the person who has to implement the activity the specification describes what the effect of the implementation should be, without going into the applications of the activity. The implementer must also base the 'correctness proof' of the program on this specification: just as with geometrical constructions he will have to show that his implementation has the desired effect.

The subtasks introduced in designing a program can be explicitly recognized in the program if they have been implemented by means of *procedures*. A procedure is a 'subprogram' that induces a certain effect on the variables to which it is applied (the *parameters*). We shall now give an example to illustrate how procedures are used for implementing subtasks.

Consider the subtask 'exchange the values of two variables'. The effect to be produced by performing this task can be described by means of statements about the condition beforehand and the condition afterwards. If for the variables u and v we have $u = a$, $v = b$ beforehand, then afterwards $u = b$, $v = a$. This effect can be obtained by means of a procedure, called *SWAP* here, which could be stated in PASCAL as follows (the components of the procedure are explained between curly brackets):

```

procedure SWAP (var  $x,y$ : item); {the procedure has two
                                parameters,  $x$  and  $y$ ,
                                of type 'item'}
var  $aux$ : item; {the procedure has an auxiliary variable,
                $aux$ , of type 'item'}
begin  $aux := x$ ; {the value of the first parameter is
                 assigned to the auxiliary variable}
       $x := y$ ; {the value of the second parameter
              is assigned to the first parameter}
       $y := aux$  {the value of the auxiliary variable
               is assigned to the second parameter}
end

```

If we want to exchange the values of two variables u and v somewhere in a program that contains this procedure declaration, we can simply write *SWAP* (u,v). What happens to u and v is described by the operations on x and y , respectively, in the procedure statement. If u has the value a and v the value b , then the result of *SWAP* (u,v) is therefore that u has the value b and v the value a .

The breakdown of tasks into activities is the most important feature of orderly software design; the programming language should support this form of abstraction. In the past the use of procedures generally led to an increase in computer time. Fortunately this has become less true, firstly because the parameter mechanism in modern languages is becoming increasingly simple and secondly because the hardware is beginning to adapt to the languages and consequently can process procedure calls efficiently.

The use of this form of abstraction in the design of software has a further aspect of great practical significance. All the details of the implementation of a certain subtask are concentrated in one section of the total program. This makes it relatively simple to replace the implementation of that subtask by another, without affecting the remainder of the program, which only makes use of the net effect of the procedure call. In developing a piece of software it is therefore possible to start off by implementing many subtasks in a simple or perhaps even inefficient way, and then look for the 'clever' solution later, if necessary. A search procedure, for example, may be implemented first as a slow linear algorithm and later as a quick tree-search algorithm.

Abstract data types

We have already seen that each variable used in a program is of a certain type that indicates the kind of values that can be stored in this variable. Programming languages contain elementary data types such as

integers, alphanumerical and Boolean quantities, and also means for building up complex types ('arrays', 'records', etc.) from elementary types.

In recent years the view has been gaining ground that types should not be characterized by the way in which their values are represented, but by the operations that can be carried out on the values. Types that are characterized in this way are called 'abstract data types'. The 'integer' type is thus characterised by the applicability of operations such as addition, subtraction, multiplication and comparison, and the 'Boolean' type by the applicability of the operations 'not', 'and' and 'or'.

This approach means that when a 'new' type is introduced a summary is given of the operations that can be carried out on a value of that type followed by a description of the features of those operations.

As an example it will now be shown how a new data type called 'stack' might be introduced. The idea is that 'stack' becomes the type for stacks made up of elements all belonging to an existing type that is called 'item'. The following operations are available for making stacks and for working on them:

CLEAR — make an empty stack,
 PUSH — add an element to a stack,
 POP — remove the top element from a stack,
 TOP — inspect the top element of a stack,
 EMPTY — check whether a stack is empty.

The definition of the new type 'stack', in addition to the name of the new type and of the existing types that have anything to do with it, includes a specification of the names of the above operations, the types of their arguments and results, and an axiomatic characterization of their properties. The specification of the type 'stack' is therefore:

new type: stack
 uses: item, Boolean
 operations: CLEAR: → stack
 PUSH: stack * item → stack
 POP: stack → stack
 TOP: stack → item
 EMPTY: stack → Boolean

characteristics:
 EMPTY (CLEAR) = true
 EMPTY (PUSH (s, i)) = false
 TOP (PUSH (s, i)) = i
 POP (PUSH (s, i)) = s

} for each *s* of type
 } stack and each *i* of
 } type item.

In this example the notation 'CLEAR: → stack' indicates that 'CLEAR' is a parameterless function that produces

a stack value (empty stack), while 'PUSH' is a function of two arguments (a stack value and an item value) that produces a new stack value. The characteristics show that a PUSH-operation on a stack value *s* and an item value *i* results in a non-empty stack value, which as a TOP element has the last added item *i* and assumes the original stack value *s* if a POP operation is applied. From the given characteristics others can be derived including the following important property: 'If *s* is a stack variable and *r* a string of operations, each of the form *s* := PUSH(*s*, *i*) or *s* := POP(*s*) for various values of *i*, and if *r* is such that in each initial section of *r* the number of POP operations does not exceed the number of PUSH operations, while for *r* as a whole these numbers are the same, then once the application of *r* is complete the value of *s* is the same as its original value'. (This property is easy to demonstrate by full induction on the length of *r*.) All that is necessary for the use of stack variables is the above specification of the type stack. Only in the implementation is it necessary to choose a representation for values of the type and program the associated operations in those terms. This programming has been correctly done if it can be shown that the operations implemented do indeed have the required characteristics.

In PASCAL notation an implementation of the type stack might read as follows:

```

type stack = ↑ aux;
    aux = record val: item; prev: stack end;
function CLEAR: stack; begin CLEAR := nil end;
function EMPTY (s: stack): Boolean;
    begin EMPTY := (s = nil) end;
function PUSH (s: stack; i: item): stack;
    var locs: stack;
    begin NEW (locs); locs↑.val := i; locs↑.prev := s;
        PUSH := locs
    end;
function TOP (s: stack): item; begin TOP := s↑.val end;
function POP (s: stack): stack;
    begin POP := s↑.prev; DISPOSE (s) end.

```

From the known characteristics of PASCAL constructions it is easy to see that the stack operations, implemented in this way, do indeed have the properties mentioned above. For example 'EMPTY (PUSH (s, i))' will certainly produce 'false', for if we have the procedure that implements the PUSH operation, we find:

```

NEW (locs) { locs ≠ nil }; ...;
{ locs ≠ nil } PUSH := locs { PUSH ≠ nil }.

```

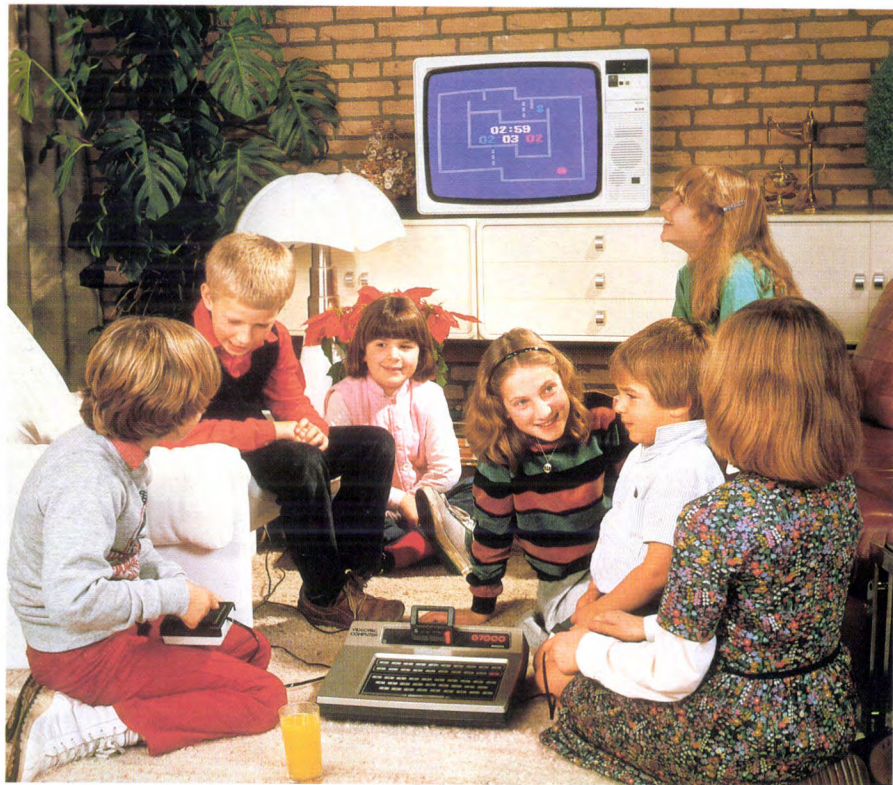
It would be desirable if the implementation of an abstract data type could guarantee that the only operations that could be carried out on it were those explicitly defined as applicable to that type. In that case the part of the program in which such a type was

defined would be properly 'encapsulated' with respect to the remainder of the program, so that no use could be made there of knowledge of details of the implementation. For the type 'stack' defined above this would mean, for example, that it would be impossible to discover the value of the element directly below the top of a stack s by making use of an expression such as $s↑.prev↑.val$. The separation of the specification and implementation would then be enforced by the programming language, with the significant advantage that it would then be possible to change to a different implementation without penalty by 'plugging in' a different software module.

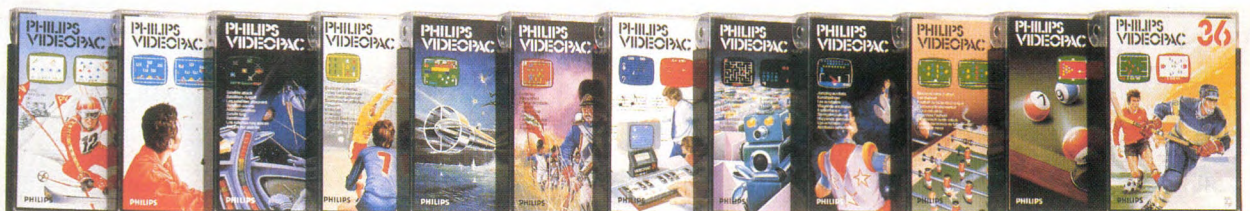
Unfortunately, neither PASCAL nor other current languages offer the possibility of encapsulating the definition of abstract data types in this way. Some experimental languages, and the new system programming languages ADA and CHILL, do offer this facility.

An attempt has been made above to illuminate some of the aspects of abstraction. It is strongly associated with a way of thinking, and an approach that plays a vital part in the design of software. Besides a certain aptitude, practice is necessary to acquire a style of working that profits from abstraction. Designers with a good mathematical background will have a very significant advantage here. The ultimate result of the application of abstraction will be order, structure, comprehensibility and correctness, and this will forge abstraction into a weapon against chaos and unreliability.

Summary. Using the correct levels of abstraction, so that less-relevant details can be ignored, is an indispensable technique in the design of clear and reliable programs. Modern programming languages offer increasingly far-reaching possibilities here. A number of examples are discussed, including the use of variables, procedures and abstract data types.



A limitless variety of television games can be played when a G 7000 'Videopac Computer' and two 'joysticks' are connected to a television set. The software for the different games is stored in the read-only memories of interchangeable cassettes.



O-BUS: a system for flexible public transport by means of on-call buses

J. H. A. Melis

Introduction

Most urban and regional public transport systems are based on fixed timetables worked out for the number of passengers expected. This type of system is not very flexible; it is slow to adjust to the fluctuations in the demand for transport and to the passengers' needs. The passengers are expected to adapt to the system, instead of the other way round.

Even the predictable disturbances to the system, as in the daily rush-hours, may mean that the passengers cannot even rely on having transport that follows the timetable to a reasonable approximation. Improvements within the existing system, such as providing more vehicles and reserving traffic lanes for public transport, are expensive and only solve part of the problem.

It would therefore seem worth finding out whether the system itself can be improved upon. The weak points in the existing system have already been mentioned. The transport pattern provided is rigid: it is controlled in a way that cannot allow for sudden unexpected changes in the numbers of passengers and the traffic conditions. More flexible control might be the solution. This does imply, however, that better means of communication are required between the vehicles, the passengers and the 'controlling' units.

The availability of telecommunication techniques and computers now enables improvements in this direction to be carried out. Research into these possibilities has been going on at the Project Centre of Philips Research Laboratories in Geldrop. The outcome of this investigation is a new system of public transport, known as the *On-call bus*, or O-BUS for short.

The characteristic feature of O-BUS is the high degree of operational flexibility it provides. This sys-

tem can be adjusted to suit transport requirements at any moment and it adapts smoothly when the operating conditions alter. O-BUS is particularly suitable for dealing with origin/destination relations that do not easily fit into a fixed pattern, either because there are many different relations or because they vary in time and place. In terms of method of operation, O-BUS comes somewhere between the taxi and the conventional bus and the intention is to optimize the service for the individual passengers. By providing vehicles in as rational a way as possible and by combining journeys, the costs are kept down and the greatest possible number of scattered destinations can be served by a limited number of vehicles. Easy local access to the transport is the aim, without introducing long waiting-times or other delays.

With flexible control of the vehicles, the available transport capacity is used to the full and on average passengers will occupy seats in the vehicle for shorter periods than with fixed-timetable operation. The result of this is that relatively small manoeuvrable vehicles can cope with a large demand for transport.

After a short description of the basic system, the control hierarchy and the communication systems will be discussed in greater detail. Then follows a description of what passenger and driver actually have to do. Finally, quantitative results will be given that were obtained from a simulation experiment based on practical transport data.

Outline of the system

In a system of public transport based on the O-BUS principle, the journeys do not follow a fixed timetable. The schedule is flexible and can be adjusted at any time to changes in circumstances. As soon as any new passengers request transport, the picture of 'pas-

Ir J. H. A. Melis was formerly with Philips Research Laboratories, Eindhoven; he is now with Philips Nederland, Eindhoven.

senger demand' alters. Then it is necessary to decide which buses can be used to transport these passengers and how to alter the routes of these buses so as to satisfy this new demand in the best way. The general procedure is as follows:

- New passengers make contact with the control centre (by a calling system) and give their destination.
- At the control centre the best vehicle to deal with the request is selected, with the best route for that vehicle. The existing pattern is changed as little as possible. Information about the 'status' of each vehicle in operation is taken into account, such as its location and number of passengers. The control centre must always be aware of this information.
- The callers are given information about the reservations made, the drivers are given their new instructions and the system controller can see how the system is behaving at that moment and in the immediate future.

If the level of service or the efficiency of operation makes it necessary, the system controller can adjust the system: he may use more or fewer vehicles, for example, or alter the optimization criteria.

Control hierarchy and communication systems

O-BUS is a complex system in which decisions have to be taken at various levels, in the long and short term, based on a number of objectives and depending on various criteria. At all levels information will come into the system from outside (for example local-authority plans, requests from passengers and prospective passengers) and all this information will have to be processed. Information will of course be continually exchanged between the various component parts of the system. The various subsystems included in the pattern of communication and control are shown in *fig. 1*.

The first control level relates to the direct control of the individual vehicles. This is called Vehicle Control. The 'control' is in human hands; the important quantities here are the status (location, number of passengers) of the vehicle itself and its immediate surroundings (the local traffic situation).

The communication system NETwork COMMunication (NETCOM) provides direct communication between the driver and the situation on the road. NETCOM has been included as a separate item in the diagram in *fig. 1*, to emphasize that here again there are electronic aids such as beacon transmitters and systems that control priority at traffic lights in response to vehicle signals picked up by inductive loops in the road surface. Although these aids can be of great importance, they are not directly characteristic

of O-BUS. The driver's personal observation of the traffic situation is also part of NETCOM.

The communication system PASsenger COMMunication (PASCOM) provides communication with the passengers (boarding, alighting, seated) for verification and assistance during the journey.

At a higher control level, Vehicle Management, the sum total of 'fleet operations' is optimized. Here the routes, stops and allocations of vehicles to prospective passengers are constantly adjusted to changing circumstances. This control level must therefore be kept informed at all times of the current situation (location of all vehicles, number of passengers, transport requests, etc.). For 'following' and 'controlling' the individual vehicles, Vehicle Control and Vehicle Management communicate to one another by mobile radio. This communication system is called VEHICLE COMMunication (VECOM).

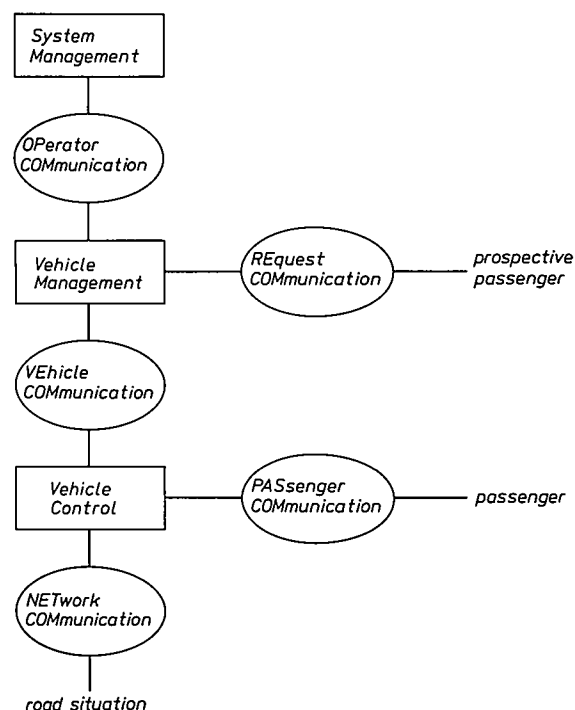


Fig. 1. Subdivision of a transport system into subsystems for control and communication. The 'boxes' in the diagram represent subsystems for control; the 'ovals' are for communication. There are three levels of control: control of individual vehicles (Vehicle Control), control of the entire fleet (Vehicle Management), and the decision-making for the operation of the entire system (System Management). System Management communicates with Vehicle Management (via Operator Communication), Vehicle Management communicates with Vehicle Control (via Vehicle Communication), and Vehicle Control communicates with the actual environment of the vehicle (via Network Communication). Passengers report to Vehicle Management via Request Communication. On boarding and during the journey the passengers communicate with Vehicle Control via Passenger Communication.

Prospective passengers request transport from Vehicle Management; this is done via the communication system REquest COMmunication (RECOM), the subsystem for calling and for reserving seats.

At still higher control levels the coordination of various transport subsystems (train, bus, boat) and longer-term adjustments (number of vehicles, expansion of the network, etc.) may be necessary. In the diagram in fig. 1 these control levels are brought together under the name System Management.

System Management must be kept informed of the activities of Vehicle Management, and must be able to make changes to them. The communication between System Management and Vehicle Management is indicated in fig. 1 by the communication system OPERator COMmunication (OPCOM).

The control systems

O-BUS is based on the existing network of roads and on existing vehicles. The individual vehicles are under human control. At the control level called Vehicle Control here, therefore, O-BUS is not very different from a conventional system. O-BUS simply has better communications, enabling the driver to do his job in the simplest and best way.

Most of the tasks in System Management are not automated, though decision-making at this level will again make increasing use of electronic aids.

The most fundamental difference between O-BUS and a conventional transport system is at the level of Vehicle Management. Here the schedules for all the vehicles are constantly revised. This is such an enormous task that it can only be done properly by a computer. The 'real-time' nature of the Vehicle Management decisions requires that an optimum operating schedule is available at all times. Since this schedule also depends on a great number of parameters and variables, a heuristic optimization method has been used (i.e. one in which the best result is constantly updated from a series of successive estimates). This means that in busy periods a quick estimate of the best transport schedule is accepted, whereas in quiet periods there is time for further optimization.

The optimization criterion implies that for every change in the schedule *the increase in the total estimated 'inconvenience' to all passengers* is minimized. The 'degree of inconvenience' is calculated from the waiting times, detour times and intermediate stop times at bus stops. The degree of inconvenience increases progressively with the length of these times. Also, various weighting factors are allocated to the various times. It is assumed, for example, that waiting time is much more inconvenient than other delays.

The weighting factor of the 'detour time' also depends on where the detour starts and on the length of the route originally planned between the start and end of the detour. (If passengers are almost at their destination and see that the bus is making a detour, they will find this extremely inconvenient.)

The transport schedule is constantly changed. It is not necessary, nor is it desirable, to keep all the vehicle drivers informed of the entire schedule at all times. For this reason each driver is only given definite instructions that concern him directly.

The communication systems

In a system like O-BUS, direct communication between the prospective passenger and the system is a primary requirement. Although O-BUS is potentially suitable for 'house-to-house' transport, with the buses called by telephone, the initial intention was for 'stop-to-stop' transport with call facilities at each bus stop.

The O-BUS subsystem RECOM allows buses to be called and confirms reservations. It can also handle the payment by the passenger to the transport company. RECOM in the O-BUS system is based on the use of fully electronic user cards, which can serve both as a ticket and as a credit card.

At the bus stop the passenger contacts control by placing his card on a 'tray' on the 'call-post'. The information is transferred between card and terminal by electromagnetic induction. This means that robust, and relatively simple terminals can be used. The transport requests, such as destination and number of passengers, are made by means of an ordinary voice link. When the prospective passenger's requests have been processed, the details of interest for the passenger are written into the card electronically. This card, which contains an electronic memory and a small display screen, then carries the necessary transport information, such as the destination, the number of reserved seats, and the number of the bus on which the reservation has been made.

As well as semi-automatic calling and reservation, the possibility of fully automatic reservation has also been provided for. The same type of bus-stop terminal is used for this. The additional facilities for fully automatic reservation are therefore provided on the card and not at the terminal. The 'automatic card' has a mini-keyboard for pre-programming it. This reduces the access time to the system and increases RECOM's processing capacity. A pre-programmable card offers several possible facilities, such as protection against unauthorized use of the card. The user can activate the card by keying in a code number, which is compared with a built-in code. This check is carried out by the card itself and does not have to be made at the terminal. The card can also be deactivated and is then of no value to unauthorized users.

Even in a transport system with regular timetables there may be delays en route. It is therefore necessary to monitor the progress of the individual vehicles so as to be able to control the punctuality and regularity of the transport. To achieve the flexibility that O-BUS is intended to provide, it is not only desirable but necessary for control to always have a good picture of the current situation in the field. Information must therefore be transmitted frequently from each vehicle to control. Conversely, control must from time to time communicate new routing instructions to each driver. The communication channel that enables this to be done is VECOM. The transfer of signals between control and the vehicles is done by mobile radio, so that contact can be maintained no matter where the vehicle is. To enable a relatively large flow of messages to be sent back and forth without overloading the drivers, a system of digital mobile radio was used. VECOM maintains contact between control and the vehicles completely autonomously. Only when there are new instructions or when changes are made to the schedule is contact made with Vehicle Management. The driver can request information relevant to him whenever he wishes. He does this from a driver's terminal in the vehicle, which he can also use for sending messages to control.

For communication with the passengers while boarding the vehicle and during the journey (PASCUM), a number of facilities have been installed to carry out electronic card inspection and to give passengers information during the journey. This is done with the aid of a liquid-crystal display inside the vehicle.

On the front and side of the outside of the bus the vehicle number is visible on a screen. This vehicle number can be changed whenever necessary. If the number of the bus corresponds to the number on the electronic card, the waiting passengers at the bus stop knows that this is the right bus. Vehicle Management can change the vehicle number of any bus. This means that if circumstances require it a 'better' vehicle can be designated right up to the last moment.

Vehicle Management is kept constantly informed of the whereabouts of the individual vehicles mainly by the drivers reporting their arrival at the stops (by pressing a button in the bus). The distance travelled between the stops is measured by a special odometer. For calibration, a number of beacons are placed at strategic points along the routes. These beacons emit signals (beacon numbers), which are received by the bus as it passes. The combination of beacon and odometer information could be used for monitoring the vehicles entirely automatically.

Vehicle Management always has complete information about the transport situation (locations of the

vehicles, number of passengers, number of passengers waiting, etc.). This information can be displayed via the subsystem OPCOM. This subsystem also processes this information statistically: the results can be used for assessing the long-term behaviour of the system. If necessary, OPCOM can then be used to adjust the operation of the system by altering the optimization criteria or by having more or fewer buses in service at certain times of the day.

O-BUS in practice

The operation of O-BUS in practice will now be described with the aid of the photographs of *fig. 2*. At the departure stop the passenger selects the number of his destination (*fig. 2a*) and calls control by placing his electronic card on the tray (*fig. 2b*). The call-post 'reads' the card and contacts control (*fig. 2c*). The passenger can hear the card being read through the loudspeakers in the call-post, so that he knows that the connection has been made. The passenger then asks for his journey by speaking into the microphone and the staff at control input the details to the system. These details are immediately shown on the electronic card. If they are correct, the passenger presses the green button on the call-post. The number of 'his' vehicle appears on the card immediately afterwards. This is the number shown on the front and side of the appropriate bus (*fig. 2d*).

At a suitable moment the driver of the bus can display a new instruction on his driver's terminal, such as the next bus stop or intermediate route-information point (*fig. 2e*). For communication with the passengers there is also a liquid-crystal display in the vehicle (*fig. 2f*). Information about the route can be shown on this display.

O-BUS simulation

To obtain some quantitative information that could indicate the extent to which a system like O-BUS would function better than a system with a regular timetable, we made a simulation experiment. The initial data was taken from figures supplied by the Philips internal personnel-transport service, which operates between the various Eindhoven plants. This service is provided by a number of minibuses, which operate to a fixed timetable.

The simulation was based on previously measured transport and network information for a service between 24 stops, spread over an area of about 20 km².

[1] J. H. A. Melis, Design aspects and potential characteristics of a demand responsive bus system, *Electronique + 5*, Coll. Int. Paris 1977, pp. 129-137.

A complete description of this experiment is beyond the scope of this article [1]. The results are illustrated in *fig. 3a*. The average 'lost' time per passenger (waiting time, intermediate stopping time and detour time) is plotted vertically to represent the 'inconvenience'. The solid lines show the average delay as a function of

the number of passengers per hour to be transported, with the number of buses available as a parameter. The dashed lines give the same relationship for operation to a fixed timetable.

Some of the details that can be read from this figure will now be given. If there are 8 buses in service,



Fig. 2. *a*) The passenger arrives at the O-BUS stop and selects the 'number' of his destination with the aid of the town map. *b*) He puts his 'electronic card' on a small 'tray' and asks O-BUS control for transport. The reservation is made and the number of 'his' bus appears on the card. *c*) An operator at control, who passes on passengers' requests to the computer system. *d*) The bus arrives at the stop. Note that one bus may have several numbers which may change from stop to stop. *e*) The driver's terminal in the bus. The driver receives route instructions from the control centre via the screen. *f*) A display screen, based on liquid-crystal technology, for showing destinations inside the vehicle.

the average delay at 200 passengers per hour is about 8 minutes with O-BUS, compared with more than 12 minutes with conventional operation. For a correct interpretation of these data, they should be viewed in relation to the average shortest journey time (ex-

cluding any delay), which is approximately 4 minutes. It can be seen from the same figure that a supply of 300 passengers, with a delay time of 10 minutes, can be dealt with in a conventional system by 15 buses. O-BUS requires 8 buses for this.

In fig. 3b, large 30 passenger buses are used for the fixed schedules, and 10 person mini-buses are used for the on-call system. The thick line shows that in this case both systems can offer the public the same service. Below the thick line, i.e. at lower numbers of passengers, transport by O-BUS is clearly faster.

In situations which are comparable with the case under examination, a system of on-call buses can therefore offer considerable advantages regarding quality of the transport and the number of vehicles required. These situations will mainly be found in urban areas with many transport connections between widely scattered stops. Sometimes, however, there are existing transport patterns, which have often evolved historically, in which mass transport is concentrated in only one or a few routes (for example in underground systems). Here, too, on-call buses can be used to advantage for transporting passengers to and from the stops on the fixed route. The control and communication methods that have been developed in O-BUS are not only useful for public passenger transport, of course. The O-BUS system, or parts of it, can be used to optimize many kinds of transport, for passengers or freight.

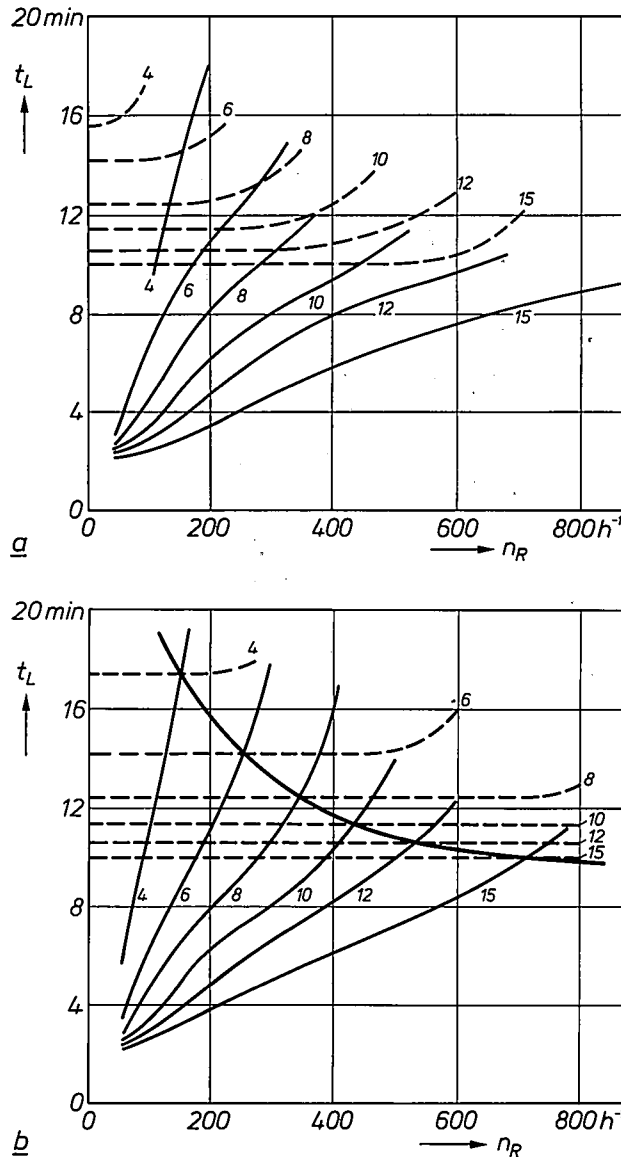


Fig. 3. a) The average delay t_L plotted against the average number of requests n_R per hour for various numbers of vehicles in service. Each vehicle can carry up to 15 passengers. Solid lines: O-BUS. Dashed lines: fixed timetable. For the same number of vehicles, O-BUS transport is faster than conventional transport. b) The average delay time plotted against the number of requests per hour. Solid lines: O-BUS, with vehicles for up to 10 passengers. Dashed lines: buses for up to 30 passengers, operating to a fixed timetable. The thick line connects the points at which both systems offer the same service. Below the thick line O-BUS is more efficient.

Summary. O-BUS is a system of public transport using on-call vehicles. Compared to conventional systems for urban and regional public transport, O-BUS is characterized by highly developed systems for control and communication. There are three levels of control: Vehicle Control, Vehicle Management and System Management. Subsystems for communication are Network Communication, Vehicle Communication, Operator Communication, Passenger Communication and Request Communication. These subsystems are described briefly. More attention is devoted to the call and reservation system, since this is a characteristic feature of O-BUS. Finally, a simulation experiment based on actual transport data is described. This showed that a system such as O-BUS can be more efficient than a conventional transport system, particularly in areas with many routes between scattered stops.

INDA, a software tool for the production engineer

D. J. Burnett



Flexible automation

On many mechanized assembly lines and in many modern machine shops the production engineers need to be able to develop different kinds of supervisory and control programs for their 'robot' machines. These programmable machines perform a variety of tasks, provided that the appropriate software is available. The engineers must be able to test the existing software and correct it if necessary, and sometimes they will feel the need to modify or extend the software. Because of their knowledge of the manufacturing process and production conditions the engineers are also the people who can most readily recognize the consequences of a change in the specifications of a product and will have to be able to adapt the software accordingly. Many production engineers, however, will not be expert at writing computer programs and they therefore have need of support in this area.

INDA — a name derived from the words 'INDustrial Automation' — is a software system that offers the

facilities for such work and enables the production engineer to solve the software problems himself. It is easy to learn to use INDA, and it does not take long. INDA is fast and flexible in a broad field of applications, particularly when it comes to testing, correcting and modifying software.

The wider use of programmable robots, an important part of 'flexible automation', is a subject now attracting keen interest^[1]. To people on the factory floor it may look quite complicated. Indeed, it does require advanced knowledge of various subjects, such as electronics and mechanical engineering, general cybernetics, picture processing, and computer programming. All this is in addition to a knowledge of and expertise in production engineering, of course. This knowledge will probably be available in a scientific and technical research institute, but is unlikely to be found in the factory — the very place where it will be regularly required.

D. J. Burnett, M.A., is with Philips Research Laboratories, Redhill, Surrey, England.

^[1] See for example J. G. van den Hanenberg and J. Vredenburg, An experimental assembly robot, Philips tech. Rev. 40, 33-45, 1982 (No. 2/3).

In the factory the best solution is to standardize as much as possible of the required knowledge in 'packages' that must be easy to use by the non-specialist engineer. Robot dynamics — for example the movement of a mechanical arm — and the picture-processing operations that the robot has to perform can usually be prepared beforehand as packages (sub-routines) in a 'library of functions', a collection of

work to generate the application programs. It may be necessary to do this very rapidly. The knowledge of a versatile yet simple programming system will then be a great advantage.

With INDA the engineer has the use of a very flexible high-level programming language. It offers him the means of adapting and adjusting ('tuning') existing programs, or constructing a completely new program.

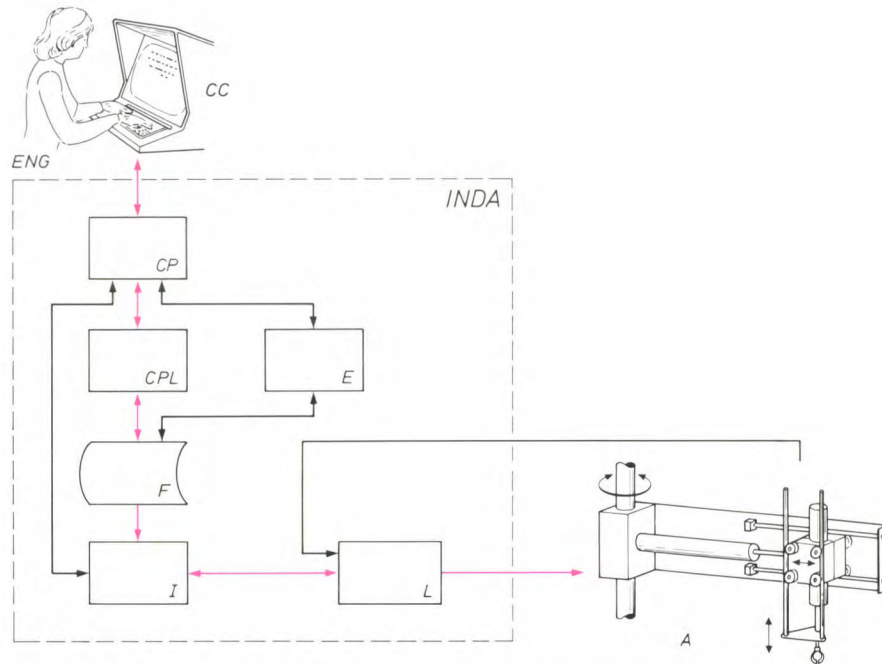


Fig. 1. Block diagram of INDA, an aid system for making, testing and modifying application programs for flexible automation. The area bounded by the dashed line contains the software modules that form INDA. *Eng* represents a production engineer, who wants, say, an automatic series of pick-and-place operations to be executed by a robot *A* equipped with a gripper arm and a television camera. *CC* control console with video screen for communication between *Eng* and *A*. *CP* command processor. *CPL* compiler. *F* data base with work files on magnetic disks and data tables in the store of the computer on which INDA runs. *I* interpreter, which executes the program and aids debugging. *L* library module, which contains the procedures that determine the selected area of application ('window to the world'). *E* editor, which helps with the correction of errors and the incorporation of the program modifications desired. *Red arrows*: main programming route. The first version of the application program is stored in *F* without employing INDA.

general programs, for use by any production engineer. Although they provide a definite improvement, standardized packages are not the complete solution. The production engineer still needs more software, e.g. to link the standardized packages together and to adjust parameters used in the packages to appropriate values for the task in hand.

Every task performed by a production machine in fact requires a *specific* program; this program defines the task by describing, for example, *which* functions are to be selected and just *how* they are to be used. Such an 'application program' must be produced for each new task to be performed by the machine. In the modern factory it is part of the production engineer's

The advantage of the high-level language and the flexibility of INDA is that the man on the shop floor can produce working programs without too much difficulty, even though he may be limited in specialist knowledge.

The part of INDA known as 'Interpreter' enables the engineer to run a program — with all the production hardware connected in, of course — while a variety of 'debugging' aids help him to check the execution. The various facilities offered by INDA enable the engineer thus to locate program errors.

Another part of INDA, called 'Editor', provides the means for piecemeal correction of errors, i.e. correcting them one at a time where they occur. Besides

correcting errors, it is also possible to make small modifications. On-the-spot editing obviates the need to go back to the very beginning of the program if its execution can be continued; this is very important. (The traditional way of debugging software has the effect of introducing a new program after each correction; this new program has to be re-compiled and rerun from the beginning.)

Why is it so important to avoid the need to go back to the original starting point? The reason is that this first run of a program (for example for an assembly operation), during which the engineer is continuously on the alert and is on the lookout for errors, often takes several hours — although later when everything has been corrected the entire work cycle may only last a few minutes. The costs of corrections can therefore be reduced considerably if there is no need to repeat parts of the work cycle that have already been corrected. This is particularly true of the many parts of the run that are involved in the first phase of the setting-up (or re-setting-up) of the machines, operations such as the positioning of manipulators, setting up stocks of parts, mounting tools, adjusting duplicating mechanisms.

The following sections describe the general structure of INDA and give a flow chart showing how to use it. Details are given of the processing of the control signals (via the command processor as the interface between the operator and the machine), the ‘Compiler’ module and the language, the ‘Interpreter’ module with debugging aids, the ‘Editor’ module for the local (definitive or temporary) execution of corrections and modifications, and the library of function procedures. To help clarify the situation, the example chosen to illustrate all these components of INDA is a camera-equipped robot that picks up rotors and puts them on to the shafts of small electric motors, as part of an assembly operation in the production of large batches of the motors (an imaginary case). The title photograph shows an experimental robot arm (the horizontal ‘box-like’ structure) built by Philips Research Laboratories at Redhill. INDA is being used here to program this robot.

General structure and operation

Fig. 1 shows the overall structure of INDA and fig. 2 gives a flow chart showing how to operate it. The area inside the dashed line in fig. 1 represents a single integrated program that can perform all of the three stages in the development of software that will run without errors. The first stage is the writing of an application program, while the second and third stages are the running of the program and the correc-

tion of the errors in it by means of the editing program. The command processor acts as an interface between the engineer at the control console and the other modules. There is a common data base with work files on disks and data tables in a store. The ‘Library’ module contains the standardized procedures for the various application areas. It provides the link with the actual application world. After writing

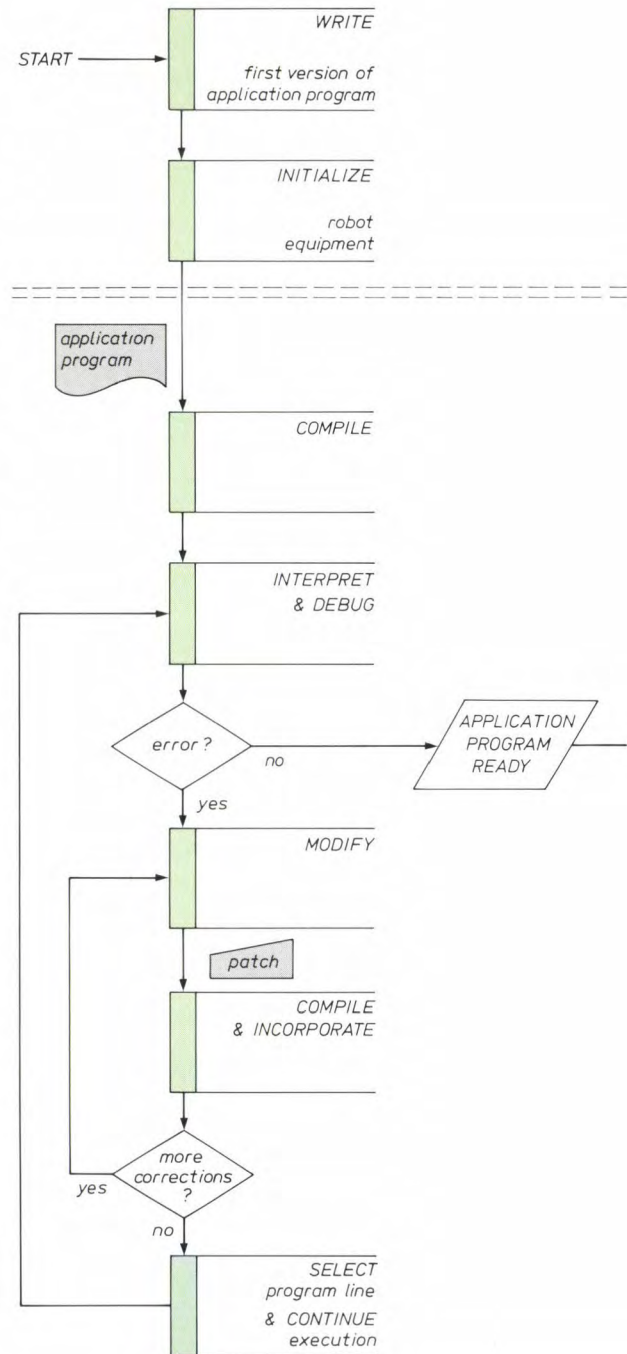


Fig. 2. Using the INDA system for making an application program. The subtasks are illustrated in a kind of flow chart of task-descriptors (the open blocks with green edges). In grey: material for the application program. Programming with INDA starts below the double-dashed line. The program can be run on a Philips P856/7 minicomputer.

and storing the first version of a new application program and setting up the production equipment, the engineer can start running the INDA software as indicated below the double-dashed line in the flow chart (fig. 2).

The engineer first calls on the command processor to run the compiler, which checks the syntax of the application program entered and translates it into machine language; the result of the compilation is added to the work files and the data tables. The interpreter is then called upon to perform the actions specified in the program. At the same time the engineer locates any errors and other deficiencies in the pro-

manent change of the program structure, is compiled and incorporated into the work files and the data tables.

Since most of the program usually remains untouched, it is both possible and sensible to continue the program either at the point where it was stopped or at some other arbitrary point within the program. Having decided where to recommence, the engineer resumes the checkout and continues the sequence of corrections until the program behaves as desired. This completes the programming work of the production engineer, at least as far as INDA can support him with the interpreter and the editor.

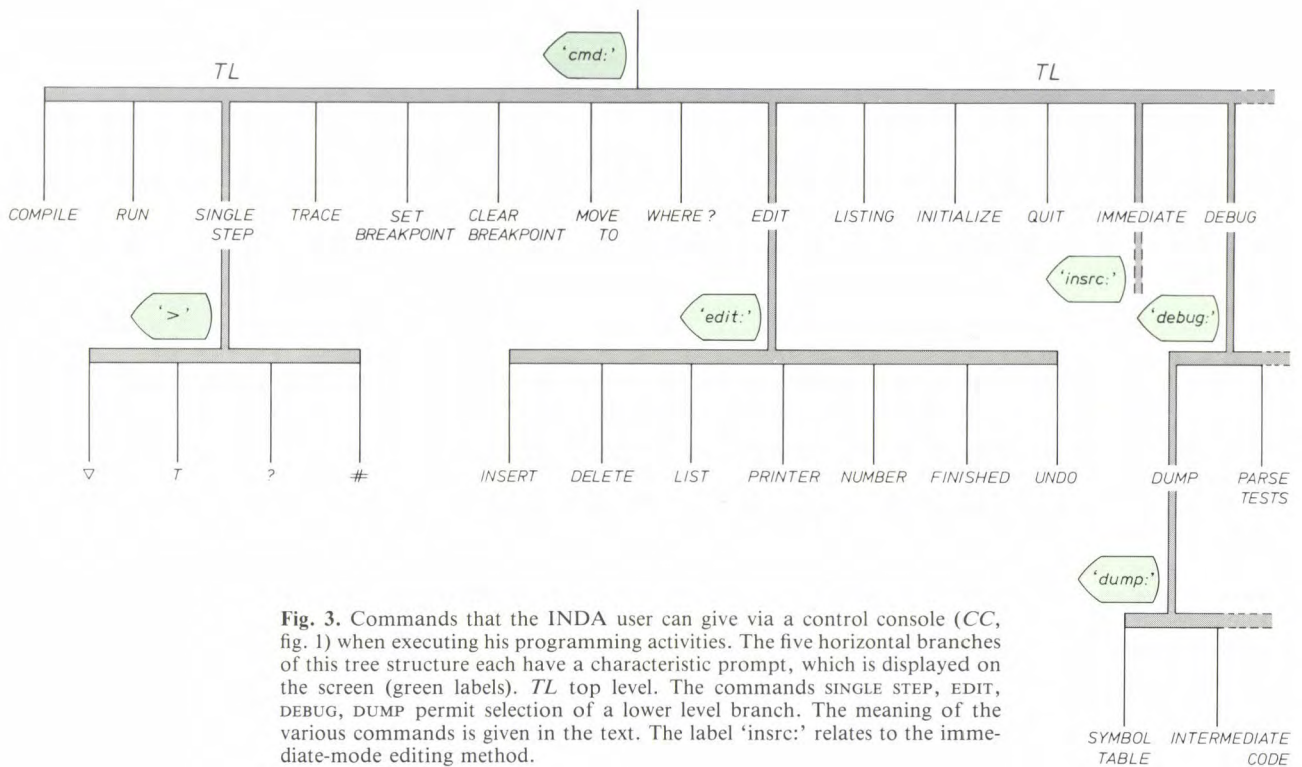


Fig. 3. Commands that the INDA user can give via a control console (CC, fig. 1) when executing his programming activities. The five horizontal branches of this tree structure each have a characteristic prompt, which is displayed on the screen (green labels). TL top level. The commands SINGLE STEP, EDIT, DEBUG, DUMP permit selection of a lower level branch. The meaning of the various commands is given in the text. The label 'insrc:' relates to the immediate-mode editing method.

gram by using the special facilities offered by the interpreter. The interpreter can perform a number of actions on its own, for example all computations. For other actions such as picture analysis and controlled movement of the robot arm the interpreter has to use the 'Library' module.

When the engineer notices an error in the performance of the application program, he stops the program and calls on the 'Editor' module to modify the program. This can be done by the 'patch' method, which is intended to provide instant correction of the smallest possible part of the program. Each correction made by the patch method, which represents a per-

The corrected program thus obtained may now be compiled by using the conventional compiler^{[2][3]} and can be permanently linked with the library procedures that are to be used. The resultant program is usually smaller — and needs less time to run — than the INDA version available to the engineer after the last correction was made. In the final stage a micro-computer may often be sufficient to run the program.

The command processor as interface

All interactions between the production engineer and INDA are routed through the command proces-

sor (fig. 1). This means that the engineer can communicate with all subordinate parts of INDA in the same way, which helps to make INDA an easy and 'natural' system to use. About 40 commands are available, distributed over the branches of a 'tree' (fig. 3). Each branch has its own characteristic 'prompt', an alerting signal that is displayed by INDA on the control-console screen as soon as activities in the branch of interest can continue. The engineer therefore always knows which branch has been activated and that INDA is ready for a subsequent command. If, for example, the screen shows the prompt 'cmd:' he knows that he can only type in a command from the topmost branch. If he does not remember what commands are available at the activated level, then he simply types in '?' and then all the relevant commands appear on the screen, i.e. the entire 'menu' at the level of interest. The processor needs only a few letters to distinguish the given command from all the other commands at the same level.

If, for example, when working at the top level, the engineer wishes to give the command `COMPILE`, then it is only necessary to type in 'co', since this is sufficient to avoid any confusion with other commands from the same level. After 'co' he simply taps the space-bar once and the processor itself will fill in the rest of the command. The processor then asks automatically for any further data such as the line number, and then executes the command.

The processor is relatively impervious to human errors. If, for example, the engineer types in an insufficient number of letters, then the processor automatically asks for more. If a subsequent command, for example `EDIT`, causes another branch of the tree to be selected, then this is immediately signalled to the engineer on the screen by the appearance of the new prompt signal. If he then wishes to return to the nearest higher level he can do this by the command `FINISHED` — at least from the branch with the prompt signal 'edit:' — or, more generally, by typing in '#'. During the execution of a `RUN` command an immediate return to the top level is possible by means of a special key, the `INTERRUPT` key. When this is pressed, the program, run by the interpreter, will stop at the end of the current line and wait for a fresh command.

The compiler and the language

The INDA system provides the production engineer with a flexible high-level programming language that is suitable for his primary task, the construction of new application programs, or for 'tuning' existing

programs. The language has a number of primitive functions that are available in the form of procedures. Some impression of the language can be gained from fig. 4, which is an excerpt from a program for the assembly of electric motors, the fictional example mentioned at the end of the introduction. It is the purpose of this excerpt of the robot program to pick up the rotor of an already half-assembled motor from a feeder station, convey the rotor to the assembly station, and mount it on the central shaft of the motor. The robot arm makes an essentially triangular movement between its parked position, the feeder station, the assembly station, and finally back to the parked position. Two sensors are used; the first (*M*) to check whether the rotor has been picked up successfully; the second (*C*) to supervise the mounting task. A proce-

L	Statement	Comment
50		% INSERT ROTOR ONTO ROD %
60	SPEED(90);	% HIGH SPEED %
70	GRIP(OPEN);	
80	MOVARM(XS,YS,ZS+4.5);	% APPROACH PICK-UP POSITION %
90	SPEED(30);	% SLOWLY WHEN CLOSE %
100	LSPICKUP;	% LABEL: BACK HERE TO RE-TRY %
110	MOVARM(XS,YS,ZS);	% GRIPPER APPROACHES ROTOR %
120	GRIP(CLOSE);	% GRAB ROTOR %
130	MOVARM(XS,YS,ZS+4.5);	% LIFT %
140	IF ROSENS(M) = 1 DO	% SUCCESSFUL? ASK SENSOR M %
150	ERROR("ROTOR FEEDER JAMMED");	
160	PAUSE();	
170	GRIP(OPEN);	
180	GOTO LSPICKUP;	% NOT TRY AGAIN %
190	END;	
200	SPEED(50);	% MEDIUM SPEED WHEN CARRYING %
210	MOVARM(XROD,YROD,ZROD);	% TO TIP OF ROD %
220	SPEED(30);	% SLOW INSERTION %
230	Z := ZROD;	
240	WHILE PICTOL(C) = 1 DO	% CAMERA MONITORS TOLERANCE %
250	Z := Z - 0.25;	
260	MOVARM(XROD,YROD,Z);	% ANOTHER SMALL INSERTION STEP %
270	REP;	
280	GRIP(OPEN);	% DONE IT %
290	MOVARM(XROD,YROD,ZROD);	% SLOWLY MOVE CLEAR %
300	SPEED(90);	
310	MOVARM(XP,YP,ZP);	% GO HOME FAST %

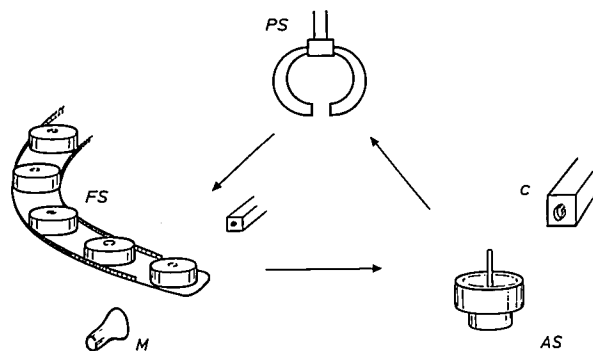


Fig. 4. An example of part of a pick-and-place program written with the aid of INDA for controlling a robot arm. *L* line number. Statements and comments can be displayed on the screen of the control console (see the title photograph). *PS* parked position of the robot arm. *FS* feeder station, for supplying the rotors. The rotors are picked up at a predefined rendezvous point, coordinates (*X*_S, *Y*_S, *Z*_S). *AS* assembly station, for mounting a rotor on the central shaft of a motor; the operation starts at the top end of the shaft; this point has the coordinates (*X*_{ROD}, *Y*_{ROD}, *Z*_{ROD}). *M*, *C* sensors, for checking rotor pick-up and mounting, respectively.

[2] The conventional RTL/2 compiler can be used successfully for the compiling.

[3] J. G. P. Barnes, RTL/2 design and philosophy, Heyden, London 1976, and BS 5904, Specification for computer programming language RTL/2, British Standards Institution, London 1980.

cedure from the library is required for checking the pick-up. This procedure, *RDSENS(M)*, signals '1' internally if the pick-up has been unsuccessful, which brings about a repetition of the operation. The second sensor, a video camera, employs the *PICTOL(C)* procedure from the library to check whether the mounting is being done at the correct depth. As long as a given tolerance is exceeded a '1' is signalled to the program, and this causes a continued increase in the depth. It will be obvious that the image analysis described here is a simplified and idealized representation of what happens; reality is more complicated. The program excerpt illustrates the effect of speed changes, arm movements, sensor input signals and error correction, albeit at a rather primitive level.

Certain points are clear. Firstly, in spite of the use of the high-level programming language, the description of even a simple task is difficult and can easily go wrong — and this makes it all the more important that the system offers a large number of facilities for locating errors. Secondly, the programming language has the control structures such as **if**, **while** and **goto** that are common in such languages and it cannot work without a library of procedures for the specific areas of application. Examples of such procedures are *MOVARM*, *GRIP*, *SPEED*, *RDSENS* and *PICTOL*. A language that is intended for application programs requires such a combination of general control structures and specific application procedures. The given program fragment is almost self-contained and can relatively easily be turned into a separate procedure in which certain parameters, for example the depth of mounting, could be specified externally. Such a program fragment could be created within the framework of the preparation of a much more complicated application program. Once again we have the familiar situation where complex problems can best be tackled by breaking them up into separate and smaller problems.

The language used in INDA is in fact a subset of the RTL/2 programming language^[9] (implementing the whole language would have created few essential extra facilities and would have taken up much more storage area and reduced the processing speed); the subset selected can evaluate algebraic expressions, execute conditional statements (**if ... then ... elseif ... then ... else ... end**), define looping activities (**while ... do ... rep**), and perform **goto** instructions. The useful data may be integers, real numbers, strings of characters, or tables of these.

The library procedures are written in the standard form of the RTL/2 language so that any piece of INDA data can be safely used as a parameter in these procedures. The compiler (*CPL*, fig. 1), which translates the program, produces three types of out-

put: symbol tables that define the data, interpretable machine code that defines the logic, and a record of the program text (the 'work file') for later processing by the editor. The compiler works incrementally, which means that it compiles each statement separately and in sequence. (This is not true for **if** and **while** statements, which are exceptions because of their composite nature.) The internal structure and layout of the compiler is conventional, with a lexical scanner, a parser, a code generator, and a symbol-table handler. In the listing produced, each line is numbered (fig. 4). These numbers, which rise in steps of 10, fulfil a function later in the process (figs 1 and 2) in the debugging by the interpreter and any possible modifications by the editor.

Running a program and error finding

The interpreter (*I*, fig. 1), which executes the programs while the production engineer employs its facilities for error finding, can be looked upon as a computer simulated in software. The 'machine language' of this simulated machine consists of *symbol tables* and the interpreter *code*, both of which result from the compilation of the programs. Because of this arrangement the original program text is followed much more closely during error finding than in the usual situations where the intermediate step of this simulated machine code with its symbol tables is not included.

About half the top-level commands (fig. 3) relate to the interpreter. The simplest command is **RUN**. This causes the program to run until either it ends of its own accord or the engineer presses the **INTERRUPT** key. In the case of such a stop the command **WHERE** enables the number of the next line to be executed to appear on the screen. The command **MOVE TO** moves the point where the program can start again to a line that the engineer can determine himself by selecting and typing in a number. The command **SET BREAKPOINT** has an effect similar to a manual interrupt; in this case, however, the program does not stop running at the current line (as in the case of **INTERRUPT**), but at a line specified by the engineer. This offers him the facility of being able to run any parts of the program separately. The command **CLEAR BREAKPOINT** is used to remove a breakpoint. The command **TRACE** causes the program to run completely, just like **RUN**, but with the additional factor that before the interpreter starts on a new statement it displays the line number of this statement on the screen. **SINGLE STEP** enables the program to be executed one line at a time under the engineer's supervision. It is indicated by the prompt signal '>'; in this situation a small separate menu of com-

mands is available. Each time the engineer taps the space-bar '▽' another line is interpreted; if, on the other hand, he taps 'T' then the line number always appears first on the screen. Thus the two facilities '▽' and 'T' are the single-step equivalents of RUN and TRACE, respectively. The command '?' corresponds here to WHERE.

The interpreter has other facilities that are useful for debugging. It can for example give a signal whenever the program is executing operations on a given variable. Some other facilities are dealt with in a later section, in the discussion on the command IMMEDIATE for transient corrections and modifications.

Modifying suspended programs

The editor is used to modify and edit a program (*E*, fig. 1). The patching method mentioned earlier permits local modifications, with the program continuing to run from the corrected point ('hot editing').

Correcting by the patching method is unusual, but is used with the language BASIC. The reason why we did not choose BASIC as the language for the INDA system is that it lacks certain language facilities (procedures in BASIC cannot have parameters, for example). Other reasons include the low speed of execution and the limited scope for input and output with BASIC.

The editing commands form a separate branch (fig. 3) just below the top-level commands. A very simple group of editing functions is sufficient here. Old programming lines can be deleted and new ones inserted. The compiler numbers the lines in steps of 10 allowing the insertion of up to 9 new lines between two successive lines. If line numbers become congested the command NUMBER can renumber all the lines, again in steps of 10. The command UNDO always offers the facility of abandoning any set of edits before they have been compiled. When the engineer has completed a given edit he issues the command FINISHED, to send the 'patched' fragment of the program to the compiler; once it has been established that it does not contain any compilation errors, the fragment of new text is written into the program at the original spot in place of the old text (which then disappears). The extent of the area to be patched must be given at the start of editing. It is wise of course to keep these patches as small as possible. Often an amendment of the program is best performed as a group of small distinct patches rather than one big one (which might require an unduly large storage area).

To make matters clearer, let us assume that the task described in fig. 4 has to be corrected. This takes place as shown in fig. 5. As can be seen the gripper carrying the rotor has moved a little too near to the shaft, causing a collision. Seeing this happen the engineer has quickly pressed the INTERRUPT key to stop the program. The program has stopped just before the execution of line 250 (as indicated by the arrow). The engineer will deduce that line 210 in the program is causing the difficulty. He then decides to increase the variable (*ZROD*), which contains the apparent length of the shaft, by 1.5 mm while the rotor is being mounted. He stores the old value of the apparent length in address *ZZ*. The inserted lines 205 and 206 illustrate this (the first patch). The original value of *ZROD* is re-inserted in line 275 as soon as this manoeuvre has been completed (the second patch). The text in the vicinity of the point of suspension has not been disturbed and execution of the program can be restarted at line 250.

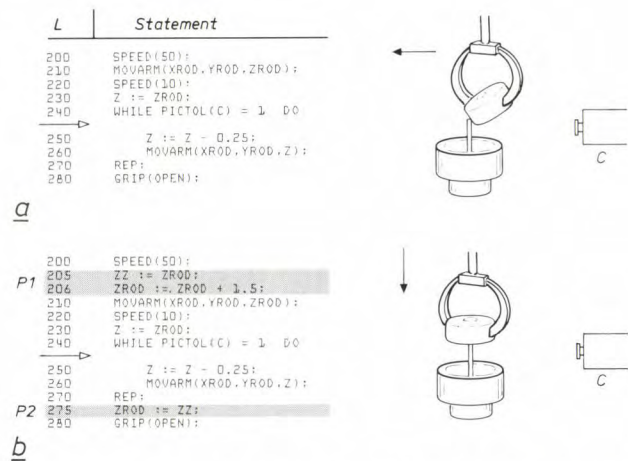


Fig. 5. An example of editing in the pick-and-place program of fig. 4 by the patching method. *a*) Original text of the program. The open arrow indicates where the operating engineer has stopped the program because he saw the rotor collide with the shaft. *b*) Edited program text. The inserted patches are *P*₁ (lines 205, 206) and *P*₂ (line 275). *L* line number. *C* video camera, checks the operation of mounting the rotor on the central motor shaft.

Transient program steps

INDA has a mode for transient modifications, a method of 'immediate' editing. In this mode a single statement or a small group of statements can be typed in. This is followed by the compilation and execution of the relevant program fragment — which is then immediately deleted, with full recovery of the original program text. Fig. 6 shows an example of this. The non-underlined characters are typed in by the engineer (the underlined characters form the computer output) when the computer is operating in the 'immediate'

mode. Such transient program fragments are used for individual examinations or data modifications, and also for directly adjusting the equipment. It is easy to see how this INDA facility could also be used for testing the modifications described in the previous section; the arm could in fact be moved in an immediate-mode statement, i.e. once only, to a safe starting position before the mounting manoeuvre is restarted.

Screen	Purpose
<u>CMD: IMMEDIATE</u> <u>INSRG: RWRITE(ZROD); 2.79 E1</u> <u>CMD:</u>	(data readout)
<u>CMD: IMMEDIATE</u> <u>INSRG: ZROD := 3.5 E1;</u> <u>INSRG: MOVARM(X.Y.ZROD+2.5);</u> <u>CMD:</u>	(data modification) (vertical displacement robot arm)

Fig. 6. Two examples of transient editing, in 'immediate mode'. The engineer types in the statements that are not underlined; those that are underlined represent computer output. *a*) The engineer wants to know the value and the address of *ZROD*. According to the computer the value is 2.79 and the storage address is *E1*. *b*) The engineer wants to change the value of *ZROD* and store the new value (3.5) in the original address. After this he wishes to shift the robot arm perpendicularly upwards a distance of 2.5 cm. Immediately after their execution the two commands (*a*) and (*b*) are automatically deleted from the program.

There are two ways, therefore, in which the immediate mode differs from the editor. The first is the non-permanent nature of the change made to the program; the second is the immediate execution of a statement after it has been typed in. (The editor does not start the execution until after all the statements to be used in a patch have been typed in and compiled.) The typing in of such immediate-mode program changes can be compared, in the matter of immediate execution, to commands given to the operating system.

The 'Library' module

The features of INDA that have been described can be used in many different fields. The 'Library' module (*L* in fig. 1) defines the areas of application. In the case discussed, the library contains the software tools for picture input and analysis, robot-arm control, and the exchange of information with the production engineer, all essential procedures for the particular application. An example of communication between the application program and the engineer is found in program line 150 in fig. 4 where the *ERROR* procedure displays the statement 'ROTOR FEEDER JAMMED' on the screen.

Reconstructing INDA with a suitable library for a new area of application is rapid and simple. The more difficult part is writing new procedures for a library.

Generally speaking, this is a task that can only be done by experienced programmers and even they will usually need considerable time.

Limitations and potential

The application programs should not be more than a few hundred lines, mainly because the available storage capacity is limited. The speed at which programs can be executed in the INDA system is relatively low. The main reason for the low speed is the use of the overlay method for program storage. (In this method the same area of store is used to contain successively different parts of the program. This increases the storage capacity but slows down the program.) Because the programs frequently relate to moving mechanisms and their components, the low speed of execution is not so critical. The treatment of INDA given in the present article is only restricted to applications in modern production methods for the purpose of explanation. At a much more general level, the INDA system can offer assistance in developing software in all kinds of projects, especially those in which complex hardware is subject to direct control and supervision by software, a situation that may also be encountered in the medical world and in many laboratory investigations.

No attempt has been made to include parallel and distributed processing within INDA. Programming based on these new ideas is still under development^{[4][5]}.

The author wishes to thank Dr W. T. Park of SRI International, Menlo Park, California, U.S.A., for his collaboration in the development of the INDA system.

^[4] See for example J. L. W. Kessels and A. J. Martin, Parallel programs, this issue, p. 254.

^[5] See for example B. L. A. Waumans, Software aspects of the PHIDIAS system, this issue, p. 262.

Summary. The INDA software system enables production engineers with little experience of programming to prepare and use application programs for automated production or modify such programs, by means of a display screen and keyboard ('menu' of about 40 instructions distributed over 5 levels). The system, which is written in the language RTL/2, offers high-level statement and processing facilities. There are two ways of dealing with errors: program patching — which rules out the need to restart at the beginning with expensive re-initialization of the production-machinery cycle — and 'immediate-mode' editing — without permanent insertion of modified program steps. INDA is suitable for setting up or 'tuning' supervisory and control programs consisting of a maximum of a few hundred lines and a library of standard procedures (for example image analysis, robot dynamics, information exchange). The system can run on a minicomputer such as a Philips P856/7. The programming of a robot arm for a pick-and-place operation is discussed as an example.

A data-base management system for CAD and CAM

W. E. Fischer

The computer is playing an increasingly important role in the industrial process. The use of the computer in the design of all kinds of technical objects (electronic products, engineering workpieces, buildings, tools) is referred to as 'Computer Aided Design' or CAD; its use in production with the aid of numerically

enabled to manufacture the product. During the design process a 'model' of a future product is thus developed and documented in a form that people can read and with which they can continue to work. This process takes place in stages, each stage leading to a more detailed specification of the model.



Fig. 1. Designing on an interactive display. The designer touches the light pen against a particular element on the display; this tells the computer that he is interested in that element. The element may be a piece of graphical information or an instruction; the result is either that data is retrieved from the data base, or that particular parts of the program are called. The designer can now take part in a dialogue with the computer.

controlled machines is referred to as 'Computer Aided Manufacturing' or CAM. Nowadays efforts are made to involve the computer in the design and production process at the earliest possible stage.

The term 'design' will be used in this article in the sense of designing a new product and then describing it in technical documents (such as parts lists and engineering drawings) in such a way that others are

Dipl.-Ing. W. E. Fischer is with Philips GmbH Forschungslaboratorium Hamburg, Hamburg, West Germany.

If we want to carry out this process entirely or in part with the support of a computer, then instead of engineering drawings we shall have to use other representations of the models that are suitable for processing by computer programs. Such representations, which are stored in the memory of a computer, are called 'internal computer representations' of the objects. On the basis of such representations the various stages of the design process can be implemented with the aid of the computer. The computer

can also be brought in for implementing other processes, e.g. for designing the tools required for manufacturing the new product.

Products are not usually designed completely independently of each other, because it is desirable to use the same components or tools if possible. It therefore makes sense to store the computer representations of different technical objects in a single data base, which can then be used time and again for different CAD/CAM applications.

Designing with the aid of the computer thus means that the designer builds up the internal computer representation of, say, an engineering workpiece in the data base and then elaborates it step by step. The dialogue with the computer that this requires is best performed graphically via an interactive display (*fig. 1*), rather like conventional drawing on the drawing board.

A system that makes this possible is the CAD/CAM system developed at Philips, called PHILIKON

(*PHILips Integriertes KONstruktionssystem*), whose principles have been described in an earlier issue of this journal [1]. The system uses a single common data base for all its CAD/CAM applications; hence the name 'integrated CAD/CAM system'. The different applications of PHILIKON utilize their own special 'dedicated' areas of the data base. There is for example a workpiece area, which is used for designing the tools required or for displaying the appropriate engineering drawings on a CRT screen or a drawing machine.

The components of a system like PHILIKON are shown in *fig. 2*. The central component consists of the 'modelling functions', which are the operations used in a particular area of application to generate, change or extend internal representations of objects. These operations are specified by an application programmer in application programs and are available to the users of PHILIKON. The users can call the required operations in dialogues with an interactive display.

The application programs that define the modelling functions are independent of the specific peripheral equipment used for data input and output and for communication with the user. All the data input to the PHILIKON system is translated by the PHIGRA program into a standard form that is used at the device-independent interface level I_1 . In the reverse direction all outputs of graphic data from the modelling functions are independent of the output equipment.

The application programs are also independent of the manner in which the internal representations of the models are stored in the data base. In *fig. 2* we therefore also see a data-independent interface level I_2 . The gap between the modelling functions and the data base is bridged by the PHIDAS system, which will be discussed in detail in this article.

The amount of data in a CAD/CAM system soon assumes such proportions that a special system becomes necessary for managing it. Previously such systems have always been individually developed for each specific CAD/CAM system. This procedure means that all application programs are dependent on the method of data storage. Improvements in data-storage methods, e.g. to achieve greater efficiency, then require so much effort that they prove to be impractical. The same applies to the integration of modules to form a larger CAD/CAM system.

There is therefore much to be said for replacing existing specific solutions by a more general-purpose system of data-base management, so that the application programmer no longer has to bother about data-management problems. He can then devote all his

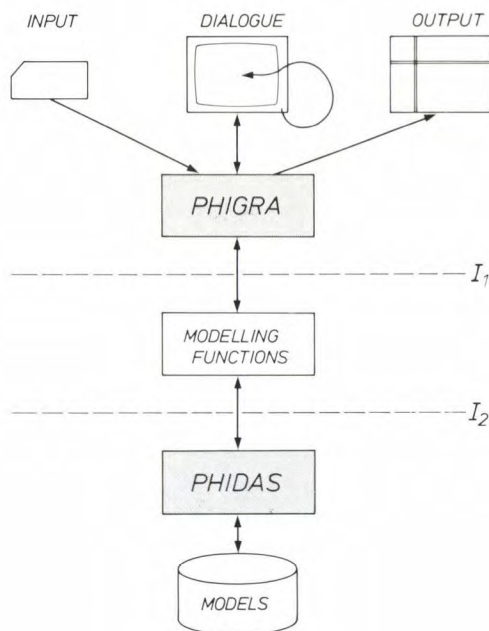


Fig. 2. The components and the data flow in PHILIKON — a system for computer-aided engineering design. The design process comprises the creation, modification and consultation of internal computer representations of engineering workpieces ('models'), stored in a data base. This is done by means of 'modelling functions', which are specified in application programs. Between the modelling functions and the input and output equipment there is the PHIGRA program module, which ensures that the modelling functions are independent of the peripherals employed. Access from the modelling functions to the models stored in the data base is obtained via the PHIDAS system. The modelling functions are therefore independent of the way in which the models are stored in the data base. The application programmer, who programs the modelling functions, is thus 'shielded' from the details of the communication with the user (above interface I_1) and from the details of the data storage (below interface I_2), which are the responsibility of the data-base administrator.

attention to the functional aspects of the operations on the model representation stored in the data base.

At Philips Forschungslaboratorium in Hamburg a system called PHIDAS (*PHilips DATenbank System*) has been developed for the special purpose of managing complex technical data bases, as used in CAD work. PHIDAS is smaller, faster and more flexible than the general-purpose systems previously available, which are mainly used in commercial applications.

In the following the use of the PHIDAS system for CAD/CAM applications will be described. We shall take a closer look at the organization of PHIDAS and at the structuring of the data stored in the data base, paying particular attention to the interface between the (data-independent) application programs and the data-base system. In conclusion the practical application of PHIDAS will be demonstrated in connection with the PHILIKON CAD/CAM system.

Data-base management by PHIDAS

Before dealing with the organization of a data base managed by PHIDAS, some relevant basic concepts will be defined.

First of all, of course, there is the *data*, i.e. the internal computer representations of the objects that play a part in the given application. A *data base* contains large amounts of data, which may be stored in many different formats. A data base is accessible to different users, who can retrieve different kinds of data from it.

A *data-base management system* is a software package that makes it possible to perform the necessary operations on a data base in a simpler manner. It enables the user, on the basis of a structural description of the data base, to add new data and to retrieve or modify stored data. A data-base management system, just like an operating system or a compiler for a programming language, is a software aid to the use of a computer.

When a data-base management system has been implemented on a particular computer, and a data base has been created that is accessible by means of the system, we have what we call a *data-base system*.

Management of the data comprises first and foremost regulating the storage of the data in the various storage devices and constantly attending to the distribution of the data in the main memory and the external memories. The system must also be able to make free memory space dynamically available for new components of an internal computer representation that appear on the display during the designer's work. All this makes fast data storage necessary and also the ability to make rapid changes in the data.

The PHIDAS data-base management system is designed to cope with these problems. It provides for central storage of the models. The data that each user needs is made available from the common data base. PHIDAS ensures that this data is transferred to each user's buffer memory and then returned to the disk system on which the actual data base is stored.

The PHIDAS system is addressed by instructions in the application programs. To make this possible the language FORTRAN has been extended by the addition of a number of instructions for data manipulation. This collection of instructions is referred to as the Data-Manipulation Language DML. We shall return to this presently.

It is of essential importance that the instructions used for the data management are independent of the specific application, i.e. that they form a common interface for all applications. Such an interface, which is independent of the specific implementation and also of the application, imposes certain requirements on the organization of the data-base system. For PHIDAS we have kept to the '3-SCHEMA concept' [2] recommended by the ANSI (the American National Standards Institute). The idea underlying this concept is the definition of the notion of 'data base' at three different levels:

- the user level,
- a 'conceptual' level,
- the data-storage and access level.

At each level the data base is described by its own SCHEMA, i.e. SUBSCHEMA, SCHEMA and INTERNAL SCHEMA. The division into separate levels enables the programs of the users to be kept largely immune to changes in the actual storage devices, in the storage structures employed and in the access paths to stored data.

Fig. 3 shows the arrangement of a PHIDAS data-base management system. The centre of the system is the conceptual SCHEMA, which specifies at a logical level *what* can be stored in the data base. The SCHEMA is described in DDL, a language specially devised for defining data (DDL stands for Data Description Language). This is done, once and for all, during the creation of the data-base system. The conceptual SCHEMA specifies the structure of the data base in so far as this follows from the properties of the kinds of objects that will be represented in it. Such a structural description does not yet specify the actual contents of the data base. This is not done until an application program is actually being run. With the

[1] P. Blume, Computer-aided design, Philips tech. Rev. 36, 162-175, 1976.

[2] D. A. Jardine (ed.), The ANSI/SPARC DBMS model (Proc. 2nd SHARE Working Conf. on Data base management systems, Montreal 1976), North-Holland, Amsterdam 1977.

aid of the instructions in the data-manipulation language DML, it is then possible to generate *records*, which are stored in the data base and are available to every user.

From outside of the data-base system the application programmer 'sees' only the part of the data base that he needs for his application. This particular part of the data base is defined by a SUBSCHEMA, also called an EXTERNAL SCHEMA (see fig. 3). This schema is also described in DDL.

The INTERNAL SCHEMA specifies the access paths and the mode of storage in the memories, the idea being to translate the logical structure embodied in the (conceptual) SCHEMA as efficiently as possible on to an actual physical storage structure. PHIDAS has a separate language for this description, called SSDL, for Storage-Structure Description Language. Changes in the INTERNAL SCHEMA have no effect on the other SCHEMAS. It is therefore said that the application programs are 'data-independent'.

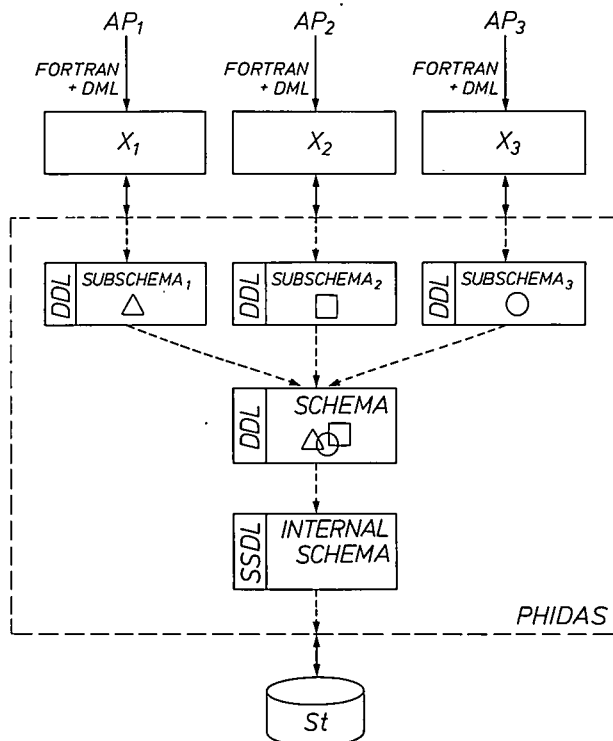


Fig. 3. Organization of a data-base management system in the 3-SCHEMA concept proposed by the American National Standards Institute (ANSI). In the (conceptual) SCHEMA the broad structure of all the stored data is described at a logical level in DDL (Data Description Language). The INTERNAL SCHEMA, using SSDL (Storage-Structure Description Language), specifies which of the possible structures have been selected for the actual physical storage of the data. Each application programmer AP_i 'sees' only the part of the data base that is described in 'his' SUBSCHEMA (a substructure of the SCHEMA). The application programmer has access to the data via an extended version of FORTRAN, which includes a Data Manipulation Language DML. PHIDAS conveys the required data from the disk store St to the buffer store X_i .

The SCHEMAS defined in the various description languages are compiled into a data catalogue that is connected to PHIDAS. From then on the data base can be filled via PHIDAS and the information stored in this way can be retrieved.

The definition of the various SCHEMAS is usually produced centrally by a specialist in basic software for CAD applications, the administrator of the data base; this does, however, have to be done in close consultation with specialists in the applications of interest.

Defining a SCHEMA for the PHILIKON CAD/CAM system

The internal computer representations of technical objects are essentially more complicated in structure than the data for numerical calculations, which can usually be represented by numbers (integer, real or complex).

To permit the structure of technical objects to be described in a SCHEMA that can be read by a system like PHIDAS, an appropriate formalism is required. The CODASYL Data Base Task Group^[3] has proposed a network concept that is suitable for this purpose. This concept, proposed as a standard for general data-base structures, is also the basis of the data structure used in PHILIKON for models and workpieces. This has been dealt with at some length earlier in this journal^[1]. Let us now consider the example of an engineering workpiece, to recapitulate the main aspects of the network structure, and show how the structure of a data base can be described in the DDL language proposed by CODASYL.

The CODASYL structure is based on the idea that a 'piece' of the world can be specified in the form of data by describing the objects present in that piece of world, together with their properties and interrelations. The objects are represented by *records*, their properties specified by means of *attributes*, and their interrelations represented by *CODASYL sets*. Objects of the same kind are represented by records that belong to the same *record type*. We shall now consider this in more detail.

The description of the SCHEMA of a workpiece model in the DDL language specifies for each type of object present a record type by specifying its name and its attributes. In fig. 4a this is done for the record type VERTEX. A record of this type represents a vertex of a solid body. In this case there are four attributes: each point has a number and three coordinates (X , Y and Z). Each record of the type VERTEX specifies a value for each of these four attributes; these values can be integers or real numbers.

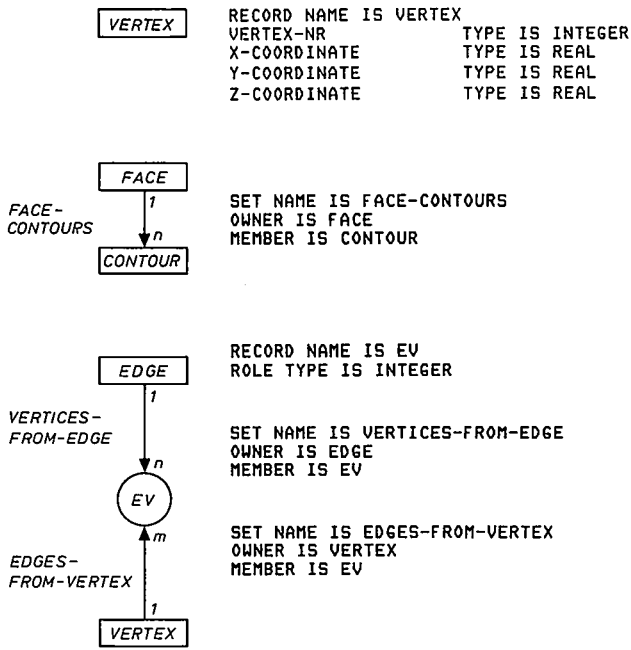


Fig. 4. Representation of the various elements for data structures as proposed by CODASYL, with the expressions for the definition of SCHEMAS in DDL. a) The record type VERTEX, with the coordinates X, Y, Z as attributes. b) The CODASYL set FACE-CONTOURS. To specify the role that each record plays in the set, the type of the OWNER records and the type of the MEMBER records have to be defined in each case. c) Relation record EV, for establishing a relation $n:m$ between record types. The CODASYL sets VERTICES-FROM-EDGE and EDGES-FROM-VERTEX have the relation records as MEMBERS. Also defined in the relation records is the attribute *ROLE*; its value indicates the role that the associated point plays in the relation with the particular line section.

Each record type is represented by a rectangle (see fig. 4a). The individual objects of a particular type, referred to as 'occurrences' of that type, are each represented by a separate record stored in the data base.

A CODASYL set serves for describing a relation between two record types; one of them is called OWNER, and the other MEMBER. A CODASYL set distinguishes a number of non-overlapping groups within the MEMBER records and assigns one OWNER-record to each group. It can thus be said that a CODASYL set represents a $1:n$ relationship between the OWNER records and the MEMBER records. A CODASYL set is represented by an arrow with a name beside it. The arrow starts at the type of

the OWNER records and its head points to the type of the MEMBER records (fig. 4b).

To describe a CODASYL set in DDL it is necessary to specify not only the name of the set but also the record type that is OWNER and the record type that is MEMBER. An occurrence of a set then gives the relation between an individual OWNER record and a number of MEMBER records.

A relation that exists between two kinds of objects is not necessarily a $1:n$ relationship. If we take, for example, a cube structure that is built up from vertices and edges, then the relation between vertices and adjoining edges is a $2:3$ relation: each edge is bounded by two vertices, while three edges meet at each vertex.

To describe such a relation by means of CODASYL sets we split it into two $1:n$ relations. This is done by introducing what are called *relation records*.

The description of an $n:m$ relation between two record types in DDL is performed by specifying the

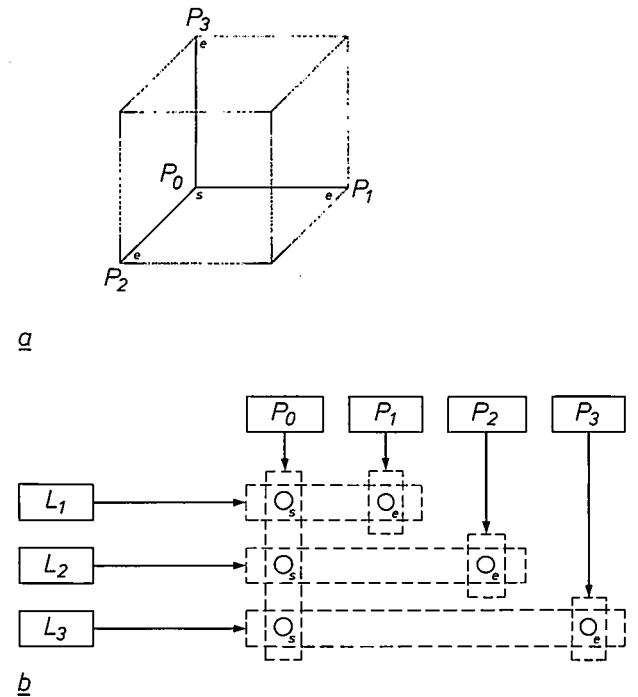


Fig. 5. Example of the use of relation records in a data base. a) The structure to be represented in the data base: three line sections, L_1 , L_2 and L_3 , meet at a vertex P_0 , and these in turn are bounded by P_1 , P_2 and P_3 respectively, and by P_0 . The points s and e are the start and the end of a line section. b) Since a number of line sections can meet at a point, while each line section is bounded by two points, we use relation records to indicate the relation between vertices and lines. For each 'vertex-line' combination that occurs there is then a relation record in the data base; such a relation record (the circles in the figure) is then a MEMBER in a link set (VERTICES-FROM-EDGE (horizontal dashed rectangles)), in which the associated EDGE record is OWNER, while at the same time the MEMBER in a link set is EDGES-FROM-VERTEX (vertical dashed rectangles), in which the associated VERTEX record is OWNER. This 'dual membership' permits the desired relation between VERTEX and EDGE records to be established.

[3] See: CODASYL Data Base Task Group, ACM report, New York, April 1971. CODASYL is an acronym formed from 'Conference on Data Systems Languages', an impartial organization formed in the interests of standardization in the area of data processing. Later modifications to the CODASYL concept are described in F. Manola, A review of the 1978 CODASYL Data Base Specifications, Proc. 4th Int. Conf. on Very large data bases, Berlin 1978, pp. 232-242. The PHOLAS system developed at Philips is based on the CODASYL proposal (see for example Introduction to PHOLAS, publication No. 5122 991 24961, Philips Electrológica, Apeldoorn, Nov. 1973). Unlike PHIDAS, PHOLAS was developed for large commercial data-base systems on main-frame computers.

name of the relation-record type (e.g. EV in fig. 4c) and defining two CODASYL sets in which that type occurs as a MEMBER.

If additional information has to be stored concerning the relation between the two record types, this can be done by means of attributes that have been defined in the relation record. If, for example, we assign a direction to the line sections that are represented by the records of the type EDGE, then we can indicate in each record of the type EV whether the associated point is a starting point or an end point of the associated line section. This can be done by defining an attribute *ROLE*, which is applicable to the relation records and can have two values: *s* for a starting point and *e* for an end point.

Fig. 5 shows the structure of the data that describes the relation between points and line sections at a corner of a cube: three edges with a common starting point (*s*) and a number of end points (*e*). There are six occurrences of the type relation record.

Fig. 6a shows by way of example an engineering workpiece for which an internal computer representation has been developed on the basis of the CODASYL structure. As illustrated, a workpiece (of the type PART) can be divided into a number of components that belong to other types, such as BASIC SHAPE, FORM ELEMENT, FACE, CONTOUR, EDGE, VERTEX, between which there are different relations. It is of course necessary to distinguish carefully between the relations $1:n$ and $m:n$ so as to be able to select the appropriate data structure relating to them, indicated by fig. 4b and 4c, respectively. Fig. 6b shows the complete diagram of the data structure for the workpiece of fig. 6a, in the form recommended by CODASYL.

As already mentioned (p. 246), for special applications it is possible to use the PHILIKON system to access specific parts of the data base; these are called 'archives'. This subdivision of a data base is based on the *area* concept.

By defining areas the data base can be divided into separate archives with different access rights, in such a way that the computer representations of different types of technical objects or of objects that are still in various stages of development are stored separately in different archives. When the SCHEMAS are specified the areas are defined at the same time. The user of the data base has access only to the contents of areas that are 'open'. The language DML is used for opening and closing the various areas. Access to these areas (a distinction is made between the right to 'read' and the right to 'modify') is obtained by means of the corresponding 'access key', which is known only by the authorized users and is specified in the SCHEMA.

Defining the INTERNAL SCHEMA

The definition of the conceptual SCHEMA depends very largely on a good understanding of the intended application, and is therefore the responsibility of the application programmer. He need not be concerned with the INTERNAL SCHEMA, however: this definition is usually left to the data-base administrator — who must of course bear in mind that the system must operate efficiently, with short access times to the data base. These depend on the 'access paths' available and the arrangement of the data in the actual memory device. Both are completely specified by the INTERNAL SCHEMA, using the storage-structure description language SSDL.

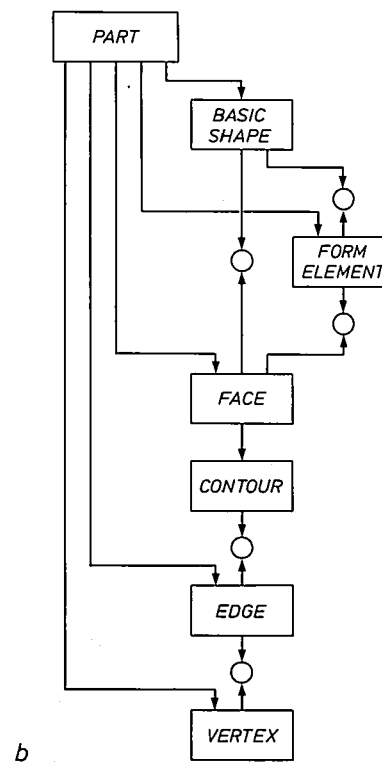
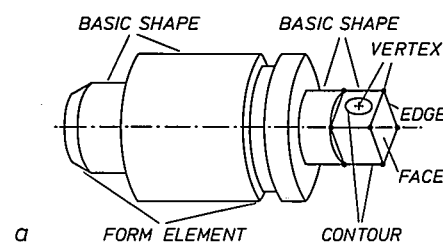


Fig. 6. a) Symbolic model of an engineering workpiece with parts that individually permit an increasingly refined subdivision into parts and types. Certain attributes can be assigned to each part, e.g. roughness to the surfaces and coordinates to the vertices. The separate components with their attributes and their interrelations form a completely defined model of the workpiece. b) The data structure of the workpiece in (a) in the form of a network model of the type recommended by CODASYL.

In choosing the actual physical arrangement of the data the administrator is guided by the applications. For example, he will arrange all the records of the internal computer representation of the same object in a single group, called a 'cluster'. When that cluster is called, all the data belonging to that particular object, which are in fact stored together, will then be quickly available and can be loaded into the main memory for executing a series of modelling steps on that object. This has the effect of reducing the number of times the external memories have to be accessed, thus saving much time. This clustering of data is rather different from the common practice in data bases for commercial data. In such applications it is useful to have all the data belonging to individual record types stored in a single file, for example all the orders in the order file, all the suppliers in the supplier file, and so on. In CAD applications such a structure would mean that all the line sections would be stored in a line file, etc. This method of grouping would slow down the system intolerably.

A 'cluster' in PHIDAS is a 'piece' of data structure with a single entry record — a record type designated as such — and an associated set of records. The records from that set all carry a mark that indicates their dependence on the entry record. When PHILIKON is used the entry record is always of the type PART. An occurrence of such a cluster is mapped 1:1 on to a virtual address space, called a segment. While a program is being run the size of a segment varies between 2 and 128 kilowords. These limits apply for computers with words of 16 bits. The segments are numbered, and a consecutive row of segment numbers always corresponds to a single file. An area may contain a number of these files. The details of the distribution of the files among the areas are decided by the data-base administrator.

The interface between application program and data-base system

Design engineers using computer assistance often write their programs in FORTRAN. To gain access to the data-base system they can use DML, the data-manipulation language mentioned earlier. This language forms an interface between the system and the application program and is independent of the specific application and the data-storage equipment. The functions of DML can be called in the same way as ordinary FORTRAN subroutines.

[4] A discussion of different versions of the interface between CODASYL and FORTRAN is given by W. E. Fischer, The interface between CAD/CAM-software and CODASYL-data base management systems, in: Eurographics 79, Proc. Int. Conf. and Exhib., Bologna 1979, pp. 178-189.

[5] CODASYL-FORTRAN-DMLC: CODASYL-FORTRAN data base facility, J. Devel. (Specification Board, Supply and Services, Canada), January 1977.

The committee concerned with the standardization of the FORTRAN-DML interface has put forward some proposals for functions of this kind. Some of these functions are used in our system; the most important ones are listed in *Table I*.

The interface we use [4] differs from the CODASYL proposals [5] in one essential point. This difference arises because of a special requirement that is imposed

Table I. Functions in DML (Data Manipulation Language) relating to the PHIDAS system for the management of an engineering data base.

* Call for a subschema and specification of areas, CODASYL types and record types within the SUBSCHEMA:

CALL INVS .. (<name of subschema>, <user key>)

CALL AREA .. (<list of parameters>)

CALL STYP .. (<as above>)

CALL RTYP .. (<as above>)

* Issue of access rights; opening and closing of areas:

CALL PRIVA .. (<name of subschema>, <name of area>, <access key>)

CALL OPEN .. (<name of subschema>, <name of area>)

CALL CLOSE .. (<as above>)

* Storage of records, and integration in a CODASYL set:

CALL STORE (<name of subschema>, <reference record>, <name of set type>, <name of record type>, <name of record>)

CALL INSERT (<name of subschema>, <reference record>, <name of set type>, <name of record>)

* Retrieval and reading of records, proceeding from reference records or from the contents of records:

CALL FIND .. (<name of subschema>, <reference record>, <name of set type>, <name of record>)

CALL GET (<name of subschema>, <name of record>, <name of record type>)

* Modification of contents of a record or of the relation to other records:

CALL MODIFY (<name of subschema>, <name of record>, <name of record type>)

CALL MOD2FY (<name of subschema>, <reference record>, <name of record type>, <name of record>)

* Deletion of records and sets:

CALL DELET .. (<name of subschema>, <name of record>, <name of record type>)

CALL REMOVE (<name of subschema>, <reference record>, <name of set type>)

In the variables <reference record> and <name of record> the associated temporary data-base key is stored or issued by the system.

Names with points (.) after them indicate that there are different versions. These are specified when the points are replaced by the corresponding characterizing letter combinations.

if a data-base system has to be accessed by means of graphical input devices such as a light pen or an interactive display: when the light pen ^[6] is pointed at the display the required data must be rapidly retrieved from the memory.

To meet this requirement, the PHIDAS system assigns unique *temporary data-base keys* to the relevant records at any moment. There is a fixed relation between these keys and the elements on the display, which are identifiable by means of the light pen while the image is present on the display. This relation is stated explicitly in a table, or it is an implicit part of the file containing the information that is used for continuously renewing the image on the display and for the drawing operations on the screen. When an element of the display is selected by means of the light pen, the graphic system assigns the associated data-base key to a FORTRAN variable, determined by the user, in the DML call that is used for putting the question to the data base.

The key establishes a correct relation between an element of the display and a record in the data base, provided that the area within which the associated model is managed is open and that the record is in fact stored. To obtain fast access to the record that represents the object selected with the light pen, the key referred to can for example be directly used as a parameter in an instruction GET (Table I). All the information connected via CODASYL sets to the record selected in this way can be directly obtained from the record. Each search process in the entire data base is thus reduced to a fast and purposeful navigation within a suitably selected subset of the data.

Applications of the PHIDAS system

The first prototype of PHIDAS has been implemented on a number of 16 bit minicomputers with the PHILIKON integrated design system. In this implementation all the modelling functions for the workpieces have the same standardized interface with respect to the data base: the data-manipulation language DML (I_2 in fig. 2). In a dialogue with the computer the designer can activate this function and thus describe new workpieces or form from parts already stored new combinations that can then serve for designing other workpieces. Stored workpieces can be retrieved. The designer can change details in a workpiece, and if he wishes he can remove a workpiece

from the data base. All this is possible by means of the DML operations made available by PHIDAS.

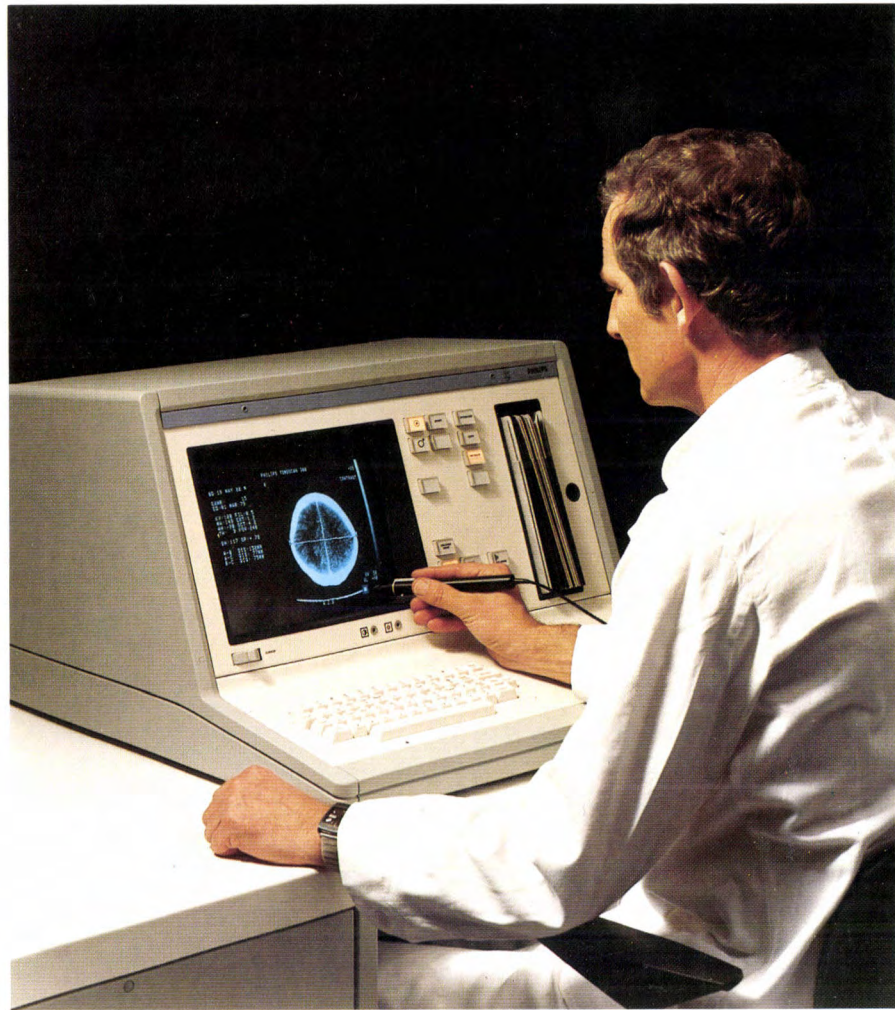
Initially PHIDAS was used in a number of single-user CAD/CAM systems. Recently it has also been implemented in situations where the data base is accessed from different displays that are used at the same time, or different interconnected minicomputers.

The particular problems associated with the use of interactive displays for CAD applications arise because every touch with a light pen triggers a large number of operations. These must then be performed 'immediately' in such a short time that the user is not inconvenienced by undue delays. PHIDAS is designed to cope with these problems. Even in subsequent extensions, it will therefore never become a competitor of existing large-scale data-base management systems, where the objectives are different. PHIDAS really comes into its own in engineering environments, in man-machine dialogues relating to complex sets of data where the response time of the system is the critical factor.

The work described in this article was carried out with financial support from the Ministry of Science and Technology of the Federal Republic of Germany, under contract No. HH PHI/105 and also from the combined German and Norwegian study group on Advanced Production Systems. The author alone is responsible for the contents of the article.

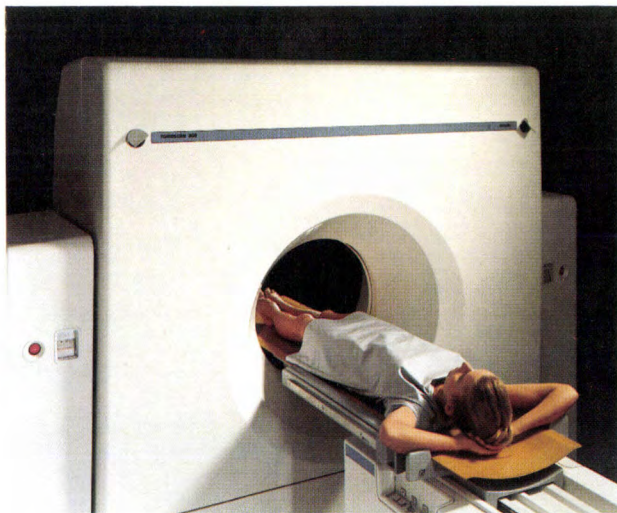
[6] Further details of the use of the light pen will be found in [1].

Summary. Data-base management for CAD and CAM dialogue systems imposes special demands on the application programmer, because of the complicated structure of the data and because the user, who works on his design via an interactive display, must have the required data available 'at the touch of a pen'. To facilitate the programmer's task, a team at Philips Forschungslaboratorium Hamburg have developed a software package for this purpose called PHIDAS (*PHilips DATenbank System*). PHIDAS is a data-base management system that is independent of the individual application and the computer employed. The system is based on a 3-SCHEMA concept as proposed by the American National Standards Institute (ANSI). The conceptual data structure used is based on the standards of the 'Conference on Data Systems Languages' (CODASYL), which enables PHIDAS to be employed with all data bases organized in accordance with these standards. The general usefulness of the PHIDAS system has been demonstrated by a number of implementations — one of them in combination with the PHILIKON integrated design system also developed at Philips — which have been implemented on various minicomputers.



X-ray pictures, taken from a large number of positions around the patient, can be used for reconstructing the X-ray attenuation in a thin slice or cross-section, so that an image of this cross-section is obtained. The reconstruction

is based on Fourier transforms and requires a great deal of computation, which is performed by a computer. The reconstructed image appears on a display screen, and the doctor using the TOMOSCAN 310 equipment can 'process' the image further and hence evaluate it. He can also use the light pen to indicate, on a more or less conventional X-ray picture, the cross-sections to be investigated. The patient is then moved automatically to the required position.



Parallel programs

J. L. W. Kessels and A. J. Martin

Advantages of parallelism

A computer program describes the way in which a complex information-processing task is broken down into elementary subtasks. The most familiar computer programs are *sequential programs*. A sequential program consists of a list of subtasks, which have to be performed one after another.

In general, different subsystems of the computer configuration are involved in the different subtasks: processing units perform arithmetical operations, input devices perform input operations, and output devices perform output operations. These subsystems are to some extent autonomous, i.e. they can perform their functions simultaneously and independently of one another for some period of time.

In the execution of a sequential program no use is made of the fact that different subsystems of the configuration can operate simultaneously: an operation is only started when the preceding operation has been completed. Sequential programs may therefore result in an inefficient use of the available equipment. This is true for conventional computer configurations, consisting of a single processor and various input and output devices. It is even more true for configurations with more than one processor.

In general there is a considerable amount of freedom in the order in which different subtasks could be executed. The execution of a task may therefore be described by a parallel program rather than a sequential one. Such a program specifies a number of *cooperating processes*; each of these processes is sequential, but the set of processes can be executed simultaneously (in *parallel*). For effective cooperation they must usually be synchronized at some points.

Ir J. L. W. Kessels and Ir A. J. Martin are with Philips Research Laboratories, Eindhoven.

Since in a parallel program it is explicitly indicated that some of the operations may be executed independently of one another, it is readily possible in this case to make use of the fact that different subsystems of the computer configuration can operate autonomously. The available equipment is then used more efficiently. For this reason parallel programming is used in operating systems and process-control systems.

Sometimes the choice of parallelism arises naturally from the environment in which the computer is used. If for example we consider a computer system with a number of terminals, the activities of the various users of these terminals are largely independent of one another. Much the same is true for the various sensors of a process-control system and for the subscribers of a telephone exchange. In all these cases there are various uncoordinated interactions between the computer system and the application environment. Operations that directly relate to these interactions have no fixed order in relation to one another; they are therefore best described as parallel processes.

Let us now look more closely at the nature of parallel programs. After we have shown how synchronization is achieved by operations on 'semaphores' we shall discuss the applications of synchronization. Finally, we shall look at two problems that can result from synchronization.

^[1] Semaphores were invented by E. W. Dijkstra. See: E. W. Dijkstra, Co-operating sequential processes, in: F. Genuys (ed.), Programming languages, pp. 43-112, Academic Press, London 1968; E. W. Dijkstra, The structure of the 'THE'-multiprogramming system, Comm. ACM 11, 341-346, 1968; E. W. Dijkstra, Hierarchical ordering of sequential processes, Acta Informatica 1, 115-138, 1971. Various possible definitions of semaphores are discussed in: A. J. Martin, An axiomatic definition of synchronization primitives, Acta Informatica 16, 219-235, 1981.

Synchronization

If a number of processes are executed on a computer configuration, these processes are generally not completely independent. At some points in the execution of a process, progress may depend on the progress of another process: the two processes have to be synchronized.

To achieve this the processes can perform special operations: the 'synchronization operations'. The most widely used synchronization operations are the P-operation and the V-operation, which operate on special variables called *semaphores* [1].

At the start of the program each of its semaphores receives an initial value — a non-negative integer. From then on the value of a semaphore can only be changed by the synchronization operations P and V, in the following way:

— the P-operation decreases the value of the semaphore by one, unless the semaphore has the value zero; if it has the value zero the operation is suspended until the moment when the semaphore attains a positive value;

— the V-operation increases the value of the semaphore by one. Therefore, if the V-operation is applied to a semaphore that other processes are waiting for in a P-operation, one of these processes will proceed.

The two following sections will show how the synchronization operations may be applied.

Mutual exclusion

A problem that frequently arises in parallel programming is the problem of preventing two program sections in different processes from being active simultaneously — i.e. mutual exclusion between these two program sections is required. This may be the case, for instance, if both sections require the same piece of hardware (e.g. an input/output device).

Synchronization to achieve mutual exclusion is easily programmed with semaphores. Consider for example the programs P1 and P2 of two cyclic processes, each consisting of two sections, which we shall call the 'non-critical section' and the 'critical section'. The two sections of P1 are called NCS1 and CS1 respectively; the two sections of P2 are called NCS2 and CS2 respectively. If the processes did not need to be synchronized, the two complete processes could be described by the program

```
P1: cycle NCS1; CS1 end
P2: cycle NCS2; CS2 end      (1)
```

If however the critical sections CS1 and CS2 must mutually exclude one another, the processes must be

synchronized in such a way that when P1 is in its critical section, P2 is guaranteed not to be in its critical section, and vice versa.

To arrange this, we include a semaphore s . The semaphore s has the value 1 if none of the processes is in its critical section, and the value zero if there is a process in its critical section. The processes must then modify the value of s in the correct way on entry to and exit from their critical sections. This can be achieved in the following way. At the start of the program each process is outside its critical section. We therefore assign to s the initial value 1. Before a process enters its critical section, it performs a P-operation on the semaphore s . If s has the value 1, the P-operation changes this value to 0, and the process continues. But if s has the value 0, the process will wait until s attains the value 1. At the end of a critical section s is made 1 again by means of a V-operation. We therefore have the following program:

```
Initial value:  $s = 1$ 
P1: cycle NCS1; P( $s$ ); CS1; V( $s$ ) end
P2: cycle NCS2; P( $s$ ); CS2; V( $s$ ) end      (2)
```

For the correctness of this program it is essential that the P-operation is one indivisible operation, which, if it establishes that the value of s is larger than zero, at the same time decreases the value of s , without another process being able to intervene.

If an attempt is made to achieve the same result by using separate test and assignment operations on 'ordinary' variables, problems arise, as can be seen in the program below. The intention is that whenever there is a process in its critical section, s has the value 0. A process may therefore only enter a critical section if $s = 1$. (Explanatory text is given between the curly brackets.)

```
P1: cycle
    NCS1;
    {the non-critical section}
    L1: if  $s = 0$  then goto L1;
    {as long as  $s = 0$ , the process is waiting}
     $s := 0$ ;
    {this statement is only reached if  $s$  was one;
    the fact that the process is now about to
    enter a critical section is indicated by making
     $s$  equal to 0}
    CS1;
    {the critical section}
     $s := 1$ 
    {this indicates that the process has left its
    critical section}
end
```

```

P2: cycle
    NCS2;
    L2: if  $s = 0$  then goto L2;
     $s := 0$ ;
    CS2;
     $s := 1$ 
end
(3)

```

This program does not guarantee correct synchronization. If, while s has the value 1, both processes execute the test ' $s = 0$ ' in the waiting cycle simultaneously or one immediately after the other, then *both* processes enter their critical section.

We see from this example that a parallel program that at first glance seems correct may in fact be incorrect. It is therefore of importance that methods exist that make it possible to design and prove such programs in a formal way.

The correctness of the program mentioned earlier, which used P- and V-operations, can for example be proved by an argument that makes use of an 'invariant' of the program. If n is the number of processes that are in a critical section, it can be shown that the value of $n + s$ remains invariant for parallel execution of P1 and P2. If s is decreased by 1 by a P-operation, then at the same moment n is increased by 1; in the same way, if a V-operation increases s by 1, the number of processes in the critical section decreases by 1.

Initially we have $n = 0$ and $s = 1$. The — constant — value of $n + s$ is therefore 1. Since s can never become negative: $s \geq 0$, and hence $n \leq 1$. This proves that the number of processes in the critical section is at most one.

This synchronization method, and its proof by the method of invariants, is also valid for more than two processes.

Partially ordered operations

A sequential program may be considered as a totally ordered set of operations. In many cases the correct execution of a task does not require such total ordering. In other words, in such cases the task may be specified by a partially ordered set of operations. The translation of such a specification into a program consisting of a number of parallel cooperating programs is in general not unique. If two operations have to be executed in a certain order, there are two different ways of describing this in a program: we can include the two operations in the correct order in a single process, but we can also include them in two different processes and prescribe the order by means of synchronization operations.

Let us assume for example that we wish to write a program that reads in a number of arguments via a punched-card reader, calculates the sine for all arguments, and outputs the results via a card punch in the corresponding order. It is clear that for each of the arguments the operations of reading the argument,

calculating the result and punching the result should be executed one after the other. This ordering of the operations in time is represented in *fig. 1* by the vertical arrows or 'edges'. The vertices in this network represent operations and the edges define an ordering relation between these operations: for each edge an operation represented by its initial vertex must be completed before the operation represented by its terminal vertex may be executed.

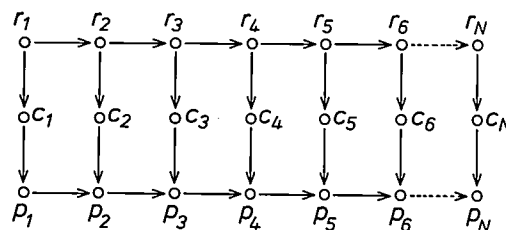


Fig. 1. Cards with function arguments are read, the function is computed for each argument and the results are punched (in the corresponding order). Different computer programs can be written for this task, since the order in which the read operations r_1, \dots, r_N , the calculation operations c_1, \dots, c_N and the punch operations p_1, \dots, p_N are to be performed is only partially specified. The partial ordering of the separate operations is represented by the arrows or 'edges' in the figure.

All the arguments are input via a single punched-card reader; this implies that the arguments are read one after another. This ordering is represented by the upper horizontal edges. Furthermore, we also wish all the results to be output via a single punch device, in the same order in which the arguments were read; this imposes an ordering on the punch operations that is represented by the lower horizontal edges. In this way the definition of the task establishes a partial ordering for the various read, arithmetic and output operations (fig. 1). The ordering is partial because it has not been determined, for example, in which sequence the operations c_1 and r_2 should be executed.

Sequential programs for partially ordered operations

A sequential program prescribes a complete ordering of its operations. If it is desired to write a sequential program for the partially ordered set of operations shown in fig. 1, there are various possibilities for choosing a total ordering within the specified partial ordering.

One obvious possibility is a program that reads an argument, performs the calculation, and punches the result, and then repeats these operations for all the following arguments. Such a program looks like this:

```

begin for i = 1 to N do begin READ(x);
                        y := SIN(x);
                        PUNCH(y)
end
end
    
```

(4)

The ordering prescribed by this program is shown in *fig. 2*. We see that the total ordering conforms to the partial ordering in *fig. 1*.

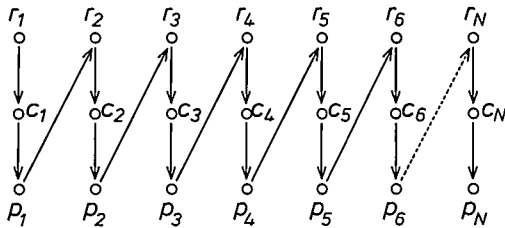


Fig. 2. The total ordering as prescribed by the sequential program (4). All arguments are dealt with completely one after another.

Another possibility is a program in which all the arguments are read first, all the calculations are executed next and finally all the results are punched. Such a program looks like this:

```

begin for i = 1 to N do READ(xi);
  for i = 1 to N do yi := SIN(xi);
  for i = 1 to N do PUNCH(yi)
end
    
```

(5)

The ordering prescribed by this program is shown in *fig. 3*. This total ordering also conforms to the partial ordering in *fig. 1*.

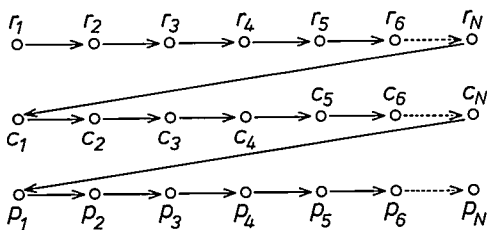


Fig. 3. The total ordering as prescribed by the sequential program (5). First all the arguments are read, then all the calculations are performed and finally all the results are punched.

two of the three subsystems that are involved in the execution of the task (the input device, the processing unit and the output device) are inactive, so that the execution of the task takes longer than it should.

Parallel programs for partially ordered operations

It is clear that a program that performs the task in the minimum time should allow some of the operations to be executed simultaneously. One way of achieving this is by subdividing the task into a number of parallel processes. This can be done in various ways; two of them are fairly obvious:

- by making use of a number of processes that each read, calculate and punch for a subset of the arguments,
- by using one process per type of operation, i.e. a read process, a calculation process and a punch process.

These programs are given below. In both programs we have used two vectors: x for the storage of arguments and y for the storage of results. If we do not introduce any further synchronization constraints, these vectors would have to be of length N (as is the case in programs). But since the execution time for the program does not become significantly shorter if more than three arguments are dealt with at the same time, we have used vectors of length three in both programs and the processes are synchronized accordingly. Both programs consist of non-terminating cyclic processes

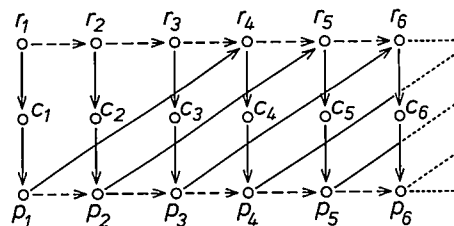


Fig. 4. The partial ordering as prescribed by the parallel program (6). The program consists of three parallel processes, each dealing with one argument at a time. Each process is represented by a sequence of solid edges; synchronization between the processes takes place as indicated by the dashed edges.

(it is assumed that the read operation is suspended when there is nothing more to read).

In the first solution the three processes proceed as indicated by the solid edges in *fig. 4*. Synchronization between the processes is indicated by the dashed edges. Each process has two semaphores, sr and sp ,

Both programs have the disadvantage that at any moment they allow only one operation to be executed, i.e. reading, or calculation or punching. From this it follows that during the execution of such a program

which indicate whether this process may execute the next read operation or punch operation. The program consists of the following processes (explanatory text between the curly brackets for each operation):

Initial values: $sr_1 = sp_1 = 1, sr_2 = sp_2 = sr_3 = sp_3 = 0$.

Process 1:

```

cycle P( $sr_1$ ); {wait until the read operation of process 1 is permitted}
    READ( $x_1$ ); {read argument}
    V( $sr_2$ ); {permit the read operation of process 2}
     $y_1 := \text{SIN}(x_1)$ ; {calculate the sine}
    P( $sp_1$ ); {wait until the punch operation of process 1 is permitted}
    PUNCH( $y_1$ ); {punch result}
    V( $sp_2$ ) {permit the punch operation of process 2}
end
    
```

Process 2:

```

cycle P( $sr_2$ ); READ( $x_2$ ); V( $sr_3$ );
     $y_2 := \text{SIN}(x_2)$ ;
    P( $sp_2$ ); PUNCH( $y_2$ ); V( $sp_3$ )
end
    
```

Process 3:

```

cycle P( $sr_3$ ); READ( $x_3$ ); V( $sr_1$ );
     $y_3 := \text{SIN}(x_3)$ ;
    P( $sp_3$ ); PUNCH( $y_3$ ); V( $sp_1$ )
end
    
```

(6)

If we choose the second approach, there are three semaphores: sr , which indicates whether the next read operation is permitted, sc , which indicates whether the next calculation operation is permitted, and sp , which indicates whether the next punch operation is permitted.

Initial values: $sr = 3, sc = sp = 0, i = j = k = 1$.

Read process:

```

cycle P( $sr$ ); READ( $x_i$ ); V( $sc$ );
    if  $i = 3$  then  $i := 1$  else  $i := i + 1$ 
end
    
```

Calculation process:

```

cycle P( $sc$ );  $y_j := \text{SIN}(x_j)$ ; V( $sp$ );
    if  $j = 3$  then  $j := 1$  else  $j := j + 1$ 
end
    
```

Punch process:

```

cycle P( $sp$ ); PUNCH( $y_k$ ); V( $sr$ );
    if  $k = 3$  then  $k := 1$  else  $k := k + 1$ 
end
    
```

(7)

The partial ordering specified by this program is shown in fig. 5. The synchronization of the read process by the punch process (the diagonal upward edges in fig. 5) is necessary because the read process and the calculation process can only store three intermediate results (x and y are vectors of length three). It is therefore necessary to prevent the read process from running too far ahead of the punch process.

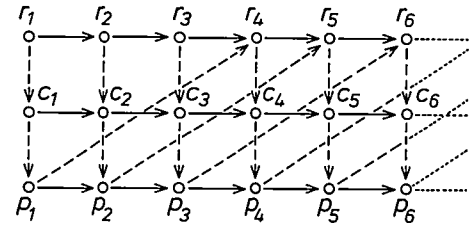


Fig. 5. The partial ordering as prescribed by the parallel program (7). This program consists of a read process, a calculation process and a punch process; each process is represented by a sequence of solid edges. The synchronization between the processes is indicated by the dashed edges.

Problems with parallel programming

In the foregoing we have indicated the attractive possibilities of parallel programming. However, it should be pointed out that this method of programming does have its own particular problems.

The first pitfall that we shall discuss is the effect known as 'deadlock', i.e. the occurrence of a situation in which the progress of two or more processes is permanently blocked because each is waiting for the other. Another problem is the prevention of 'individual starvation'; this will also be discussed shortly.

An example

As an example let us consider an airline reservation system. In such a system a central data base containing information about the number of vacant seats for each flight is consulted and modified from various terminals. The following two operations on the data base are of interest in our example:

VACANT SEATS(A) is a test that yields 'true' if there is a vacant seat on flight A, and 'false' if this is not the case.

RESERVE SEAT(A) reserves a seat on flight A. As a result of this the number of vacant seats for flight A is decreased by one.

'Overbooking' is not permitted; this means that the number of seats on a flight may not become negative. In reserving a seat a check is made first to see if there is still a vacant seat, and only if this is the case is the actual reservation made. A program for this purpose cannot however merely take the following form:

```

begin if VACANT SEATS(A)
  then RESERVE SEAT(A)
end (8)

```

If programs of this form were initiated simultaneously from two terminals and there is only one vacant seat, then these processes could both establish that there is still a vacant seat. Both could therefore make reservations and the number of vacant seats would become negative.

The solution is obvious: testing for the number of vacant seats and making a reservation should be a single indivisible operation, i.e. an operation that cannot be interfered with by other processes. We can achieve this by including both operations in a critical section of the process. We have already seen how we can use semaphores to ensure that different critical sections are not executed simultaneously (mutual exclusion). We therefore include a semaphore for every flight, e.g. s_A for flight A, with the initial value 1. If the reservation procedure now has the form of the program below, then it is ensured that there is never more than one process dealing with the seat reservations for a particular flight.

```

begin P( $s_A$ );
  if VACANT SEATS(A)
  then RESERVE SEAT(A);
  V( $s_A$ )
end (9)

```

Deadlock

Straightforward extensions of this solution can be used if we wish to reserve seats on two connecting flights A and B, with the condition that there are vacant seats on *both* flights. The interaction between different programs may lead to problems, however. As an example, imagine that the following two programs P1 and P2 are initiated simultaneously from two different terminals.

```

P1:
begin P( $s_A$ );
  if VACANT SEATS(A)
  then begin P( $s_B$ );
    if VACANT SEATS(B)
    then begin RESERVE SEAT(A);
      RESERVE SEAT(B)
    end;
    V( $s_B$ )
  end;
  V( $s_A$ )
end

```

```

P2:
begin P( $s_B$ );
  if VACANT SEATS(B)
  then begin P( $s_A$ );
    if VACANT SEATS(A)
    then begin RESERVE SEAT(B);
      RESERVE SEAT(A)
    end;
    V( $s_A$ )
  end;
  V( $s_B$ )
end (10)

```

In the parallel execution of P1 and P2 a 'deadlock' may occur, since it may happen that P1 has completed P(s_A) and is waiting at P(s_B), because P2 has completed P(s_B) and is held up at P(s_A). Both processes therefore wait for one another, so that access to the data for the flights A and B is permanently blocked.

Avoiding the occurrence of such deadlocks is one of the most difficult problems in writing parallel programs. In the case given above we have to ensure that in all processes the P-operations on different semaphores are performed in the same order.

Individual starvation

Another undesired effect that can occur in the parallel execution of programs is known as 'individual starvation'. Here again a process is permanently suspended in a synchronization operation, but the cause of this is not the same as in 'deadlock'.

Let us consider for example a program that consists of three parallel processes A, B and C, which each include a critical section; these critical sections are protected in the usual way by semaphores:

```

Initial value:  $s = 1$ 
A: cycle NCSA; P( $s$ ); CSA; V( $s$ ) end
B: cycle NCSB; P( $s$ ); CSB; V( $s$ ) end
C: cycle NCSC; P( $s$ ); CSC; V( $s$ ) end (11)

```

Let us assume that at a particular instant A is in its critical section and that B and C are held up in the operation P(s). If A leaves the critical section, s becomes equal to 1 and process B or C can therefore now proceed. Assume that B receives its turn and that A starts its P-operation before B has completed its critical section. At a particular instant B will leave its critical section and will increase the semaphore, so that A or C can continue. If A now receives its turn and B starts its P-operation before A performs its V-operation, the system is back in its initial state. If this course of events keeps on repeating itself, C will never receive a turn: C 'starves to death'.

The danger of 'individual starvation' puts a constraint on the way in which the P- and V-operations may be implemented. The implementation of the semaphores must ensure some degree of 'fairness', i.e. each process that is waiting in a P-operation may only be overtaken by another process a finite number of times.

Distributed programs

It has been assumed so far that parallel processes share a common memory, so that common variables, the semaphores, can be used for the synchronization. We call this *centralized* multiprocessing.

The term *distributed processing* is used when the parallel processes do not share a common memory. If distributed processes have to be synchronized, this cannot be done by making use of semaphores. Instead of P- and V-operations, communication operations are used both for the exchange of information between the different processes and for synchronization.

Various kinds of *communication primitives* have been proposed. One example is discussed in this issue in the article about the PHIDIAS distributed system^[2]; a more detailed discussion is to be found in the literature^[3].

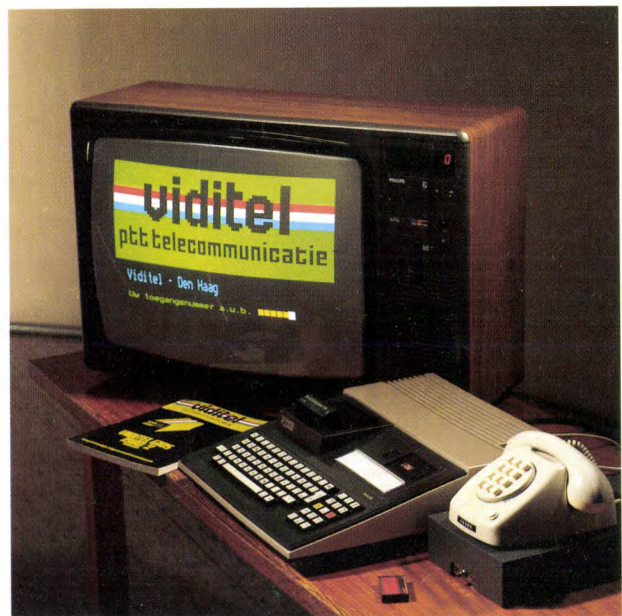
^[2] See B. L. A. Waumans, Software aspects of the PHIDIAS system, this issue, p. 262.

^[3] J. L. W. Kessels, The soma: a programming construct for distributed processing, IEEE Trans. SE-7, 502-509, 1981, and the article by A. J. Martín under [1].

Summary. Parallel programs have several advantages over the conventional sequential programs. They can be extremely simple and comprehensible, and offer possibilities for more efficient execution. A parallel program consists of various largely independent processes that have to be synchronized at some points. It may be necessary, for example, to ensure that certain parts of different processes are never executed at the same time. Synchronization is achieved by means of special operations on special variables (e.g. the P-operation and the V-operation on semaphores). In the design of parallel programs particular problems arise, such as the avoidance of situations in which certain processes remain permanently waiting for one another ('deadlock'), or in which some processes never receive a turn ('individual starvation'). The development of techniques for preventing such situations from arising is an important area of research.



The computer comes into the home. The user of the P 2000 personal computer can program it in BASIC for all kinds of administrative and educational activities; the same computer can also be used for various kinds of games — it can even play chess. The P 2000 makes use of the same set of symbols as the videotex systems in a number of countries ('Viditel' in the Netherlands) and can function as a videotex terminal.



Software aspects of the PHIDIAS system

B. L. A. Waumans

Introduction

Distributed computer systems

A conventional computer system consists of a central processing unit, or processor, one or more storage units or memories, and peripheral equipment for data input and output.

With the advent of VLSI (Very-Large-Scale Integration) the costs of actual information processing have decreased substantially in recent years. The costs of data communication, however, have remained relatively high. It is therefore often worth trying to keep data communication in the system to a minimum by performing calculations at the location where the required information is available or where the results are wanted.

This can be achieved with *distributed* computer systems, that is to say systems consisting of a number of processors each with its own memory. A system of this type may be considered as a set of interconnected computers that exchange messages with each other so that they can cooperate on a problem.

Other advantages of a distributed system are the following:

- Reliability. In a distributed system a component that fails can often be effectively isolated from the other components, which can then continue to operate normally. It may even be possible for the task of a defective component to be taken over by another component ('graceful degradation').
- Specialization. A distributed system can have components that are specialized in specific tasks. Such components may be cheaper and more efficient than general-purpose processors.
- Expandability. The size of a distributed system can be optimally chosen for the actual user environment. If a need should subsequently arise for a larger system, the system can be extended by adding more components.
- Speed. Independent operations can be performed simultaneously.

Ir B. L. A. Waumans is with Philips Research Laboratories, Eindhoven.

The PHIDIAS system was designed by a team consisting of Ir J. C. A. Boekhorst, G. Bos, Ir J. R. Brandsma, Ir M. A. Deurwaarder, Ir F. H. J. Feldbrugge, Ir P. G. Jansen, Ir J. L. W. Kessels, Ir W. J. H. M. Lippmann, Drs E. F. M. Steffens and the author.

Distributed programs

In many user environments (e.g. process control, telephone exchanges) the task of a computer system can be most clearly described by means of a collection of more or less independent processes that exchange messages. In such environments (in which of course distributed equipment will often be used) there is a need for a programming language that enables the programmer to write his programs in the form of more or less autonomous processes. It is of course advantageous if that language is a high-level programming language, so that the programs can be independent of the specific characteristics of the equipment associated with a particular computer system.

The PHIDIAS system

In view of the growing need of languages for distributed programming and the attractive features of distributed computer architectures, it was decided to set up a study project at Philips Research Laboratories to investigate the hardware and software aspects of distributed systems. In the course of this project an experimental system has been designed to which the name PHIDIAS has been given (for PHILIPS DIStributed ASynchronous System). In this article this system will be discussed with particular reference to the software aspects.

The PHIDIAS system can be considered at different levels of abstraction. At the highest level PHIDIAS can be described in terms of the facilities it offers to the user; at the lowest level it can be described in terms of the details of the hardware. Let us now briefly consider the levels with which we shall be concerned in this article.

1. The programming language. PHIDIAS is a general-purpose computer system; it does not fulfil a specific function, but enables a user to describe, by means of an application program, the execution of a task. The programming language that PHIDIAS offers for writing these programs has been designed so that it is particularly suitable for the distributed architecture of the system.
2. The operating system. Like any ordinary computer, PHIDIAS has an operating system, that is to say a

program — permanently running on the system — that controls the execution of the application programs. In view of the distributed organization, however, we have to make a distinction between the Global Operating System of PHIDIAS and the Local Operating Systems of the various components.

The Global Operating System (GLOS) determines which activities are to be carried out by which processors. It also ensures that the activities are appropriately redistributed among the processors if one processor fails.

In addition each processor runs a Local Operating System, LOS. This is a program that provides for the exchange of messages with other processors and enables its own processor to be used for more than one process (multiprocessing).

3. The architecture. At this level a description is given of the functions performed by the components from which the system is built up, and of the interconnections between these components. PHIDIAS is characterized by its distributed architecture; it is a network of intercommunicating processors each with its own private memory. The system does not have a memory that can be directly accessed by more than one processor. The processors are therefore called 'primes', a name formed from 'PRocessor with Individual MEmory'.

4. The components. PHIDIAS is built up from two types of components: the primes and the modules that ensure communication between the primes. Some primes serve for controlling specific system components (e.g. peripherals); others have the task of executing application programs or system programs. The communication modules are used to construct a network without central control.

We shall now look in somewhat more detail at the various levels of PHIDIAS, paying particular attention to the software levels. After dealing with the programming language, we shall consider the Global and the Local Operating Systems. The article will conclude with a brief discussion of the hardware (architecture, components).

The programming language

If optimum use is to be made of the facilities of a distributed computer system, the application programs to be executed should also have a 'distributed' character — that is to say, the program should specify a number of independent processes that interact solely through the exchange of 'messages' via communication channels. The design of a programming language for writing distributed programs was an essential complement to the design of the PHIDIAS hardware.

A programming language of this type must contain operations for the exchange of messages between the various processes. These operations are called 'communication primitives'. The choice of the communication primitives is often largely derived from the structure of the hardware employed. In PHIDIAS, on the other hand, we have tried to choose the communication primitives in such a way that the programmer using the language need not be aware of this structure. His programs could be run just as well on a conventional system with a single processor as on a system with a number of processors that have access to a common memory, or on a completely distributed system of processors each with its own private memory.

The programming language must also permit all kinds of structures for the communication paths between the processors. Pairs of processors can each have their own communication path, or different processors may use a common path. A path may or may not be capable of temporary information storage, and may or may not have to bridge a geographical distance.

The only assumptions that the PHIDIAS language makes about the communication network are that the network is capable of transferring messages to given destinations, and that a series of messages passed from one processor to another remain in the same sequence.

The soma

A program in the PHIDIAS language describes a number of somas (SOftware MAChines) ^[1]. A soma is a process that can inspect and change its own internal state by means of read and write operations on local variables, and which can communicate with other somas by sending and receiving messages. A soma is described by a sequential program, usually cyclic. For communication with the other somas the program uses operations on 'mailboxes'. Each soma has a number of mailboxes, in which it receives messages from other somas (see *fig. 1*). The number of mail-

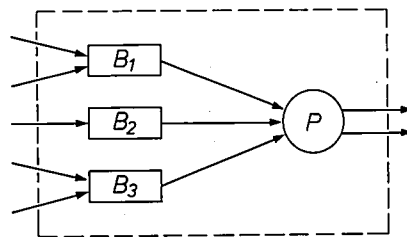


Fig. 1. A soma ('SOftware MAchine'), consisting of a process P with mailboxes B_1 , B_2 , B_3 . The process can send messages to other somas. Messages from other somas arrive in the mailboxes and are fetched from them by the process P .

[1] J. L. W. Kessels, The soma: a programming construct for distributed processing, IEEE Trans. SE-7, 502-509, 1981.

boxes in a soma, how the mailboxes are identified and how large they are must be specified by the programmer in the description of the soma.

A soma can collect messages from each of its own mailboxes and can send messages to any mailbox of any other soma. Whenever messages are fetched or sent the name of the mailbox involved is indicated. The programmer may reserve the various mailboxes of a soma for different categories of messages. In this way a distinction can be made between messages from different senders, between messages received in answer to the soma's own request and messages received on the initiative of another soma, and so on.

The mailboxes can also serve as buffers, since they may consist of different memory locations each of which can contain a message. If a message in a mailbox is not immediately fetched by the soma, it is still possible to send to the soma a limited number of subsequent messages. The buffer function is in this way completely separate from the transfer of messages. The communication path between the somas serves only for transfer, and not for the temporary storage of messages. Since a mailbox is a fixed part of a single soma, in this system it is not possible for different somas to use a buffer space collectively. This obviates a major cause of 'deadlock', since sending somas never have to wait until common buffer space is available.

If a soma tries to fetch a message from an empty mailbox, this action is postponed. The soma then waits until a message arrives in the mailbox, whereupon the action is then executed. A sending action, however, is *always* executed; if the mailbox is full, this results in an error message to the operating system, and as a result of this the execution of the program is terminated. The programmer must therefore ensure that messages are never sent to full mailboxes, and existing techniques enable him to see that this is not done.

Language constructions for specifying somas with their mailboxes and for describing the sending and fetching actions might in principle be added to any arbitrary sequential programming language. Ideally, an elegant language like PASCAL should be used for this purpose. Since the RTL/2 language was available on the hardware we were using, we decided to add the soma constructions to that language and to define the PHIDIAS language in this way.

Even when there is a good programming language available, writing correct distributed programs still poses its own particular problems. For example, the difficulties that are typical of parallel programs occur with distributed programs as well^[2]; in the synchronization between the different processes, patterns can

arise that prevent some processes from ever taking place ('individual starvation') or cause a state of deadlock because processes are constantly waiting for each other.

The programmer who designs distributed programs therefore needs techniques for making a careful investigation of the properties of such programs. It has been found that programs like those written in the PHIDIAS language can be effectively modelled with the aid of Petri nets^[3]. This formalism has been used, for example, to study the properties of the Global Operating System of PHIDIAS, which is also written in the PHIDIAS language.

The Global Operating System

The GLObal Operating System (GLOS) is responsible for controlling the peripherals and the operator's console and for allocating programs to primes.

Since GLOS has the typical aspects of a process-control program, it is readily programmed in the PHIDIAS language.

For some applications of distributed systems there is no need to load the system with different programs. The computer controlling a telephone exchange, for example, always runs the same program. In such cases the Global Operating System can be much simpler.

Control of peripheral equipment

Each peripheral should have its own prime on which the programs that control it are executed. Such programs are also written in the PHIDIAS language; by receiving and sending messages they communicate with the other programs that run on the system. In principle, primes that control peripherals are not used for other functions.

Partitions

One of the tasks of GLOS is allocating the programs to the primes. Each prime is loaded with one or more somas belonging to the same program. A soma is considered to be an 'atomic' element and therefore it may not be divided up among different primes. A collection of somas loaded on the same prime is called a *partition*.

In the implemented version of PHIDIAS the programmer himself decides on the partitioning of his program (see *fig. 2*).

The reasons for dividing programs in a given way into partitions are related to the architecture of the system, the size of the somas, the intensity of the interactions between the somas, the response times and the reliability required.

Loading the programs

The loading of programs on primes, starting and stopping the execution of the programs and the 'granting' of primes are functions of a part of GLOS called the 'Job-Supervision System', 'job' meaning the execution of a program. This JSS system itself consists of several parts: the Global Job Supervisor and a number of Local Job Supervisors (one for each prime on which partitions can be loaded). The Local Job Supervisor of a prime is the communication partner of the Global Job Supervisor for messages meant for this particular prime.

Instructions from the operator relating to the loading of programs, etc. are received by the Global Job Supervisor. There are four kinds of instructions, denoted by the key words GRANT PRIME, RECLAIM PRIME, START APPLICATION and STOP JOB.

The consequence of the instruction GRANT PRIME (p) is that the prime p is added to the set of available primes belonging to the system. (If jobs are queued up that could now be executed on this additional prime, they are *not* automatically loaded on to it.)

The effect of the instruction RECLAIM PRIME (p) is the opposite of the effect of GRANT PRIME (p). Its effect is to remove prime p from the set of available primes, unless p is in use at that moment for the execution of a job. In that case the Global Job Supervisor advises the operator accordingly. The operator can then stop the job and remove the partition (possibly allocating it to another prime). A subsequent RECLAIM PRIME (p) will now be successful, and has the effect that p is in fact removed.

The instruction START APPLICATION (a) causes the Global Job Supervisor to find out whether the primes and the peripherals required by the application program a are available. If they are not, the operator is informed accordingly. If the hardware required is in fact available a job is created and the hardware required is assigned to this particular job. A job identification is then communicated to the operator. (Dif-

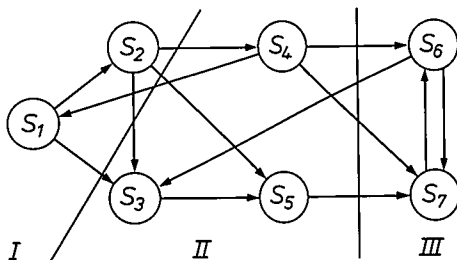


Fig. 2. A program consisting of seven somas, divided into three partitions. Partition I consists of the somas S_1 and S_2 ; partition II consists of the somas S_3 , S_4 and S_5 ; partition III consists of the somas S_6 and S_7 . Messages can be exchanged between somas from the same partition, and also between somas from different partitions.

ferent instances of the same application program can run simultaneously on the system, each instance being identified by its own 'job number'.)

The Global Job Supervisor creates a partition table, in which it records which prime has been granted to each of the partitions of the application program. During the execution of the program this table is required for communication between the somas, since it is obviously necessary to ensure that messages go to the correct primes.

In the next phase of the execution of the START APPLICATION instruction, the partitions are loaded on to the primes. This is done by considering each partition as a 'message' that has to be sent to the appropriate mailbox of the Local Job Supervisor in the appropriate prime. This 'mailbox' covers almost all of the memory of the prime. When all the partitions have been loaded, the program is executed.

The instruction STOP JOB (j) terminates job j . This action consists of a number of steps so that the termination proceeds in an orderly manner.

First, all the Local Job Supervisors are instructed to prevent their partitions from sending messages. All current write and read operations are completed, however. As soon as a partition terminates its activities, the Local Job Supervisor sends a message to the Global Job Supervisor. When the latter has received such messages from the Local Job Supervisors of all the primes involved in the execution of this job, the job can be stopped.

Error control

In a conventional computer system the failure of a component may cause the breakdown of the entire system. In a distributed system this can often be avoided. GLOS checks frequently to determine whether all the primes are still operating correctly. If it turns out that some are not, programs can still be correctly executed by transferring them to other available primes: this is known as 'reconfiguration' of the system.

The part of the operating system that performs these tasks and similar ones connected with error control is called the Error-Supervision System. The structure of this system is comparable with that of the Job-Supervision System. It embodies, for example, a Global Error Supervisor and a number of Local Error Supervisors, one for each prime.

Continuous tests are made to check that the messages are being transmitted correctly: the receiving

[2] See J. L. W. Kessels and A. J. Martin, Parallel programs, this issue, p. 254.

[3] W. Brauer (ed.), Net theory and applications (Lecture Notes in Computer Science 84), Springer, Berlin 1980.

prime 'echoes' every message received. If the Local Error Supervisor of the sending prime detects a mismatch between the message sent and the message received, it sends the message again. If the echo is still different from the original message, this is reported to the Global Error Supervisor, which then puts the malfunctioning prime out of action.

Primes are tested by sending certain messages at regular intervals to all primes. Each prime must then respond adequately within a specified period. If a prime fails to do so, it is treated as defective. The prime is then removed from the list of available primes.

Before a prime is removed, the job making use of this prime is terminated in an orderly manner. After removal of the prime the job may perhaps be continued or restarted on the processors that are still available.

If the Global Job Supervisor expects messages from a prime that has been found to be defective, the Error-Supervision System sends 'substitute messages' to answer that expectation. If this were not done, the failure of a prime would have repercussions on the entire system. Its effect would be to block the Global Job Supervisor, since this would be waiting for a message that would never arrive.

Possible extensions

In deciding on the design of the Global Operating System of PHIDIAS we often had a number of options available to us. In such cases we usually took the simplest solution. It is worth mentioning, however, that a somewhat more complicated mechanism might occasionally have led to an even more attractive system. Some examples follow below.

— In a PHIDIAS program the somas are grouped in a predetermined way to form partitions (see p. 264). In a version that exploited the full potential of the soma concept, the grouping of somas into partitions might be done automatically by GLOS; even the distribution of somas among the primes could be done dynamically. The result would be a system of greater flexibility.

— In PHIDIAS a program can only be executed if all the partitions have been loaded on to primes (see p. 265). It is also possible to consider a procedure in which partitions are loaded and removed during the execution of the program. Precautions will then have to be taken to ensure that somas never send messages to other somas that are not loaded.

— When a new prime is added to the set of available primes, waiting jobs are not put into execution automatically (see p. 265). It would be useful if this could be done.

The Local Operating Systems

On each prime a program runs that we call the Local Operating System (LOS) of this prime. This program LOS consists of two layers. The 'lower' layer of LOS, called LOS layer 0, is responsible for multi-processing on the prime, that is to say it allocates the available processor time among the different processes that run on the prime. (The Local Operating System of a prime that only controls peripherals can of course be simpler.) LOS layer 0 also makes available the P- and V-operations^[1] required for the synchronization of these processes.

The 'upper' layer of LOS, called LOS layer 1, is responsible for the sending and fetching of messages by the somas.

LOS layer 0

We shall now consider the way in which LOS layer 0 allocates the processor time between the various processes loaded on a prime. These processes have three different states, referred to as 'busy' (*B*), 'ready' (*R*) and 'waiting' (*W*) (see fig. 3).

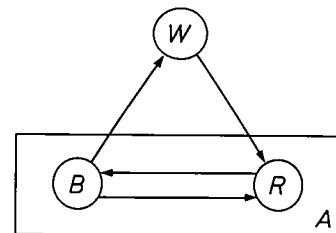


Fig. 3. The possible states of a process. *W* waiting. *B* busy. *R* ready. *A* active. The categories *B* and *R* together constitute the category *A*. The arrows indicate how a process can change from one state to another.

At any particular moment no more than one of the processes running on a prime is 'busy' — in other words only one process is being executed by the processor of the prime at any particular moment. A process that is not 'busy', but which is not delayed by a synchronization operation and could thus become 'busy' as soon as the processor is free, is referred to as 'ready'. Processes that are 'busy' or 'ready' are called 'active'. Not all processes are continuously 'active'. It can also happen that a process can only proceed after some external condition has been satisfied. Such a process is then referred to as 'waiting'.

Various 'ready' processes can have a different priority in the system. All 'ready' processes are recorded in order of their priority in a list called the 'Ready-to-Execute Queue' (REQ). Two operations can be performed on this list:

— the process with the highest priority can be retrieved from it, and

— a given process can be entered in the list at the location corresponding to its priority.

The 'busy' process must always have at least the same priority as the process with the highest priority in REQ. When a process that was previously waiting is added to REQ, the 'busy' process is interrupted and also added to REQ. The process that now has the highest priority is then selected from REQ and executed by the processor.

LOS layer 1

In LOS layer 1 the mailboxes and the associated send and fetch operations of the PHIDIAS language are implemented by means of operations that can be performed by LOS layer 0. There is a difference between the send and fetch operation, and this has a very important bearing on the implementation. The implementation of the fetch operation is completely independent of the properties of the communication paths between the somas, because the fetch operation takes place within a single soma. The implementation of the send operation, however, which results in the transfer of a message from one soma to another soma, does depend on the communication network and the input and output facilities of the primes.

A *fetch operation* is performed by a soma to retrieve information from one of its mailboxes. Each of these mailboxes consists of a number of locations. The state of the mailbox is described by a list of the filled locations, a list of the empty locations, and a semaphore^[2] that indicates the number of filled locations. A fetch operation on the mailbox starts with a P-operation on this semaphore. If the number of filled locations is zero, the fetch operation has to be suspended. After the P-operation has been successfully executed a message is copied from the first location in the list of filled locations; this location is then removed from the list of filled locations and entered into the list of empty locations.

The *send operation*, i.e. the transmission of a message from a soma on a prime to a soma on the same or another prime, might be done in various ways. The optimum choice depends to a great extent on the architecture of the system.

In PHIDIAS all primes can communicate with one another via a common network (see fig. 4). Each prime is connected to this network by an input port and an output port. A prime has one control unit available for sending messages and one for receiving messages. Steps are taken to ensure that only one message can be sent at a time and only one message received at a time. A single prime can however send a message at the same time as another message is being received.

It is not difficult to prevent different messages from being sent simultaneously by a given prime. The send actions used by the output port of a particular prime are all initiated by somas on that prime. The LOS of the prime thus has available all the information necessary to ensure the mutual exclusion of these send actions.

Mutual exclusion at the input port of a prime is a more complicated problem, because the arrival of messages at an input port is a consequence of independent decisions of somas on other primes.

To ensure that messages from different sources do not arrive simultaneously in the same prime, a prime wanting to send a message to another one must give advance notice of this to the receiving prime by sending a control message. Steps have been taken to ensure that these control messages are safely received. On the basis of the incoming control messages the receiving prime selects only one candidate sender at any given moment and this sender is then allowed to send the complete message.

The hardware of PHIDIAS

Since we are primarily concerned in this article with the software of the PHIDIAS system, we shall be very brief on the subject of its hardware.

Architecture

The general architecture of the PHIDIAS system is illustrated in fig. 4. The main components are the primes and the switching (or intercommunication) network. As already mentioned, there are two types

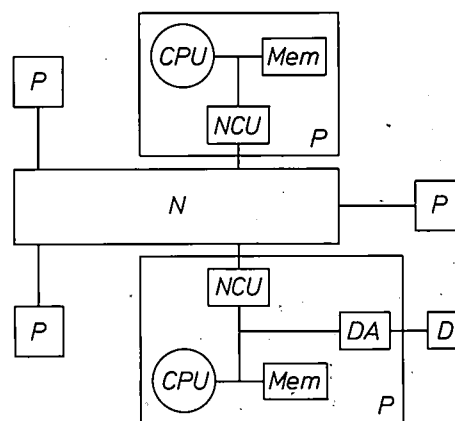


Fig. 4. The architecture of the PHIDIAS system consists of a number of primes P , which are interconnected by a common switching network N . A prime that serves for actual information processing consists simply of a central processing unit CPU and a memory Mem , which are connected to the PHIDIAS network N via a 'Network Control Unit' NCU . A prime that serves for controlling a particular hardware module D is connected to that hardware by means of a device adapter DA .

of primes. One type serves for controlling a specific hardware module. The other type serves for the actual information processing; the partitions of the application programs are executed on these primes.

The advantages of such an architecture were briefly mentioned at the beginning of this article.

Components

The components of the PHIDIAS system are various computers of the Philips P850 family: a P857 minicomputer, three P856 minicomputers, a P851 'single-board' computer and a FAST microprocessor, all with similar instruction sets. The direct-access memory of these computers has a storage capacity varying from 16 to 64 kbyte.

The *switching network* is of major importance in the system architecture. Measures are necessary to ensure that the network does not become blocked or that messages are not distorted by errors in the network. The network control units therefore run 'protocols' that keep a tally of the number of messages in the network; they also provide error-correction facilities by means of redundancy.

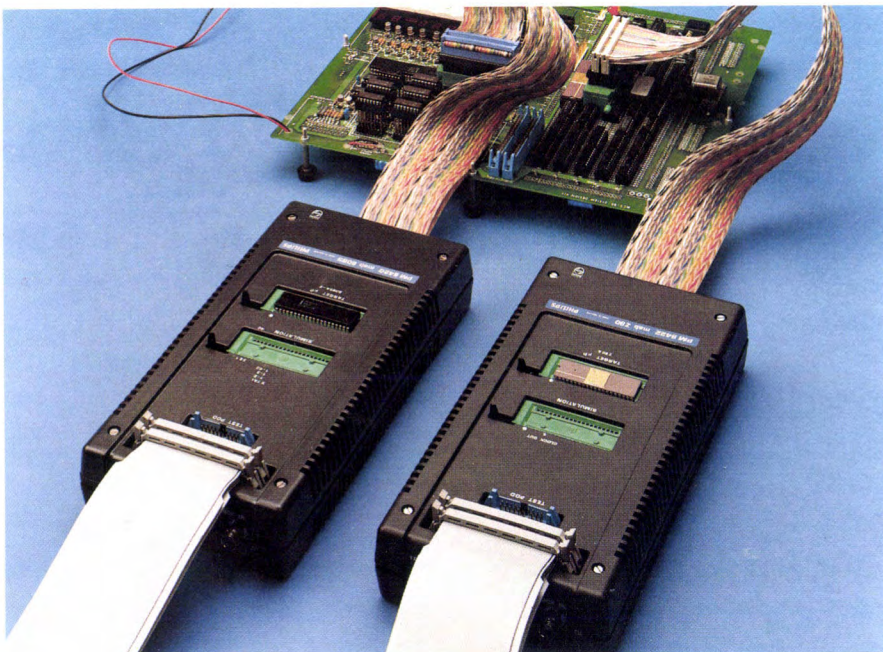
In the design of the switching network the emphasis has been placed on high information throughput rather than on short response times. The network is based on the DIMOND module^[4]. A network composed of DIMOND modules is completely decentralized and requires no central processing unit.

^[4] P. G. Jansen and J. L. W. Kessels, The DIMOND: a component for the modular construction of switching networks, IEEE Trans. C-29, 884-889, 1980.

Summary. The PHIDIAS distributed computer system is built up from 'primes' (processors with individual memory), which exchange messages by means of a common communication network. The system does not have a common memory. PHIDIAS executes programs that themselves have a distributed character. To enable programs to be written that are independent of the specific architecture of a particular computer system, an existing programming language was extended to include the facility for building up programs from independent processes (called 'somas') that exchange messages with one another. The operating system of PHIDIAS comprises a Global Operating System and a number of Local Operating Systems for the different primes. The Global Operating System can put defective primes out of action in the event of errors and redistribute the programs among the remaining primes. The Local Operating Systems ensure that a number of somas can run on one prime.



In developing the software for a microcomputer appropriate equipment is indispensable, such as the PM4421 Philips Microcomputer Development System shown here, which provides the software designer with aids such as a text editor, an assembler and compilers, which enable him to use high-level languages in the design procedure. Microprocessors made by various manufacturers can be included in the system by using special adapters connected by ribbon cables to the circuit being designed.



Distributed computations on arrays of processors

A. J. Martin

Introduction

The processor of a conventional computer is a sequential automatic device that performs its elementary operations one by one. A program intended for such a machine therefore specifies a fully ordered sequence of operations.

In many cases a complete specification of the order of the operations that constitute a computation is unnecessarily strict: many operations can often be performed in random order or simultaneously. If several computations can be carried out simultaneously, the computer time required for executing an algorithm can be reduced. To operate in this manner the computer must possess several processing units, and the program should indicate which operations can be carried out 'in parallel'.

Research on parallel machines and parallel programs is not new. 'Supercomputers' with several processors operating in parallel have been in existence for some time, and for some problems these machines are an order of magnitude faster than sequential machines.

The spectacular progress in VLSI (Very-Large-Scale Integration) has led to a rapid increase in interest in parallel computations. It is now possible to make computers that are as small and as cheap as transistors were in the sixties. The idea of constructing computer systems in which a large number of these microcomputers cooperate is therefore fairly attractive.

The processors have also become much faster. This means that the delay that occurs during the transfer of data is now a factor that must be taken into account in the design of a computer system. The original computers operated sequentially because 'wires were cheap and switching elements were expensive' [1]. It was more economical to transfer the data from all parts of the system for processing in a central processing unit. This is not necessarily true now: it is often more economical for a system to have processors at various locations to process data arriving or created there. Such a system is a *distributed* computer system.

Parallel programs

Programs to be run on a distributed computer system must indicate explicitly which parts of the program can be executed simultaneously and independently of one another.

Ir A. J. Martin is with Philips Research Laboratories, Eindhoven.

There are a number of problems associated with the design of parallel programs. In the first place it is sometimes difficult to identify the possibilities of parallelism in the solution of a given problem. We are used to thinking in terms of sequential algorithms, but the best sequential algorithm for a particular problem is not always the best starting point for the design of an optimum parallel algorithm.

A second problem arises from the large number of possibilities for interaction between the various sub-computations that make up a parallel computation; this makes it difficult to determine the correct behaviour of a parallel computation.

Finally, the third difficulty is that the efficiency of a parallel program is often strongly dependent on the structure of the machine that executes the program. This is probably the greatest stumbling block in the development of general-purpose parallel machines.

The success of the conventional sequential machines is largely attributable to the simplicity of the interface between hardware and software: in programming such a machine very little knowledge of its structure is required. The instructions that can be executed on a sequential machine, and their efficiency, are defined sufficiently well to enable the programmer to assess the efficiency of his algorithms — in both time and storage capacity. For parallel computations the situation is usually much more complicated.

The parallel activities to be carried out must be distributed over the interconnected processors that constitute the computer system. If the distribution is not specified by the program, but has to be determined dynamically, then the processors have to perform an additional task: so that the actual algorithm can be carried out it is necessary to keep a tally of the processors available and the processes still to be carried out. This can be a complicated task, taking up a great deal of processor time.

The interconnection structure of distributed systems

There are two essentially different ways of interconnecting the various processors of a distributed system; each has its own special problems.

In one approach there is a common communication channel between the processors — a switch, bus or memory system with all the processors connected to it.

Each processor can then communicate with any of the others. If the processors communicate with one another frequently (as they often do), the common communication channel can easily become a bottleneck. The article about PHIDIAS^[2] earlier in this issue gives an impression of the software aspects of computer systems with this kind of distributed architecture.

common memory, communication channel, or clock. Such systems are said to be *fully distributed*.

The system of fig. 1 is used for a particular type of computation: *recursive* computations. It will be shown in this article why recursive computations are ideally suited for execution by a distributed computer system with a network structure. The properties of different network structures will then be discussed.

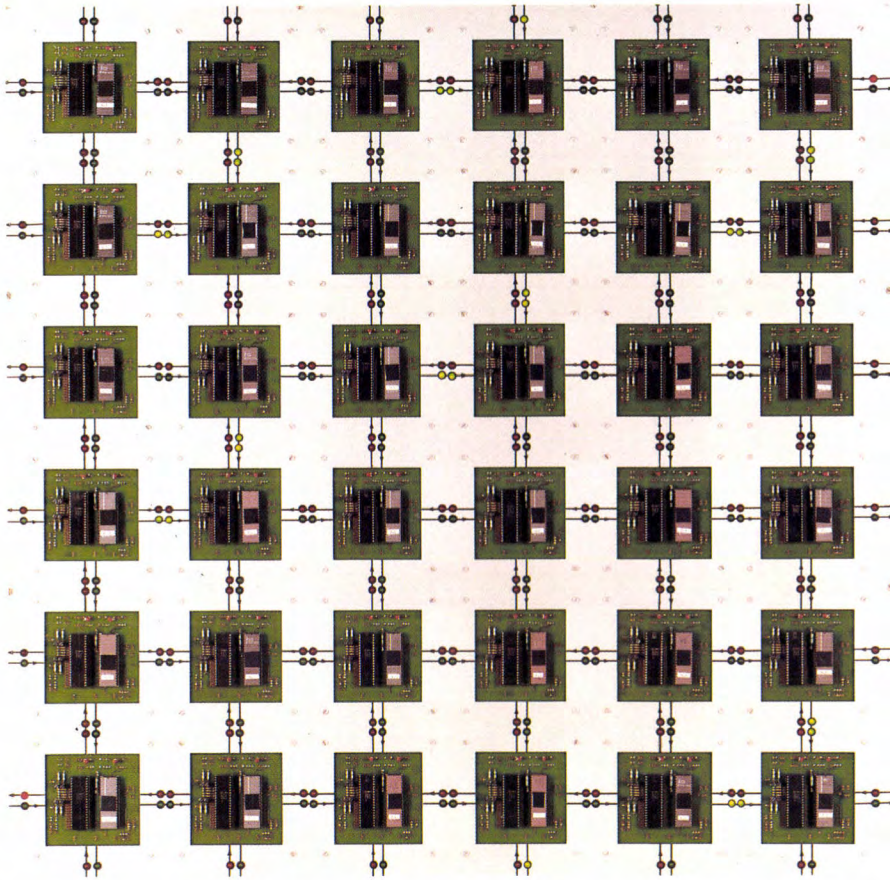


Fig. 1. Photograph of a distributed computer, a 'twisted torus' consisting of 36 'cells'. Each cell is built up from two INTEL chips (a processor with a 1 kbyte read-only memory, and a 256 byte random-access memory chip). The system is 'fully distributed'; there is no central clock. Neighbouring cells communicate with one another via a semi-duplex bit-serial protocol. The system was built at Philips Research Laboratories in early 1979 by G. Bos, M. A. Deurwaarder, W. J. Lippmann, G. A. Slavenburg and the author.

A completely different interconnection structure for the processors is produced if each communication channel only connects two processors. If the total number of connections is to remain within reasonable bounds, each processor can only be connected to a limited number of neighbours.

Fig. 1 shows an example of such a system. This system consists of an array of 'cells'. A cell is a processor with an associated memory unit. Each cell is directly connected only to its 'neighbours' in the array. The different cells in this system do not have a

Implementation of recursive computations

Recursive computations

A parallel program to be run on a distributed computer system does not specify one single sequential computation, but a whole set of subcomputations, which can mostly take place simultaneously and independently of one another. The subcomputations can be executed in many different orders. Because of the

^[1] I. E. Sutherland and C. A. Mead, *Microelectronics and computer science*, *Sci. Amer.* **237**, Sept. 1977, pp. 210-228.

^[2] B. L. A. Waumans, *Software aspects of the PHIDIAS system*, this issue, p. 262.

baffling complexity of this situation, it is important to apply programming techniques that will make it possible to assess the correctness of the program while considering all the possible interactions between subcomputations.

Considered from this perspective, recursive methods form an interesting category. A recursive computation creates subcomputations that are identical to the original computations. Because of this structure such computations are ideally suited to distributed systems.

Let us consider as an example the binomial coefficient $\binom{n}{k}$, defined by

$$\binom{n}{k} = n! / (n - k)! k! \text{ for } 0 \leq k \leq n, \text{ and}$$

$$\binom{n}{k} = 0 \text{ for other values of } k.$$

By making use of the properties

$$\binom{n}{0} = \binom{n}{n} = 1 \text{ for } n \geq 0 \text{ and}$$

$$\binom{n}{k} = \binom{n - 1}{k} + \binom{n - 1}{k - 1} \text{ for } n \geq 1,$$

we can set up a recursive procedure BC that calculates the value of $\binom{n}{k}$ for arbitrary n and k :

BC(n, k) =
 if $k < 0$ or $k > n \rightarrow 0$
 $k = 0$ or $k = n \rightarrow 1$
 $k > 0$ and $k < n \rightarrow \text{BC}(n - 1, k) + \text{BC}(n - 1, k - 1)$
 fi

It can be seen that the computation of BC(n, k) generally requires the computation of BC($n - 1, k$) and of BC($n - 1, k - 1$). These two computations are independent of one another and can therefore be performed simultaneously. They each use, in the same way, the results of two independent subcomputations,

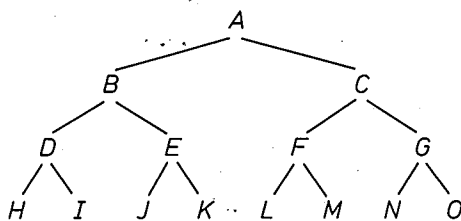
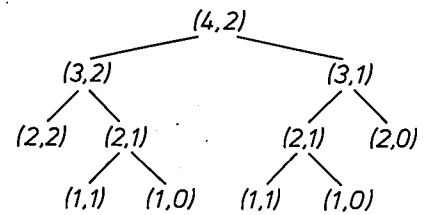


Fig. 2. A binary tree and its stepwise construction. The original computation A generates branches to the subcomputations B and C , B generates branches to D and E , C generates branches to F and G , and so on, until nodes are generated that represent 'primitive' subcomputations (in this case H, I, \dots, O).

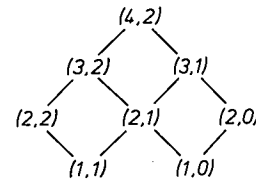
In the execution of the computation the tree shrinks again: computation I is executed, for example, passes its result to D , and disappears. If the same thing has happened to H , D can take its turn, and so on. Note that the order of these events is usually not fixed and that computations can occur simultaneously. 'Son' computations must always be executed before their 'father' computation, however.

and this process continues until we eventually arrive at subcomputations that can be considered 'elementary' (the cases for which the program specifies the result 0 or 1).

The entire computation has the structure of a binary tree (see fig. 2). The nodes of the tree represent the subcomputations, while a branch from node A to node B indicates that the result of subcomputation A depends on the result of subcomputation B . The subcomputation of BC(4,2), for example, can be represented by the tree:



The tree structure above is clearly not optimum. It would be preferable if the program led to the following network structure of subcomputations:



The generation of such network structures will not be discussed further here. In the rest of this article it is assumed that the computations have a tree structure.

As noted earlier, this article deals with the problems associated with the execution of recursive computations on distributed computers. Most attention will be given to recursive computations represented by binary trees; the situation in which each computation has only one subcomputation will be considered a special case.

Implementation networks

Let us now consider what would be a suitable structure for a distributed machine to be used for recursive computations. The kind of machine that comes to mind consists of a network of processors, each with its own memory unit. Each processor plus memory unit can be considered as a sequential machine, which we call a cell. Communication is limited: each cell is only directly connected to its neighbours. The system is fully distributed: there is no common memory, common communication channel or common clock.

The problem now is to implement a recursive computation on the distributed machine described above. The network structure of this machine, with cells as nodes and communication lines as branches, is called the 'implementation graph'. We can reformulate the problem: How is the computation tree mapped on to the 'implementation graph' — where the implementation graph is given a priori, and is fixed and finite, while the computation tree grows and then shrinks again, and we do not know beforehand how large it will become.

The number of nodes of the computation tree can thus become larger than the number of nodes of the implementation graph, so that it is necessary to assign more than one node of the calculation tree to one and the same node of the implementation graph. Various subcomputations are then performed by the same processor, and it is necessary to ensure that they are executed in the correct order. We should also consider whether in a given case the structure of the implementation graph is sufficiently well matched to the structure of the computation tree to ensure that the subcomputations are evenly distributed over the processors (optimum use of the distributed machine).

For our further considerations we start from an implementation graph in which one particular cell is in contact with the outside world; we call this cell the *root* of the graph. The initial node of the computation tree is then created from outside at the root of the implementation graph. The computation tree is mapped in recursive fashion on to the implementation graph, starting from the root cell: if node (n, k) , created in cell C , requires the subcomputations $(n-1, k)$ and $(n-1, k-1)$, these two nodes are created in two 'neighbours' of C . (If C only has one neighbouring cell both nodes are created in this cell. It is clear, however, that it is economical to use an implementation graph in which the number of 'sons' of each 'father'-cell is equal to the number of 'sons' of each node in the calculation tree — two in this case.)

In this way two nodes that communicate with one another (i.e. a 'father' and a 'son') are always mapped on to two cells that are directly connected. The communication between two cells then operates very simply: the father-cell passes data for processing to the son-cell, and the son-cell returns the result to the father-cell. This method therefore has the attractive feature that messages from one cell to another are never sent indirectly, via intermediate cells. Processors do not therefore have to waste time finding out how messages can reach their correct destinations.

Another attractive feature is that the subcomputations are distributed over the cells without any need for central control. Each cell calls on its own neigh-

bour for its subcomputations. The distribution pattern of the computations over the cells becomes more homogeneous as the structure of the implementation graph is more closely matched to the structure of the computation graph. We shall return to this later.

The danger of 'deadlocks' occurring in the parallel execution of computations has already been mentioned in this issue [3]. This problem must also be considered here, for if different subcomputations are assigned to a single cell, the order of execution must be made such that, if subcomputation A requires the result of subcomputation B , B can be performed before A . This can be arranged simply and efficiently [4].

Another problem is the 'housekeeping' of the messages that the processors send to their neighbours. There is no set upper limit to the number of messages present in the system at any moment; this number generally depends on the parameters of the computation being performed. Since it is therefore not possible to indicate the size of the buffers required, classical buffer methods cannot be used. A generalization of the stack mechanism can be used instead.

Tessellations

A first requirement in the choice of structure for an implementation graph is that an arbitrary binary tree can be mapped on to it, regardless of the size of the tree and the graph. It must also be possible to distribute the nodes of a tree over the cells of the graph in an optimum fashion.

Because of the first of these requirements we can say that we wish to 'simulate' an infinite network on a finite network. We can start with the question: If we could use an infinite network, what structure would we then choose? Networks that can be represented by regular and 'dense' patterns in a plane seem to have attractive properties. If we turn our attention to such networks, there is a choice of three regular planar tessellations: square, triangular and hexagonal. We decided to use the square tessellation, or 'grid', although the hexagonal one is also interesting.

An infinite grid is a network, consisting of nodes (i, j) , such that for $i \geq 0$ and $j \geq 0$, node (i, j) is connected with node $(i+1, j)$ and with node $(i, j+1)$.

It is fairly obvious how a binary tree can be mapped on to such a grid. The initial node of the tree is mapped on to cell $(0, 0)$. If a node of the tree is mapped on to cell (i, j) ; its right-hand son is mapped on to element $(i, j+1)$, and its left-hand son on to cell $(i+1, j)$. A difficulty with this method is that a congestion problem arises because the cells in the middle

[3] J. L. W. Kessels and A. J. Martin, Parallel programs, this issue, p. 254.

[4] A. J. Martin, A distributed implementation method for parallel programming, in: Information processing 80 (Proc. IFIP Congress 80, Tokyo/Melbourne 1980), pp. 309-314.

of the grid are assigned more nodes than the cells at the edges (see *fig. 3*.) The number of nodes at cell (i, j) is $(i + j)!/i!j!$.

In the remainder of this article we shall try to find out how we can best use finite means to approximate the behaviour of an infinite grid, or put another way, how we can best *simulate* an infinite grid on a finite network of elements.

The simulation of an infinite grid

The straight torus

We shall now consider several ways of simulating an infinite grid on a symmetrical grid consisting of M by M 'vertices'. (We shall reconsider this choice later.) The simplest way of producing an infinite grid structure with such a finite grid is by interconnecting the corresponding elements of the first and the last column and likewise the corresponding elements of the first and the last row (see *fig. 4*). This means that the vertex (x, y) is now always connected with the vertex $(x, (y + 1) \bmod M)$ and $((x + 1) \bmod M, y)$. The 'surface' thus produced has the topology of a straight torus (or 'anchor ring').

If we now consider an arbitrary vertex (i, j) of the infinite grid, and the vertex (x, y) of the finite torus on to which (i, j) is mapped, then from the above

$$i = x + k * M$$

and

$$j = y + l * M.$$

These relations describe the way in which the infinite grid is built up from 'tiles' of $M * M$ vertices: if a vertex in the infinite grid has the coordinates (i, j) , then it has the coordinates (x, y) in tile (k, l) . An example with $M = 4$, $k = 3$ and $l = 2$ is shown in *fig. 5*.

The congestion problem that was mentioned earlier can be solved for the infinite grid in the following way. If a cell is occupied by a node N of the computation tree, the cell accepts no other nodes until N (and the sons of N) have completed their activities. It is easy to prove that this cannot produce a deadlock on the infinite grid. But this solution cannot be directly applied to the torus without the danger of introducing a deadlock. Let us assume for example that a cell of the torus is occupied by node $N1$ and that a new node $N2$ is not accepted by the cell. If $N2$ happens to be a son or 'descendant' of $N1$, a deadlock situation has arisen. To avoid this, the tile it belongs to is quoted for each node of the computation tree. If a cell is occupied by a node $N1$, it may refuse node $N2$ if $N1$ and $N2$ belong to the same tile, but not otherwise. (A node and one of its descendants in the same tile are always mapped on to two different cells.)

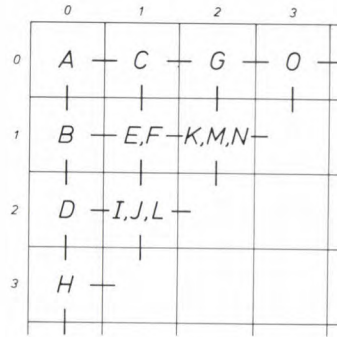


Fig. 3. The binary tree shown in *fig. 2*, mapped on to an infinite grid. The initial node of the tree is located at cell $(0,0)$, a left-hand son always goes one place to the right of its father. We see that the cells in the middle of the grid are assigned more nodes than those at the edge.

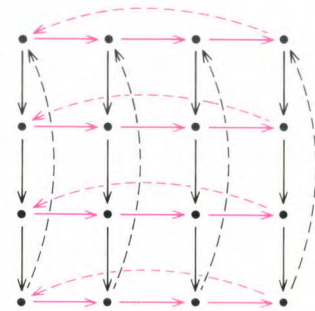


Fig. 4. Simulation of an infinite grid with a finite grid. The bottom element in each column is connected with the top one, and the element at the right-hand end of each row is connected with the element at the left-hand end. The surface obtained in this way, a straight *torus*, is unbounded in both the 'horizontal' direction (red arrows) and the 'vertical' direction (black arrows).

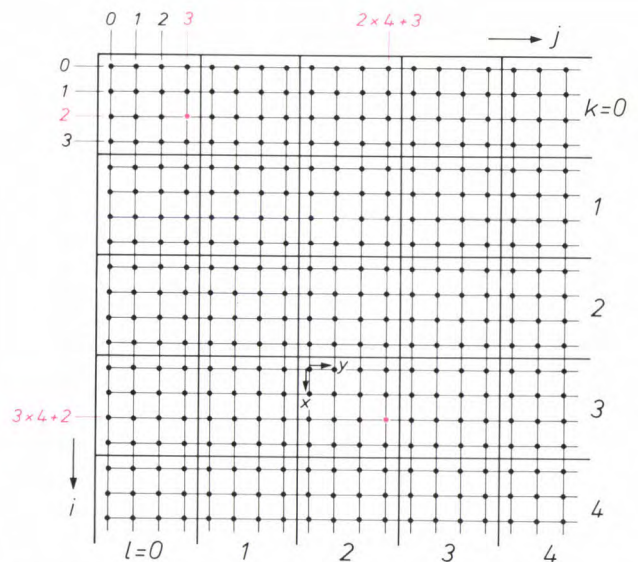


Fig. 5. Simulation of an infinite grid on a 'straight torus' of 4 by 4 elements. Every element in the infinite grid corresponds to an element in the torus. In an infinite grid we reach (say) the point $(3 * 4 + 2, 2 * 4 + 3)$ from the origin by taking a certain number of steps (downwards and to the right) from the origin $(0,0)$. By taking the same steps in the torus we reach the point $(2,3)$. The infinite grid can thus be thought of as built up from 'tiles' of 4 by 4 elements.

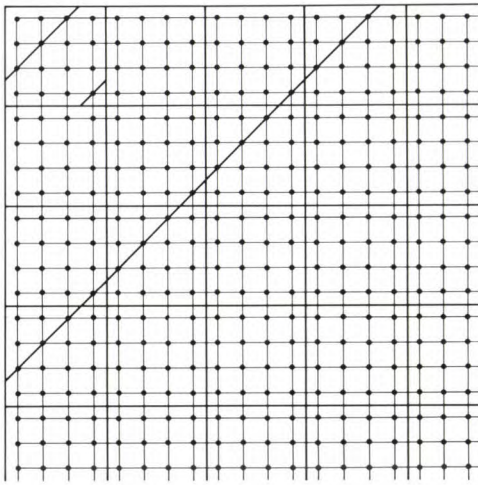


Fig. 6. Because of the way in which the tiles are interconnected, a diagonal line in the infinite grid is mapped on to only one or two diagonal lines in the torus.

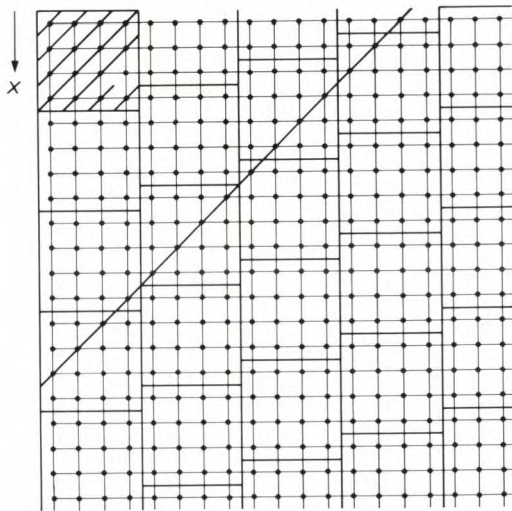


Fig. 7. Simulation of an infinite grid on a 'twisted torus'. The infinite grid is now built up from tiles that are always displaced by one position in the x-direction. A diagonal line in the infinite grid is then mapped on to a large number of diagonal lines in the torus.

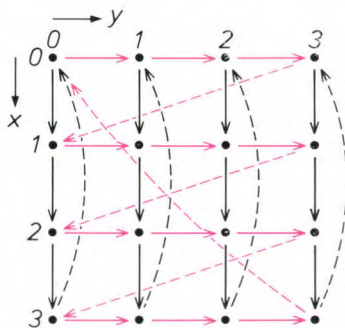


Fig. 8. A 'twisted torus' of 4 by 4 elements. The red arrows indicate the 'horizontal' connections; element $(x,3)$ is not connected with $(x,0)$ as in the straight torus, but with $((x + 1) \bmod 4, 0)$. The black arrows indicate the 'vertical' connections; element $(3,y)$ is connected with $(0,y)$ just as in the straight torus.

In considering the properties of various network structures further it is important to obtain a picture of the pattern of activities that arises in a network during a computation. It is easier to obtain such a general picture if a few simplifications are introduced.

We shall assume that all the cells and all the sub-computations have the same properties and that the propagation velocities are the same in both directions. First of all (when the nodes are assigned to the cells) there is a homogeneous expansion of the computation tree and then (when the cells are ready with all their nodes) it contracts again, almost homogeneously. During this process of expansion and contraction there is always a 'wavefront' of active nodes at a maximum distance from the root of the grid. In the infinite grid this wavefront approximates a diagonal, described by a relation $i + j = \text{constant}$. The wavefront is therefore called 'the active diagonal'. (It is only approximately a diagonal. Since different cells on a diagonal have to carry out different numbers of subcomputations, some are ready before others, so that the contraction process is not truly homogeneous.)

Fig. 6 shows that a diagonal in the infinite grid is mapped on to at most two diagonal lines in the finite grid. This has the result that at most M of the M^2 elements of the torus are active at the same time. If the wavefront has the form of a diagonal, the straight torus therefore leads to a poor distribution of the computation over the available processors.

The twisted torus

The reason for the unsatisfactory distribution of the computations for the straight torus is to be found in the symmetry of the tiles of fig. 6 about the axis $i = j$. To improve the situation we have to disturb this symmetry. One way of doing this is by displacing successive columns of tiles by one place, as shown in fig. 7. We see that the same active diagonal $(i + j) = 14$ is now mapped on to a large number of different diagonals in the finite grid. It is in fact possible to show that the distribution of the active diagonal over the finite grid is now optimum. If the active diagonal does not have more than M^2 nodes, these nodes are mapped on to different elements of the 'twisted' torus.

The tile pattern of fig. 7 is described by the relations

$$i = x + k * M - l$$

and

$$j = y + l * M,$$

where $M = 4$. A map of the infinite grid on a finite grid in accordance with this pattern is produced by connecting the cells $(x, M - 1)$ in an M by M square grid to $((x + 1) \bmod M, 0)$ and connecting the elements $(M - 1, y)$ to $(0, y)$. For the case $M = 4$ see fig. 8.

The effect that we produced by twisting the torus could also have been produced by using a rectangular torus consisting of M by P elements, where M and P are relative primes, instead of a square torus.

The twisted torus has an interesting advantage, however: a horizontal chain of nodes is mapped on to a cycle that contains *all* the elements of the torus (hence a cycle of length M^2). In the rectangular torus, however, such a chain is mapped on to a single row of the torus — a cycle of length M . The twisted torus therefore has advantages here, since the behaviour of the torus for horizontal chains is of practical interest. Such chains originate from computation trees that branch only to the right. The behaviour of the torus for vertical chains (trees branching only to the left) is however equally of interest. Now the twisted torus is no better than the rectangular one: both map such a chain on to a single vertical column.

The doubly twisted torus

There is a fairly obvious way of improving the situation outlined above: the torus should be twisted in both horizontal and vertical directions. To prevent both twists from introducing a new symmetry, we make them in opposite directions, e.g. $+1$ for the rows and -1 for the columns (see fig. 9).

There is found to be a significant difference between the straight torus and the singly twisted torus on the one hand and the doubly twisted torus on the other. A

have reached the same point by taking the same steps in a different order. A network with this property is called a ‘planar’ network. A doubly twisted torus does not have this property.

This can be illustrated with the aid of the simple example of the doubly twisted torus of three by three elements; see fig. 10. From A a ‘horizontal’ step (red line) followed by a vertical step (black line) leads to point B . In the reverse order, however, these steps lead to point C .

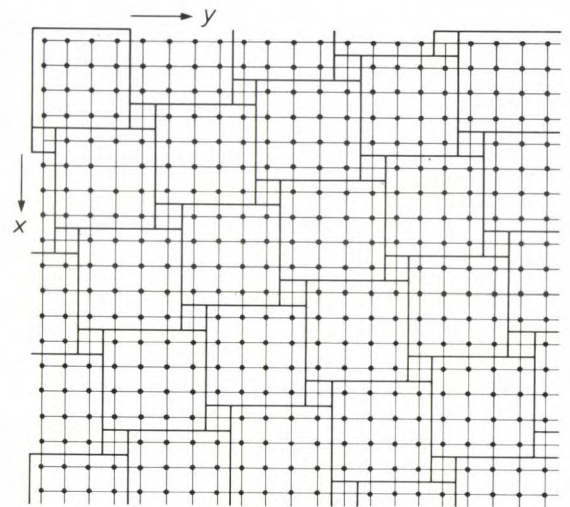


Fig. 9. If the tiles are twisted in both x - and y -directions, but with opposite signs, ‘holes’ are produced in the coverage of the infinite grid by the tile pattern.

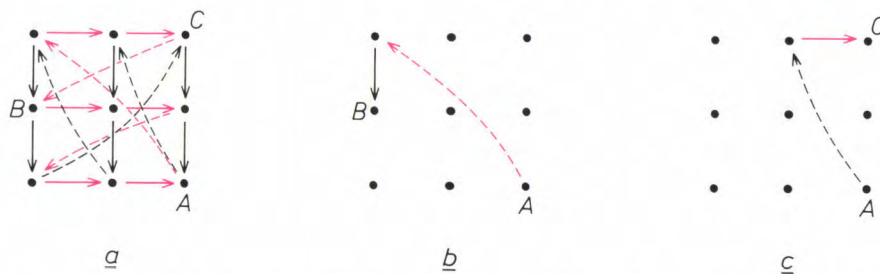


Fig. 10. a) A doubly twisted torus of 3 by 3 elements. The ‘horizontal’ connections are indicated by red lines, the ‘vertical’ ones by black lines. b) A ‘horizontal’ step from element A , followed by a vertical step, leads to element B . c) The same steps in another sequence lead to another element. A ‘vertical’ step from element A , followed by a horizontal step, leads to element C .

pattern of undisplaced tiles or tiles displaced in one direction covers the plane surface, but this is not true for a pattern of doubly displaced tiles: there are ‘holes’ in it, or overlaps, depending on the direction of displacement.

This is associated with an effect that arises on tracing a path through the grid. If we take a path from point A to point B in a straight or singly twisted torus by taking a certain number of horizontal steps and a certain number of vertical steps, we could also

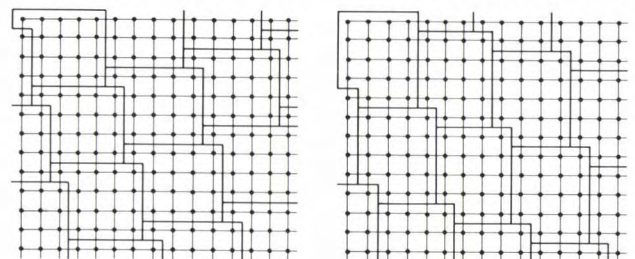


Fig. 11. Some patterns made with non-rectangular tiles, which fill the entire plane.

If the computation graphs that we wish to implement are always trees, this property is of no interest, since a tree always has a single path between a node and each of its sons, grandsons, etc. If on the other hand we wish to implement computation graphs that are not trees, then it is important that the torus should be planar. If because of the advantages mentioned earlier we still want to make use of the idea of the double twist, then we find we need patterns with non-rectangular tiles. Two such patterns are shown in *fig. 11*.

It is assumed above that each processor has its own memory module with direct access to it. It can only reach other modules by communicating with its neighbouring processors. When a processor consults a variable in the memory of a neighbouring processor, this results in a much more complicated activity than if it consults a variable in its own memory.

This discontinuity can be 'smoothed away' by using a slightly different hardware structure, so that we obtain a more uniform distribution of the computations over the surface. This is done by arranging the processors and memory modules in a network in such a way that each processor has access to at least two memory modules. A network of this type is called a 'continuous-processing surface' [6].

The future

The distributed systems discussed in this article are essentially based on the use of recursion as a programming technique. This is not found to be a serious limitation, however.

In the first place, recursive techniques have many applications. The divide-and-rule technique, which we

illustrated by the computation of binomial coefficients, is used for many numerical problems (e.g. computations on matrices and polynomials, and techniques such as the Fast Fourier Transform), and also for non-numerical problems (e.g. sorting and searching). Another widely used recursive technique is 'backtracking' — a technique for investigating all the 'candidate' solutions for a problem in a systematic way. This technique is used in compilers (for the syntactical analysis of the program being translated), and in many programs in Artificial Intelligence studies.

In the programming language LISP recursive technique play a leading part. This language has mainly been used in the university environment. However, interest in such recursion-oriented languages for commercial and industrial applications also appears to be increasing [6]. This indicates good prospects for the kind of distributed computation techniques that have been considered in this article.

[6] Details of this system arrangement and its consequences for the software are to be found in A. J. Martin, The torus: an exercise in constructing a processing surface, Proc. 2nd Caltech Conf. on VLSI, Pasadena 1981.

[6] M. Marshall, LISP computers go commercial, Electronics 53, Nov. 20, 1980, pp. 89-90.

Summary. A recursive computation can be represented by a tree structure. 'Fully distributed' computers for performing such computations can be designed in the form of a network of processors, each with its own memory. In executing a computation the 'computation tree' is mapped on to the 'implementation graph'. In this mapping process it is important to obtain a uniform distribution of the subcomputations over the available processors. Various network structures with different properties in this respect are compared.

Transformation methods for improving parallel programs

M. Sintzoff

A parallel program has the best chance of being executed efficiently if its components work as far as possible on their own local variables and interact with each other as little as possible. This article describes a method for converting a given program in a systematic way into a program that has better properties in this respect.

A similar method can also be used for dealing with a well-known problem encountered in the design of parallel programs: the occurrence of 'deadlocks' [1].

To describe these methods it is convenient to represent parallel programs by means of programs of a different kind, called *non-deterministic* programs, which have not yet been discussed in this issue. In a non-deterministic program, just as in a sequential program, only one operation is performed at a time. The difference is that the order in which the operations are performed is largely left open. All problems arising in parallel programs from the indeterminacy of the order in which the various instructions are executed therefore occur in the same way in non-deterministic programs.

First, a simple example shows how a parallel program can always be 'modelled' by a non-deterministic program. Next, a technique is presented that can be used for transforming parallel programs represented in this way. The transformation contains two stages. The result of the first stage is that the interaction between the components is limited to inspection of each other's variables. The result of the second stage is that the frequency of these inspections is minimized.

We then go on to show how similar techniques can be used for systematically eliminating the possibility of 'deadlock'. The article concludes with a brief discussion of the usefulness of the transformation techniques considered.

Parallel and non-deterministic programs

A non-deterministic program can be represented by means of a list of instructions, each preceded by the condition that must be fulfilled if a particular instruc-

tion is to be selected for execution. An instruction with its condition is referred to here as a 'guarded command'. The notation of a program of this kind is as follows:

$$\begin{aligned}
 P = & \text{ do } P_1: B_1 \rightarrow S_1 \\
 & // P_2: B_2 \rightarrow S_2 \\
 & \dots \\
 & // P_n: B_n \rightarrow S_n \\
 & \text{ od}
 \end{aligned} \tag{1}$$

P is the name of the complete program, and P_1, \dots, P_n are the names of the various guarded commands it contains (the *components* of the program). B_1, \dots, B_n are the conditions relating to P_1, \dots, P_n respectively; S_1, \dots, S_n are the commands.

The commands are carried out one by one, but the order is not fixed. Every time a command has to be designated for execution, an arbitrary choice is made from all the commands whose condition is fulfilled. As long as such commands still exist, execution of the program continues. The same command can thus come up for execution more than once.

It will now be shown how a parallel program may be represented by means of an 'equivalent' non-deterministic program. Given a parallel program one can construct a corresponding non-deterministic program without too much difficulty. For each operation of each process we can introduce a 'guarded command'. The 'guard' then contains conditions that correspond to P-operations, and 'if'-clauses in the original program, and also conditions that indicate the location and the process to which the command in question belonged.

Let us consider as an example the following parallel program consisting of two cyclic processes:

$$\begin{aligned}
 Q: & \text{ cycle } A; P(a); B; C; V(a) \text{ end} \\
 R: & \text{ cycle } D; P(a); E; F; V(a) \text{ end}
 \end{aligned} \tag{2}$$

The program uses one semaphore, a , which has the initial value 1.

A non-deterministic program that corresponds to this parallel program is the following:

Prof. M. Sintzoff, formerly with Philips Research Laboratory Brussels, Brussels, is now at the Université Catholique de Louvain, Belgium.

```

P = do Q1: CQ = 1      → A; CQ := 2
    // Q2: CQ = 2 and a = 1 → a := 0; CQ := 3
    // Q3: CQ = 3      → B; CQ := 4
    // Q4: CQ = 4      → C; CQ := 5
    // Q5: CQ = 5      → a := 1; CQ := 1
    // R1: CR = 1      → D; CR := 2
    // R2: CR = 2 and a = 1 → a := 0; CR := 3
    // R3: CR = 3      → E; CR := 4
    // R4: CR = 4      → F; CR := 5
    // R5: CR = 5      → a := 1; CR := 1
od

```

Two counters have thus been introduced, C_Q and C_R , which are used for maintaining the order of the operations within each of the two processes; they have the initial value 1. The semaphore a from the parallel program has been replaced by an identically named 'ordinary' variable, with initial value 1. The guarded commands Q_2 and R_2 correspond to the P-operation on a ; Q_5 and R_5 correspond to the V-operation.

Such a non-deterministic program describes a particular method of executing the original parallel program. This may be understood in the following way.

A parallel program consists of different processes that are carried out simultaneously at arbitrary speeds; each of these processes consists of a series of operations that are to be performed sequentially. One of the correct ways of carrying out such a parallel program would be to perform only one operation of one process in each case, while freezing all other processes; after each execution of an operation the next operation of an arbitrary process is then performed.

This 'stepwise' mode of execution is not efficient, of course, but it is theoretically interesting. This method of executing a parallel program is exactly what is described by the corresponding non-deterministic program. In such a description, problems that might arise from the truly simultaneous execution of operations from different processes are ignored. All problems that arise from the indeterminate nature of the order of execution of the operations (e.g. 'deadlock' and 'individual starvation') remain, however. Such problems can therefore be studied on the basis of non-deterministic programs.

If parallel programs are represented by non-deterministic ones, the possibility is created of applying exact methods of reasoning to their correctness. This is an important advantage of such a representation. We shall now consider a simple example to give an impression of this type of reasoning.

Example of a proof of correctness

Consider a non-deterministic program that consists of two components. One component describes the

behaviour of a producer of portions of information, the other describes the behaviour of a consumer of this information. The information is exchanged via a buffer, in which there is space for a maximum of N portions of information. The program has the following structure:

```

P = do P1: x < N → {production of a portion;
                    placing of this portion in
                    the buffer}
                    x := x + 1
    // P2: x > 0 → {retrieval of a portion from
                    the buffer; consumption of
                    this portion}
                    x := x - 1
od

```

The value of the variable x indicates the number of portions present in the buffer at each moment. The initial value of x is 0.

For this program to be correct, both components must maintain the validity of

$$0 \leq x \leq N.$$

This relation, which we shall call R , must be an *invariant* of the program.

To show that R is invariant, we make use of a general property of the instruction that assigns a value to a variable (the 'assignment statement'): the validity of a statement R is not upset by a guarded command of the form

$$B \rightarrow x := f(x),$$

if

$$R \text{ and } B \Rightarrow R[f(x)/x],$$

that is to say if $R \text{ and } B$ implies the statement that is obtained from R by systematically replacing each occurrence of x by $f(x)$.

For P_1 this does indeed apply: R is the relation $0 \leq x \leq N$, and $f(x)$ is $x + 1$; $R[f(x)/x]$ then stands for $0 \leq x + 1 \leq N$. The validity of R is therefore not affected by P_1 if

$$0 \leq x \leq N \text{ and } x < N \Rightarrow 0 \leq x + 1 \leq N,$$

which is indeed the case. In the same way it follows that the validity of R is not affected by P_2 if

$$0 \leq x \leq N \text{ and } x > 0 \Rightarrow 0 \leq x - 1 \leq N,$$

which is also the case. The validity of R is therefore not upset by either P_1 or P_2 : R is an invariant of the program P .

[1] See J. L. W. Kessels and A. J. Martin, Parallel programs, this issue, p. 254.

Increasing the independence of components

In programs with parallelism it would be advantageous to try to have each component executed by its own processor, and as far as possible simultaneously. The extent to which this can be done is seriously limited, however, as long as frequent interactions take place, such as the use of common variables or the exchange of messages, since such interactions require synchronization. We want to organize our components in such a way that for most of the time they are active only locally, i.e. perform operations only on private variables, and such that interaction between components is necessary only in exceptional cases.

In our model, interactions are reflected by operations on common variables. We shall try in a systematic way to reduce the frequency of such interactions. First we shall consider a method of transforming a program P into a program P' in which the value of each variable is modified only by one component, although the conditions of all components can still use the values of all variables. Secondly we shall indicate how the components thus obtained can be split into two parts, one that operates entirely independently of the other components, and another that is responsible for the necessary interactions [2].

Limitation of the write-operation to one component per variable

In the following the symbol x stands for 'the variables of the program'; it may thus be necessary to treat x as a vector. The following notation may now be given for the program to be transformed:

$$\begin{aligned} P = & \text{do } P_1: B_1(x) \rightarrow x := f_1(x) \\ & // P_2: B_2(x) \rightarrow x := f_2(x) \\ & // \dots \\ & // P_n: B_n(x) \rightarrow x := f_n(x) \\ & \text{od} \end{aligned}$$

We shall now try to represent the value of x by the values of n new variables u_1, u_2, \dots, u_n , one for each component. This means that a function T has to be found such that we can always write:

$$x = T(u_1, u_2, \dots, u_n).$$

This representation function must enable us to transform the program P into a program P' that uses, instead of the variable x , the variables u_1, u_2, \dots, u_n , and has the following form:

$$\begin{aligned} P' = & \text{do } P'_1: C_1(u_1, \dots, u_n) \rightarrow u_1 := f_1(u_1) \\ & // P'_2: C_2(u_1, \dots, u_n) \rightarrow u_2 := f_2(u_2) \\ & // \dots \\ & // P'_n: C_n(u_1, \dots, u_n) \rightarrow u_n := f_n(u_n) \\ & \text{od} \end{aligned}$$

The idea here is that the i^{th} condition C_i of the new program P' is only satisfied if the i^{th} condition B_i of the old version P is satisfied, and that the execution of the corresponding command $u_i := f_i(u_i)$ in P' has the same effect on $x = T(u_1, \dots, u_n)$ as the execution of $x := f_i(x)$ in P. For this it is sufficient that the following shall apply for each condition C_i :

$$C_i(u_1, \dots, u_n) \equiv B_i(T(u_1, \dots, u_n)),$$

and that the representation function T has for each value of i the property that

$$\begin{aligned} T(u_1, \dots, u_{i-1}, f_i(u_i), u_{i+1}, \dots, u_n) \\ = f_i(T(u_1, \dots, u_n)), \end{aligned}$$

or, in abbreviated notation:

$$T[f_i(u_i)/u_i] = f_i(T).$$

We shall illustrate the formula with the 'model' mentioned earlier for the finite information buffer:

$$\begin{aligned} P = & \text{do } P_1: x < N \rightarrow x := x + 1 \\ & // P_2: x > 0 \rightarrow x := x - 1 \\ & \text{od} \end{aligned}$$

Since $f_1(u_1) = u_1 + 1$ and $f_2(u_2) = u_2 - 1$, we seek a representation for x , $x = T(u_1, u_2)$, such that:

$$\begin{aligned} T(u_1 + 1, u_2) &= T(u_1, u_2) + 1, \\ T(u_1, u_2 - 1) &= T(u_1, u_2) - 1. \end{aligned}$$

It is not difficult to see that $T(u_1, u_2) = u_1 + u_2$ satisfies this condition, and thus the program P' becomes:

$$\begin{aligned} P' = & \text{do } P'_1: u_1 + u_2 < N \rightarrow u_1 := u_1 + 1 \\ & // P'_2: u_1 + u_2 > 0 \rightarrow u_2 := u_2 - 1 \\ & \text{od} \end{aligned}$$

We can interpret u_1 as the total number of portions delivered by the producer to the buffer, and $-u_2$ as the total number of portions which the consumer takes from the buffer. Then $u_1 + u_2$ is the number of portions delivered but not yet taken, and this number must of course lie between the limits 0 and N.

In the new program P' we see that for each guarded command P'_i the command changes only the value u_i , and that the new value of u_i depends only on the old value of u_i and on none of the other variables. The entire interaction is limited to the evaluation of the conditions of the guarded commands; in these conditions all variables may still occur. We shall make use of this property of P' later. Another advantage of components that do not disturb the values of each others' variables is that certain methods of proof for parallel programs are more readily applicable [3].

Splitting of the components

When the components of a program are executed in parallel, the actions on the variables of the program will within certain limits be intermingled in an unpredictable manner. A certain latitude is necessary here. In the example of the finite information buffer, this latitude arises because the number of places for portions in the buffer is greater than one; in a buffer for only one portion there will be strict alternation between production and consumption of portions of information, and there will then be no question of any unpredictable merging of producer and consumer activities.

The idea of the following transformation is to distribute the existing latitude among the various components from time to time in such a way that each of them can continue on their own for a while. Not until one of the components has used up its part of that latitude is it necessary to establish contact again with a view to a redistribution.

Formally, this amounts to saying that, in addition to the existing variables u_1, u_2, \dots, u_n , new variables u'_1, u'_2, \dots, u'_n have to be introduced, which must satisfy the following requirements:

— All u'_i are future values of the corresponding u_i , that is to say after the operation $u_i := f_i(u_i)$ has been performed j times, u_i assumes the value u'_i .

— Each component P'_i may carry out its operation $u_i := f_i(u_i)$ unhindered until u_i has assumed the value u'_i , that is to say the u'_i are chosen such that the condition $u_i \neq u'_i$ is sufficient to guarantee that the condition $C_i(u_1, u_2, \dots, u_n)$ is satisfied, provided at least that the other P'_j go no further than to $u_j = u'_j$.

Use is made here of the property of program P' , obtained by the transformation described earlier: the values that the variable u_i assumes do not depend on the values of any other variable; the interaction between the processes only determines the rate at which u_i runs through its range of values.

As soon as the variable u_i in a component P'_i assumes the value u'_i , an attempt must be made to find a new set of values for u'_1, \dots, u'_n that will satisfy the two requirements above, and for this an interaction is necessary. The transformed program is therefore approximately of the form:

```

P'' = do P''_1: u_1 ≠ u'_1           → u_1 := f_1(u_1)
      // Q''_1: u_1 = u'_1 and D_1(u_1, ..., u_n) →
                                     define new u'_1, ..., u'_n
      // ...
      // P''_n: u_n ≠ u'_n           → u_n := f_n(u_n)
      // Q''_n: u_n = u'_n and D_n(u_1, ..., u_n) →
                                     define new u'_1, ..., u'_n
od

```

The conditions D_i indicate that for given values of u_1, \dots, u_n new values have to be found for u'_1, \dots, u'_n that satisfy the two requirements and also the extra requirement that $u'_i \neq u_i$, so that after redistribution the component P_i can certainly continue.

To describe the requirements formally we introduce the abbreviated notation $K(m_1, m_2, \dots, m_n)$ for the statement that each of the commands $u_i := f_i(u_i)$ can certainly be executed m_i times unhindered:

$$K(m_1, \dots, m_n) \equiv \forall \mu_1, \dots, \mu_n \forall i: 0 \leq \mu_1 \leq m_1 \text{ and } \dots \text{ and } 0 \leq \mu_n \leq m_n \text{ and } \mu_i < m_i \Rightarrow C_i(f_1^{\mu_1}(u_1), \dots, f_n^{\mu_n}(u_n)).$$

The expression in the right-hand side of this definition states that the condition C_i for the command $u_i := f_i(u_i)$ always remains satisfied provided that this command has not yet been carried out m_i times (i.e.: $\mu_i < m_i$) and all other commands $u_j := f_j(u_j)$ at the most m_j times (i.e. $\mu_j \leq m_j$). The notation $f^\mu(u)$ gives the value that is obtained by applying the function f successively μ times — first to u , then to $f(u)$, etc. (The formal definition is thus: $f^0(u) = u$, and $f^{n+1}(u) = f(f^n(u))$.)

Although for brevity it has been omitted in the left-hand side of the definition of K , the value of $K(m_1, \dots, m_n)$ is dependent on the values of u_1, \dots, u_n . However, $K(m_1, \dots, m_n)$ does not depend on u'_1, \dots, u'_n .

The statement that a set of values u'_1, \dots, u'_n satisfies the two formulated requirements can now be expressed in the new notation as

$$K'(u_1, \dots, u_n, u'_1, \dots, u'_n) \equiv \exists m_1 \geq 0 \dots \exists m_n \geq 0: [u'_1 = f_1^{m_1}(u_1) \text{ and } \dots \text{ and } u'_n = f_n^{m_n}(u_n) \text{ and } K(m_1, \dots, m_n)].$$

In the example of the finite information buffer the expressions for K and K' yield simple and easily interpretable answers:

$$K(m_1, m_2) \equiv \forall \mu_1, \mu_2: (0 \leq \mu_2 \leq m_2 \text{ and } 0 \leq \mu_1 < m_1 \Rightarrow u_1 + \mu_1 + u_2 - \mu_2 < N) \text{ and } (0 \leq \mu_1 \leq m_1 \text{ and } 0 \leq \mu_2 < m_2 \Rightarrow u_1 + \mu_1 + u_2 - \mu_2 > 0).$$

The value of $u_1 + \mu_1 + u_2 - \mu_2$ is as large as possible when μ_1 is at a maximum and μ_2 at a minimum. Analogously we can indicate when $u_1 + \mu_1 + u_2 - \mu_2$

[2] For details see M. Sintzoff, Principles for distributing programs, in: G. Kahn (ed.), Semantics of concurrent computation (Lecture Notes in Computer Science 70), pp. 337-347, Springer, Berlin 1979.

[3] S. Owicki and D. Gries, An axiomatic proof technique for parallel programs I, Acta Informatica 6, 319-340, 1976.

becomes as small as possible. On the basis of these relations, K can then be reduced to:

$$\begin{aligned} K(m_1, m_2) &\equiv u_1 + (m_1 - 1) + u_2 - 0 < N \\ &\quad \text{and } u_1 + 0 + u_2 - (m_2 - 1) > 0 \\ &\equiv u_1 + u_2 + m_1 \leq N \text{ and } u_1 + u_2 - m_2 \geq 0. \end{aligned}$$

The first part here indicates that the producer can place m portions unhindered in the buffer if $m_1 \leq N - u_1 - u_2$, the number of empty locations in the buffer, while the consumer can take n portions from the buffer provided $m_2 \leq u_1 + u_2$, the number of portions already present in the buffer.

From the expression for K' we have in this case:

$$\begin{aligned} K'(u_1, u_2, u'_1, u'_2) &\equiv \exists m_1 \geq 0 \exists m_2 \geq 0 [u'_1 = u_1 + m_1 \\ &\quad \text{and } u'_2 = u_2 - m_2 \text{ and } K(m_1, m_2)] \\ &\equiv u'_1 \geq u_1 \text{ and } u'_2 \leq u_2 \text{ and} \\ &\quad u'_1 + u'_2 \leq N \text{ and } u_1 + u'_2 \geq 0, \end{aligned}$$

which indicates that u_1 may rise freely to a value u'_1 that does not lie above the total number of locations in the buffer plus the total number of portions already taken from the buffer; something similar applies for u_2 .

In terms of K and K' we can now specify the condition and the operation of each guarded command Q_i'' :

$$\begin{aligned} Q_i'' : u_i = u'_i \text{ and } K(0, \dots, 0, 1, 0, \dots, 0) \rightarrow \\ \text{define new } u'_1, \dots, u'_n \text{ such that} \\ K'(u_1, \dots, u_n, u'_1, \dots, u'_n) \text{ and } u'_i \neq u_i. \end{aligned}$$

In other words, for D_i we have taken K , using the value 1 for m_i and 0 for the other m_j , indicating that the operation $u_i := f_i(u_i)$ must unconditionally be performed at least once (otherwise it would be pointless to give u'_1, \dots, u'_n new values). The program envisaged for the finite information buffer thus becomes:

$$\begin{aligned} P'' = \text{do } P_1'' : u_1 \neq u'_1 &\rightarrow u_1 := u_1 + 1 \\ // Q_1'' : u_1 = u'_1 \text{ and } u_1 + u_2 < N &\rightarrow u'_1 := N - u_2 \\ // P_2'' : u_2 \neq u'_2 &\rightarrow u_2 := u_2 - 1 \\ // Q_2'' : u_2 = u'_2 \text{ and } u_1 + u_2 > 0 &\rightarrow u'_2 := -u_1 \\ \text{od} \end{aligned}$$

In the condition for Q_1'' we have no longer explicitly mentioned the condition $u_1 + u_2 \geq 0$, an invariant of the program. In the operation of Q_1'' we have chosen the maximum of all permissible values for u'_1 , which run from $u_1 + 1$ up to $N - u_2$ (see the requirements $u'_1 > u_1$ and $u'_1 + u_2 \leq N$), without for the time being adjusting the value of u'_2 (which we were able to choose from the range from $-u_1$ up to u_2). By making u'_1 as large as possible, we give the producer

the maximum latitude without it being at the expense of the consumer; by not adjusting u'_2 , we cause the minimum of disturbance to the consumer. Note, incidentally, that our formalism has imposed only boundary conditions. Within those constraints we may ourselves make a 'strategically optimum' choice for u'_1 and u'_2 .

A more general example

We shall briefly discuss another example: a finite information buffer with two producers and one consumer. The example can be interpreted as a rudimentary reversion system, in which there are two users of places, e.g. aircraft seats or hotel rooms (the 'producers' of reservations), and one party (the 'consumer') who makes reserved places available again. We model this example in the first instance as:

$$\begin{aligned} P = \text{do } P_1 : x < N \rightarrow x := x + 1 \\ // P_2 : x < N \rightarrow x := x + 1 \\ // P_3 : x > 0 \rightarrow x := x - 1 \\ \text{od} \end{aligned}$$

The first transformation uses the representation function $x = u_1 + u_2 + v$ and leads directly to:

$$\begin{aligned} P' = \text{do } P_1' : u_1 + u_2 + v < N \rightarrow u_1 := u_1 + 1 \\ // P_2' : u_1 + u_2 + v < N \rightarrow u_2 := u_2 + 1 \\ // P_3' : u_1 + u_2 + v > 0 \rightarrow v := v - 1 \\ \text{od} \end{aligned}$$

For the splitting of components (the next transformation) we have to determine K and K' . For these we find:

$$\begin{aligned} K(m_1, m_2, n) &\equiv u_1 + m_1 + u_2 + m_2 + v \leq N \text{ and} \\ &\quad u_1 + u_2 + v - n \geq 0, \end{aligned}$$

$$\begin{aligned} K'(u_1, u_2, v, u'_1, u'_2, v') &\equiv u'_1 \geq u_1 \text{ and } u'_2 \geq u_2 \text{ and} \\ &\quad v' \leq v \text{ and } u'_1 + u'_2 + v \leq N \\ &\quad \text{and } u_1 + u_2 + v' \geq 0, \end{aligned}$$

where $K(1, 0, 0) \equiv K(0, 1, 0) \equiv u_1 + u_2 + v < N$, while $K(0, 0, 1) \equiv u_1 + u_2 + v > 0$.

In Q_1'' , if $u_1 = u'_1$ and $u_1 + u_2 + v < N$, we must choose new values for u'_1 and u'_2 in such a way that we satisfy $u'_1 > u_1$, $u'_2 \geq u_2$ and $u'_1 + u'_2 + v \leq N$. For this purpose we divide the amount of latitude we have, $s = N - u_1 - u_2 - v$, as fairly as possible between u'_1 and u'_2 . In other words, both producers receive the same number of unoccupied places to reserve. If s is odd we give the extra place to u'_1 . In this way we at once guarantee that $u'_1 > u_1$, even if $s = 1$. In Q_1'' the following operations are therefore performed:

$$\begin{aligned}u'_1 &:= u_1 + (s+1) \text{ div } 2; \\u'_2 &:= u_2 + s \text{ div } 2,\end{aligned}$$

and in this way, because $s > 0$, we have indeed satisfied all the requirements. (The operation **div** is a division operation that rounds down to integers.)

This strategic choice is optimal if both producers on average provide a customer just as frequently. If in Q_1'' $s = 1$, i.e. if there is only one unoccupied place left that is in the possession of the second producer, it is transferred to the first producer, who can clearly reserve it immediately.

This finally brings us to the following version of the program:

```
P'' = do P1'': u1 ≠ u1' → u1 := u1 + 1
      // Q1'': u1 = u1' and u1 + u2 + v < N
              → u1' := (N + u1 - u2 - v + 1) div 2;
              u2' := (N - u1 + u2 - v) div 2
      // P2'': u2 ≠ u2' → u2 := u2 + 1
      // Q2'': u2 = u2' and u1 + u2 + v < N
              → u1' := (N + u1 - u2 - v) div 2;
              u2' := (N - u1 + u2 - v + 1) div 2
      // P3'': v ≠ v' → v := v - 1
      // Q3'': v = v' and u1 + u2 + v > 0
              → v' := -u1 - u2
od
```

Avoidance of deadlocks

Programs can be divided into two kinds: on the one hand programs that define a terminating computation and whose purpose is thus to be completed, and on the other hand programs relating to continuous control or administration, whose purpose is thus to be continued indefinitely. In the execution of a program in the second category it is necessary to avoid situations in which the 'guard' is not satisfied in any of the commands. Such situations are called 'deadlocks'.

The following artificial example will serve to illustrate this situation:

```
P = do P1: x < b → x := x + 1
      // P2: x < b and x ≥ a → x := x - 1
od
```

(a and b are two numbers such that $a < b$; the initial value of x is $a - 1$).

In the execution of this program a situation may be reached in which x has the value b . In that case neither the condition of P_1 nor the condition of P_2 is satisfied. The program then stops: P is not free from deadlock.

Using the method that we shall now discuss, it is possible to derive in a systematic way from a program P with deadlocks another program P' in which dead-

locks can no longer be reached. This is done by substituting for the conditions B_i of P new and more stringent conditions C_i in P' , which rule out any move to deadlock. These C_i 's can be systematically computed from the B_i 's.

For a program of the form

```
P = do P1: B1(x) → x := f1(x)
      // P2: B2(x) → x := f2(x)
od
```

we can set up the following equations that the stronger conditions C_1 and C_2 must satisfy:

$$\begin{aligned}C_1(x) &\equiv B_1(x) \text{ and } [C_1(f_1(x)) \text{ or } C_2(f_1(x))], \\C_2(x) &\equiv B_2(x) \text{ and } [C_1(f_2(x)) \text{ or } C_2(f_2(x))].\end{aligned}$$

It is clear that conditions C_1 and C_2 that satisfy these equations are restricted versions of the original conditions B_1 and B_2 : the equations have the form ' $C_i(x) \equiv B_i(x) \text{ and } \dots$ '.

If in the program P we substitute the stronger conditions C_1 and C_2 for B_1 and B_2 , we obtain the following program:

```
P' = do P1': C1(x) → x := f1(x)
       // P2': C2(x) → x := f2(x)
od
```

We can now show that in the execution of P' the statement ' $C_1(x) \text{ or } C_2(x)$ ' is always satisfied. A relation R is an invariant of P' if it is an invariant of both P_1' and P_2' , i.e. when:

$$\begin{aligned}R \text{ and } C_1 &\Rightarrow R[f_1(x)/x], \\R \text{ and } C_2 &\Rightarrow R[f_2(x)/x].\end{aligned}$$

From the equations for C_1 and C_2 it immediately follows that these statements are satisfied if we choose $R \equiv C_1 \text{ or } C_2$, which is thus an invariant of the program P' . It thus continues to be true that the condition of at least one of the components is satisfied; deadlock cannot therefore occur.

Solving the equations for the conditions

The question remaining is how to find the solutions of the equations for C_1 and C_2 . A standard technique consists in generating a series of increasingly better approximations for C_1 and C_2 by means of the recurrent computation

$$\begin{aligned}C_i^0(x) &\equiv \text{true} \\C_i^{n+1}(x) &\equiv B_i(x) \text{ and } [C_1^n(f_i(x)) \text{ or } C_2^n(f_i(x))].\end{aligned}$$

If we are lucky, the C_i^n 's no longer change in value after a certain index in the series. The value then

reached evidently satisfies our equations. In our example this situation does indeed arise:

$$\begin{aligned}
 C_1^0(x) &\equiv \text{true} \\
 C_2^0(x) &\equiv \text{true} \\
 C_1^1(x) &\equiv x < b \text{ and } [\text{true or true}] \equiv x < b \\
 C_2^1(x) &\equiv x < b \text{ and } x \geq a \\
 C_1^2(x) &\equiv x < b \text{ and} \\
 &\quad [x+1 < b \text{ or } [x+1 < b \text{ and } x+1 \geq a]] \\
 &\equiv x < b-1 \\
 C_2^2(x) &\equiv x < b \text{ and } x \geq a \text{ and} \\
 &\quad [x-1 < b \text{ or } [x-1 < b \text{ and } x-1 \geq a]] \\
 &\equiv x < b \text{ and } x \geq a \\
 C_1^3(x) &\equiv x < b \text{ and} \\
 &\quad [x+1 < b-1 \text{ or } [x+1 < b \text{ and } x+1 \geq a]] \\
 &\equiv x < b-2 \text{ or } [x = b-2 \text{ and } x \geq a-1], \\
 &\equiv x < b-2 \text{ or } x = b-2 \\
 &\quad (\text{since } b > a, b-2 \geq a-1) \\
 &\equiv x < b-1 \\
 C_2^3(x) &\equiv x < b \text{ and } x \geq a \text{ and} \\
 &\quad [x-1 < b-1 \text{ or } [x-1 < b \text{ and } x-1 \geq a]] \\
 &\equiv x < b \text{ and } x \geq a \text{ and} \\
 &\quad [x < b \text{ or } [x < b+1 \text{ and } x \geq a+1]] \\
 &\equiv x < b \text{ and } x \geq a.
 \end{aligned}$$

It can be seen here that $C_1^3 \equiv C_1^2$ and $C_1^2 \equiv C_2^2$, so that $C_1 \equiv x < b-1$ and $C_2 \equiv x < b \text{ and } x \geq a$ is the solution required.

In many examples this 'convergence' occurs after a very small number of steps. By way of illustration we shall take a classic example of deadlock: a program in which two components S and T compete with each other for the use of two resources a and b . These resources must together be available sometimes for S and sometimes for T. In the program for this situation S consists of three commands: S_1 to fetch resource a , S_2 to fetch resource b , and S_3 to return both resources (after use). The component T is similarly composed of three commands. T fetches the resources in the reverse order, however: first b , then a . The program then looks like this:

```

P = do S1: s = 1 and a = 1 → s := 2; a := 0
      // S2: s = 2 and b = 1 → s := 3; b := 0
      // S3: s = 3           → s := 1; a := 1; b := 1
      // T1: t = 1 and b = 1 → t := 2; b := 0
      // T2: t = 2 and a = 1 → t := 3; a := 0
      // T3: t = 3           → t := 1; a := 1; b := 1
od

```

In this program the variable s functions as a kind of command counter for the component S, which ensures that S_2 is only selected after S has executed the command S_1 . The counters a and b indicate the avail-

ability of the resources a and b . If the program starts in the state $s=t=1$ and $a=b=1$, and first S_1 and then T_1 are selected for execution, the state $s=t=2$ and $a=b=0$ arises. In this state none of the conditions is satisfied, and there is thus a deadlock.

The new conditions can be derived by means of the iterative procedure given above. The new condition for S_1 will have to satisfy the relation: new condition for $S_1 \equiv s = 1 \text{ and } a = 1 \text{ and } [\text{the disjunction of the new conditions for } S_1 \text{ up to } T_3 \text{ with } a \text{ systematically replaced by 0 and } s \text{ by 2}].$ The iteration then yields the following conditions:

after 0 iterations: all conditions true,

after 1 iteration: the old conditions,

after 2 iterations:

for S_1 : $s = 1 \text{ and } a = 1 \text{ and } [b = 1 \text{ or } t = 3]$

for S_2 : $s = 2 \text{ and } b = 1$

for S_3 : $s = 3$

for T_1 : $t = 1 \text{ and } b = 1 \text{ and } [a = 1 \text{ or } s = 3]$

for T_2 : $t = 2 \text{ and } a = 1$

for T_3 : $t = 3,$

after 3 iterations: as after 2 iterations.

So once again we see rapid convergence.

The solution can be interpreted such that S only fetches resource a if resource b is already free ($b = 1$) or is guaranteed to become free ($t = 3$). At the same time resource b is then reserved for S because, in the state with $a = 0$ and $s = 2$, T is not able to fetch resource b . The mechanism presented here for eliminating mutual blockage does not discover that $t \neq 3$ follows from $a = 1$.

Examples are also known in which the successive approximations do not become constant after a finite number of steps. In these cases a kind of limit transition is necessary. The treatment of such cases^[4] is beyond the scope of this article, however.

A stricter requirement than the avoidance of deadlocks is that each condition shall be satisfied within a finite number of steps, so that each command can be selected for execution within a finite number of steps. This requirement can also be satisfied by systematically substituting new conditions for the original conditions B_i ^[4].

Discussion

Because of the large number of possible execution paths in parallel programs, such programs are more difficult to design and analyse than deterministic sequential programs. For this reason exact design methods specifically devised for parallel programs are indispensable. In this article a number of mathematical

techniques are presented that may lead to a more satisfactory design process. Such techniques are necessary because design is the essential feature of programming, and subsequent verification is more difficult and subsequent testing usually incomplete.

Nevertheless, the approach described has its pros and cons. The first difficulty is that the techniques used are based on logic and mathematics: this requires a complete and rigorous specification of the problem as well as the manipulation of logical expressions and the application of mathematical induction. A quite different problem is the choice of the order of the various design decisions: should we, for example, first eliminate deadlock, or should we give priority to relative independence of the components? We lack the means for demonstrating the consistency of such strategic decisions, and we have no formalism for representing them in a clear and structured manner. Nor do we have any good criteria for choosing between various techniques that are available, for example to avoid 'individual starvation'.

One may wonder whether the exact design methods are in fact sufficiently applicable: all too often they are illustrated by trivial examples.

All in all, the situation is not too bad. The simplicity of the examples is in fact misleading. They are abstractions of characteristic types of problems; the simplicity was deliberate, and is a consequence of a successful abstraction process. Since the programs for these examples are arrived at in a systematic manner, they are by definition foolproof; in this way, moreover, new solutions often come to light.

Another use of rigorous design methods is to be found in the elucidation of ingenious but incomprehensible programs. If the concepts underlying them can subsequently be made explicit, such programs can be better understood, reconstructed and communicated.

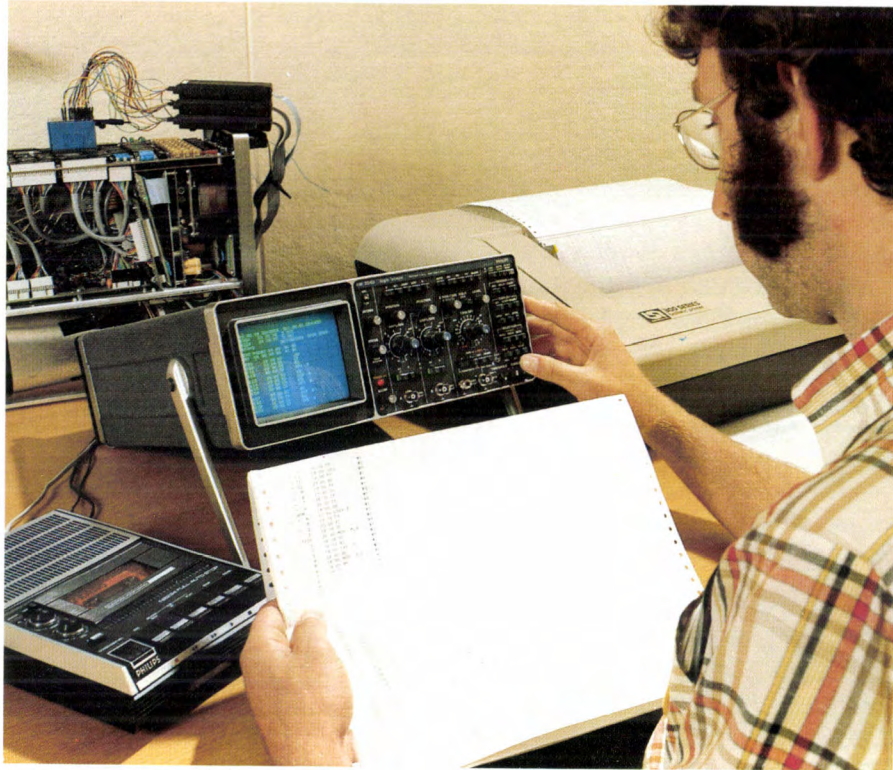
In actual cases it may happen that the methods dealt with here fall short of the requirements. Even then, however, they may be of some value: the corresponding design principles can still serve as guidelines in the development of programs. For example, when the assignments to a variable are limited to one component, it may be difficult to find a representation function that meets the specifications. In this case, however, we can use results from the field of representations for abstract data types, for which the equations given constitute a specification. Similarly, the ideas for reducing the interaction of the components, however difficult they may be to work out, are seen to be directly applicable to the design of more efficient implementations for distributed data bases^[5].

We may conclude that both informal methods and formal techniques have their value for the controlled design of parallel programs. A great deal of in-depth scientific research is still necessary, however, to develop this challenging area of programming.

[4] A. van Lamsweerde and M. Sintzoff, Formal derivation of strongly correct concurrent programs, *Acta Informatica* 12, 1-31, 1979.

[5] G. McLean Jr., Comments on SDD-1 concurrency control mechanisms, *ACM Trans. Database Syst.* 6, 347-350, 1981.

Summary. This article discusses a number of techniques that can be used in designing parallel programs when such programs are represented by means of non-deterministic programs. Parallelism in the components of a program is simpler the fewer interactions occur between these components, that is to say the less use is made of common variables. The article considers a method that makes it possible to derive from a given program an 'equivalent' program in which the necessary interactions are separated as far as possible from other operations, so that they do not unnecessarily restrict the possibilities of parallel execution of the program. In designing parallel programs it is necessary to ensure that no deadlocks can occur in their execution. The article introduces a method that makes it possible, on the basis of a program that might give rise to deadlocks, to systematically derive a program that is guaranteed to be free from deadlock.



The logic analyser is an aid in checking the operation of a logic circuit built around a microprocessor. The instructions received and sent out by the microprocessor are continuously read by the logic analyser; when a trigger signal is received, the last 256 instructions are stored in the memory, translated into a readable code and displayed on the screen. The correct operation of that part of the program can then be checked step by step. The PM3543 shown here is an oscilloscope as well as a logic analyser, and can also be used to display the pulse shape of the logic signals. The logic analyser is therefore useful for checking the actual logic hardware, whereas the Microcomputer Development System illustrated elsewhere in this issue is mainly an aid to the development of software.

Recent United States Patents

Abstracts from patents that describe inventions from the following research laboratories, which form part of or cooperate with the Philips group of companies:

Philips Research Laboratories, Eindhoven, The Netherlands	E
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	R
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	L
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	A
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	H
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	B
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	N

4 343 038

Magnetic bubble domain structure

U. E. Enz

E

A magnetic bubble domain structure comprising a bubble domain layer supported by a nonmagnetic substrate. A control layer is superimposed on the major surface of the domain layer which remote from the substrate. The control layer has an easy axis of magnetization in the plane which defines at least one bubble domain propagation path. The unique control layer according to the invention may be a garnet layer which comprises two sublayers: a first, continuous, sublayer, and a second, discontinuous, sublayer. The discontinuous sublayer defines the required propagation patterns. An advantage of the invention is that the in-plane rotary field to be applied for propagating the bubble domains can be considerably weaker than in bubble domain structures which are equipped with nickel-iron propagation patterns.

4 344 012

Anode disc for a rotary-anode X-ray tube

H. Hübner

B. Lersmacher

H. Lydtin

R. Wilden

A

The focal path of an anode disc in accordance with the invention is provided on a pyrographite ring which is oriented so that the surfaces of higher thermal and electrical conductivity extend parallel to the axis of rotation of the anode disc. As a result, suitable removal of heat can be ensured without thermal overloading of the bearings.

4 343 079

Self-registering method of manufacturing an insulated gate field-effect transistor

P. J. W. Jochems

E

A method of manufacturing an IGFET device in an entirely self-registering manner, in which on the semiconductor body a narrow silicon nitride strip is formed which covers only the active region of the body and the width of which is substantially equal to that of the transistors to be manufactured and possibly other circuit elements. This nitride strip is used as a mask for providing the channel stopper zone and as an oxidation mask for providing a first oxide layer. The nitride strip is then etched in which the strip is locally removed over its entire width and only parts remain above the channel region and contact regions which form a second oxidation mask and, in cooperation with the first oxide layer, a doping mask. The source and drain zones of the transistors and possibly further zones, for example underpasses, are formed via said doping mask after which by oxidation a sunken oxide pattern is formed over the whole surface with the exception of the channel regions and the contact regions. After the oxidation the remaining nitride may be removed by means of a maskless etching treatment after which the gate dielectric with the gate electrodes can be provided in a simple manner above the channel regions and the contacts, possibly preceded by the provision of contact zones, can be provided in the contact regions.

4 344 031

Method and device for verifying signals, especially speech signals

M. Kuhn

E. Bunge

H

The characteristic values are obtained from the signal to be verified, e.g. in the case of a voice signal the total energy in the individual frequency spectrum ranges, and these characteristic values are quantified. The frequency with which the individual partial ranges of the quantified characteristic values occur is found during a sampling phase, and these frequencies are stored. In the test phase, the signal is analyzed in the same way and the stored frequencies of the partial ranges into which the characteristic values of the signal to be verified have fallen are multiplied and compared with a threshold. In an arrangement for implementing the method, there is a characteristic store with a store address for a multi-digit binary number for each partial range of each characteristic. This store is addressed by a counter, the positions of which are allocated to the individual characteristics or spectrum ranges of the voice signal, and by the numbers of the partial ranges of the characteristics found during the analysis into which the characteristic values of the signal to be analyzed have fallen. In the sampling phase the content of each addressed store address is raised by 1 and in the test phase the frequencies obtained in the store addresses addressed are multiplied together. When a voice signal is processed for verification of the speaker, with the stepwise formation of the long-term spectrum, the intermediate values are compared with a range limit of a partial range read out from a range limit store and increased by 1 on exceeding a number allocated to this characteristic in an intermediate store. In this manner, it is possible at the end of the voice



signal, immediately and without any further processing steps to obtain the quantified characteristics in the form of the numbers of the partial ranges which can then be used to increase the frequencies in the characteristic store. A summation over all spectrum values is made to standardize the incoming voice signals. Subsequently, to all spectrum values is added a given fraction of their value, performed in the simplest manner by position shifting, until the sum of all spectrum values exceeds a predetermined constant. Thus the determination of the quantified characteristics or the partial range sums is immediately effective. The multiplication of the frequencies may be replaced by the addition of the logarithmic frequencies.

4 344 228

Shaving apparatus

F. Haes

C. M. Reynhout

G. M. P. G. Hermes

E. Boiten

There is provided a shaving apparatus comprising an external cutting element and an internal cutting element cooperating with and drivable relative to the external cutting element with a reciprocating movement. The internal cutting element is secured to the external cutting element by means of a connecting element which is flexible in the direction of driving and is rigid in directions extending substantially transversely of the direction of driving. The external cutting element is formed as a U-shaped part from a sheet material. The internal cutting element together with the connecting element constitutes a strip also made of a sheet material and positioned between the legs of the U-shaped part. The internal cutting element is situated at the end of the strip near the transverse portion of the U-shaped part; and the other end of the strip is connected to the legs of the U-shaped part.

4 345 012

Electrophotographic method of generating electrostatic images on two sides of an insulating foil

H. J. Hirsch

H. Dannert

Electrostatic charge images of identical shape but opposite sign are generated on both sides of a transparent, highly insulating foil. Subsequently, pigment is deposited on both sides of the foil by means of oppositely charged developers. The optical density of an electrophotographic image on a transparent insulating foil is increased, as compared to densities achieved in the past, for a given surface charge density by establishing a charge exchange between one side of the foil and an electrode. On the other side of the foil a charge image is generated and the foil. The electrode are separated from each other prior to development.

4 346 262

Speech analysis system

L. F. Willems

L. L. M. Vogten

In a formant speech analysis synthesis system, formant extraction to control a recursive digital all-pole filter encounters the problem that pole-pairs are not orderly arranged and that real poles may occur which are not representative of formants. The problem is solved by transforming the coefficients of the second-order sections of the filter to coefficients which can be easily ordered and by means of which it is simple to assign formants to the real poles.

4 346 267

Hybrid circuit

E. C. Dijkmans

A hybrid circuit comprising a current amplifier connected between the receive path and the transmission path of a four-wire transmission path, the common output of the current amplifier being connected to a two-way transmission path and the balancing impedance and an impedance transformer being provided between the other end of either the balancing impedance or the two-way transmission path, which results in a simple hybrid circuit.

4 347 288

High-voltage cable having a polythene containing insulation sheath which is provided with means to avoid or impede the formation or the growth of watertrees, the means comprising a metal complex

L. Minnema

G. N. van Boekel-Mol

A high voltage cable comprising a conductor and a polythene containing insulation sheath is provided with a metal complex of a diketone, salicylic acid optionally substituted with one to two lower alkyl groups, or a Schiff's base formed from an amine and salicylaldehyde optionally substituted with one to two lower alkyl groups, in order to prevent or impede the growth of watertrees during use.

4 347 441

Dual probe interferometer for object profile measuring

J. G. Dil

J. C. Driessen

W. Mesman

An apparatus is described for accurately measuring the profile on an object. The apparatus comprises two probes for simultaneously scanning the object to be measured and the reference object, the objects being rotatable about the same axis. Each of the probes is provided with a reflecting element, which are respectively incorporated in the measuring arm and reference arm of an interferometer.

4 347 610

Control circuit for the drive current of a laser

L. J. Meuleman

Control circuit for controlling a laser drive current. In known control circuits certain levels of the optical output signal are kept constant. This renders it difficult to satisfy certain requirements during the life of the laser as regards the linearity of the modulation of the optical output signal. In the described control circuit the intermodulation of two signals present in the modulation signal with substantially constant amplitudes is measured, an error signal which controls a component of the drive current being derived therefrom.

4 347 994

Magnetic tape drive arrangement

J. de Boer

A magnetic tape drive arrangement without pressure roller, having a signal generator which generates a signal, derived from the capstan motor current and the capstan speed. The signal is proportional to the instantaneous torque exerted on the magnetic tape by the capstan, which signal controls the supply and/or take-up motors in such a way that this signal is minimized, so that a slip-free magnetic-tape drive without pressure roller is possible.

4 348 484

Precision pressed glass object, glass, method of preparing glass

H. J. M. Moorman

H. Verweij

J. Haisma

Glass objects, lenses in particular, can be pressed with precision from a glass which contains 45-55 mol.% P₂O₅, 15-40 mol.% BaO, 5-15 mol.% Li₂O, 5-35 mol.% PbO, 0-2 mol.% Al₂O₃ and 0-6 mol.% F.

4 348 610

Camera tube with graded tellurium or arsenic target

J. Dieleman

E

J. H. J. van Dommelen

P. J. A. M. Derks

A camera tube target which is to be scanned on one side by an electron beam comprises a selenium-containing vitreous layer also containing tellurium and arsenic. In order to improve various target characteristics and properties in accordance with the invention, the concentration of at least one of the elements tellurium and arsenic in a first sub-layer of the selenium-containing layer, on the side to be scanned, increases towards the side to be scanned to a value at which the sum of the concentrations of tellurium and arsenic is at most 30 at.%. The arsenic concentration everywhere in the layer exceeds 1.5 at.%. In a second sub-layer adjoining the first sub-layer the concentration of at least one of the elements arsenic and tellurium is smaller than its concentration in an adjoining third sub-layer. A fourth sub-layer may also be present between the third sub-layer and the signal electrode. The concentration of arsenic and/or tellurium in the third sub-layer is larger than its concentration in the fourth sub-layer.

4 348 706

Edge guide for a magnetic tape

B. P. Videc

E

A tape recorder having a tape guide member with a guide edge against which the tape is loaded or biased by pressure members. To press the tape against the supporting edge with a transverse force approximately constant per unit of length, pressure members are arranged beyond the ends of the tape-supporting edge, and the tape supporting edge has a curvature convex toward the tape, with a peak substantially midway between tape-support entry and exit points having a smaller radius of curvature than that of the remainder of the tape-supporting edge.

4 348 795

Method of manufacturing cooling blocks for semiconductor lasers

A. H. Deunhouwer

E

H. G. Kock

A method of manufacturing cooling blocks for semiconductor lasers, in which the rounding-off radius of the line of intersection between two surfaces of the cooling blocks must have a very small value. In the method, two bodies to be formed into cooling blocks are each provided with a flat surface and these bodies are secured together with their flat surfaces by means of a curable adhesive. One side of the bodies connected together is subjected, transverse to the two surfaces connected together, to a machining treatment so as to obtain a further flat surface, in which machining treatment deformation and burring of the bodies near the line of intersection to be formed is avoided due to the presence of the cured adhesive, and a line of intersection having a rounding-off radius of only a few microns is formed.

4 349 893

Memory with current-controlled serial-to-parallel conversion of magnetic field domains

N. J. Wiegman

E

K. E. Kuijk

A device for magnetic domains contains at least two, meander-shaped, current conductors for the purpose of driving parallel domains along these conductors by means of respective currents in them that alternate cyclically. At at least one end the current conductors are connected to a third current conductor. The latter is either also meander-shaped or it acts in the same way as a meander-shaped conductor with respect to the domains. Conversion is possible between parallel drive along the first and second conductors and serial drive along the third conductor owing to the fact that a loop of the latter also forms part of the other conductors.

4 350 920

Dispenser cathode

T. C. J. M. Bertens

E

A compressed dispenser cathode having 1-15% by weight of $Ba_2Sc_2O_6$ in a porous metal cathode body has an emission having a larger current density (exceeding $6 A/cm^2$) and a longer life (more than 3000 hours).

4 352 015

Anti-contamination diaphragm for an electron beam apparatus

A. Jore

E

G. G. P. van Gorkom

N. H. Dekkers

In order to reduce object contamination during the examination or machining of the object in an electron beam apparatus, an anti-contamination diaphragm is provided in place of the customary diaphragm. The anti-contamination diaphragm has a central aperture which corresponds to a customary diaphragm aperture, and a concentric annular aperture for transmitting a non-paraxial hollow beam which irradiates a ring around the paraxial focus on the object. The conical hollow anti-contamination beam forms a barrier against contaminating residual gas molecules, notably hydrocarbon molecules.

4 352 041

Rotary anodes for X-ray tubes

H. Hübner

A

B. Lersmacher

H. Lydtin

R. Wilden

An intermediate layer comprising several sub-layers is sandwiched between the support and a target layer of a rotary X-ray anode. The sub-layer of the intermediate layer which contacts the support and the sub-layer of the intermediate layer which contacts the target layer both consist of pure rhenium. Interposed between these two sub-layers is a further sub-layer consisting of a rhenium alloy containing at least one carbide-forming metal, for example tungsten, tantalum or hafnium. This construction of the intermediate layer provides a barrier against carbon diffusion, which barrier has substantially the heat conduction properties of metals and which offers a sufficient protection against the penetration of carbon into the target layer, even at temperatures above 1500 K.

4 352 131

Memory disc addressing device

A. van Herk

E

J. A. M. Buis

J. H. Klijnstra

J. van Staalduinen

A device for determining a complete track number of a track which forms part of a set of tracks on a disc-shaped record carrier, enabling direct control by means of signals detected from servo sectors during each phase of positioning, and making a calculated estimate of a new track number on the basis of the displacement speed of the positioning member. The device reads a first track number in the group and estimates a second group number. It then estimates a complete track number on the basis of the displacement speed of positioning member body, determining the difference between the complete track number read and the estimated complete track number and comparing the magnitude of the difference obtained with a value determined by half a signal being produced if said magnitude is smaller than half, said signal validating the complete track number present in the first determination. A complete track number is determined by the sequential use of said estimates after each servo sector. Various embodiments of the circuits used for estimating a track number in this process are presented.

4 352 954

Artificial reverberation apparatus for audio frequency signals

N. V. Franssen

E

An artificial reverberation apparatus for audio frequency signals, comprises a first delay device preceded by an adder. A feedback circuit couples the output of the delay device to an input of the adder to give a loop signal gain of less than unity. The adder is preceded by a second delay device with the same delay time as said first-mentioned delay device. The signal to be delayed is applied to the input of said second delay device and also to the adder via a transmission path. The ratio of the signal gain of the transmission path to the signal gain of said second delay device is equal to but of opposite sign to said loop signal gain. Preferably the signal gain of the transmission path is equal to but of opposite sign to the signal gain of said feedback circuit. Furthermore, a plurality of apparatuses can be connected in cascade with the delay devices of the different apparatuses all having different delays.

4 353 062

Modulator circuit for a matrix display device

J. H. J. Lortetje

G. Warrink

H. W. Schneider

E

A modulator circuit for a matrix display device, the modulator circuit having pulse width and pulse amplitude control. The current amplitude for a selected picture element (row address-column address) varies during the pulse width of an excitation pulse in correspondence with the number of counting positions of a counting circuit used for determining the pulse width.

A column conductor has an associated column excitation circuit in which a counting circuit for determining the excitation pulse width controls excitation switches which pass selectively currents from current source S to the column conductor. The initial setting of the counting circuit determines the excitation pulse width, and the count positions of the counting circuit select the currents to be passed to the column conductor as the counting circuit is stepped by clock pulses to measure the duration of the excitation pulse.

4 353 130

Device for processing serial information which includes synchronization words

M. G. Carasso

J. G. Nijboer

E

A data stream is received from a medium. This stream consists of synchronization words and data words. The synchronization words are either identical or one the inverse of the other when they are correctly received. Between two synchronization words a fixed number of n data words is present. The information received is always consecutively stored in a buffer memory. The buffer memory has connected to it a detection device for generating an instantaneous synchronization signal by way of a majority decision on at least three correctly received synchronization words. Preferably, the buffer memory is a shift register comprising a data output which is situated approximately $\frac{1}{2}n$ data words $+\frac{1}{2}$ synchronization word beyond the center of the shift register when n has an even value.

4 353 623

Leadthrough for electric conductors

H. F. L. Maier

H

While avoiding the quartz fusion technique, an accurate and vacuum-tight electric leadthrough is provided which in addition is mechanically rigid. The electric conductor is fixed at the interior of a cell through an aperture by means of a ring of polytetrafluoro-

ethylene in such a manner that the conductor is held in the center of the aperture. Over the ring and inside the aperture, a layer of cyano acrylate talcum cement is provided. Inside this layer in the aperture, a layer of a granular filler which is soaked with cyano acrylate adhesive is provided. Finally, the outermost part of the aperture is filled and sealed with a layer of cyano acrylate talcum cement.

4 353 782

Method of solid-state pyrolysis of organic polymeric substances

B. Lersmacher

A

The yield of good-quality bodies consisting of vitreous carbon is increased to substantially 100% when the pyrolysis is performed in a reaction vessel in which a collecting container or a collecting disc for the condensable volatile decomposition products which is open at its upper side is arranged above the polymeric substances to be pyrolyzed. The collecting disc divides the interior of the reaction vessel into two parts, a passage for gases remaining between the upper part and the lower part. During the pyrolysis, a spatial temperature distribution is adjusted in the reaction vessel which results in a pulsating evaporation and condensation of the decomposition products, so that the decomposition products are completely removed.

4 354 103

Optical focusing device with two controls

K. A. Immink

M. P. M. Bierhoff

J. P. J. Heemskerck

E

A focusing system for focusing a radiation beam onto an object, in particular an optical disc. The focusing system comprises focusing-error detection means and two control elements for correcting the focusing depending on the detected focusing error. The detection means are constructed so that via one of the two control elements a feedback servo system is obtained while via the other control element a servo system without feedback, i.e. a feed-forward system, is obtained.

4 354 157

Method of and device for determining a nuclear spin density distribution in a part of a body

L. F. Feiner

E

It has been found that devices for determining nuclear spin density distributions in an object by means of nuclear spin resonance measurements produce artefacts in the calculated density distributions. The artefacts occurring are dependent of the gradient magnetic fields which are used during the measurements and which influence the excited nuclear spin. The invention concerns two filters for elimination of the artefacts in dependence of the kind of modulation of the gradient magnetic field.

4 354 186

Picture display device for displaying a binary signal generated by a picture signal generator as a binary interlaced television picture

H. H. H. Groothuis

E

In an interlaced picture display device by means of which signal patterns are displayed which are uniform from field to field as obtained from character generators vertical jitter and flicker phenomena are greatly reduced by attenuating the uppermost picture element of a sequence of subjacent bright picture elements in one field in combination with attenuation of the bottommost picture element in the other field or in combination with the addition of an attenuated picture element of the bottommost picture element of one field.

Manufacture of LaserVision video discs by a photopolymerization process

H. C. Haverkorn van Rijsewijk, P. E. J. Legierse and G. E. Thomas



After its introduction in the United States, towards the end of 1978, the Philips LaserVision system is now on the market in Europe. The players are mass produced in the Philips factory at Hasselt, Belgium, and video discs are made in the Mullard LaserVision Disc Centre at Blackburn, England. The disc manufacture at Blackburn is based on an entirely new process: the information is transferred by curing a liquid polymerizable lacquer on the mould or 'stamper' with ultraviolet light rather than by using the classical method of impressing into a plastic. In 1974 an investigation was started at Philips Research Laboratories into the possibilities

Ir H. C. Haverkorn van Rijsewijk is with the Philips Audio Division, Eindhoven; P. E. J. Legierse is with the Philips Plastics and Metalware Factories (PMF), Eindhoven; Dr G. E. Thomas is with Philips Research Laboratories, Eindhoven.

of applying this process to the manufacture of video discs. After the first fully playable disc had been made in this way in 1976, a team including a large number of colleagues from various product divisions investigated the prospects for mass production. The outcome of this intensive cooperation was the start of pilot production in 1980. After successful completion of the acceptance tests, full production started in 1981 at Blackburn. The European introduction of the LaserVision system seemed to us an appropriate occasion for including two articles on the new process. In the first article the process is compared with other methods and it is explained how it is used in current manufacture. The second article deals in some detail with the investigation of photopolymerizable lacquers for LaserVision discs.

Introduction

The 'video disc' in the Philips LaserVision system is the size of an audio long-play record. Picture and sound information is recorded on it as a succession of small pits of variable length and repetition frequency. The information is read out optically by the player, so that the read-out system does not come into direct contact with the disc. The operation of the player and the system, and the possible applications, have been dealt with by many Philips authors, both in this journal [1] and elsewhere [2].

Not much has yet been published, however, on the manufacture of video discs. This is because until recently they were produced with the pressing techniques used in the manufacture of audio long-play records. The video discs introduced in the United States were also made with these techniques. The information on these discs is recorded to suit the NTSC colour transmission system, which makes them unsuitable for use in the European countries where the PAL system is used [3]. For the production of video discs in Europe Philips are employing a new method, based on a photopolymerization process known as the '2p process' (from *photopolymerization*). In this process a liquid polymerizable lacquer is cured on an information-carrying mould (sometimes called a 'stamper') by exposure to ultraviolet light.

The first ideas on this application of the 2p process, and actual trials, came from the Philips Research Laboratories in Eindhoven. As often happens, the introduction of an entirely new technology into mass production posed a number of unforeseen problems. Solving the problems that arose on the introduction of the 2p process required the professional skills of many people from several Philips product divisions and the Research Laboratories. The combined effort of this team led first to a successful period of pilot production, followed by the introduction of the process in the production of PAL video discs for the European market, now carried on at Blackburn.

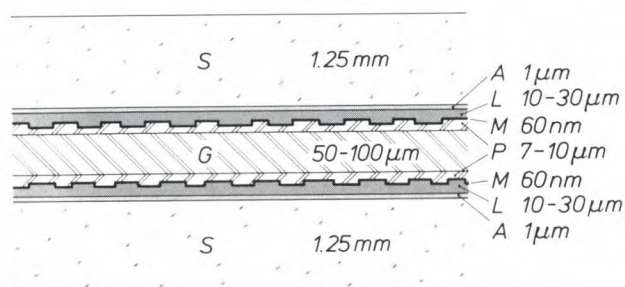


Fig. 1. Diagram showing the configuration of a double-sided LaserVision disc, as now made at Blackburn. The hole at the centre is not shown. *S* transparent substrate. *A* primer layer. *L* lacquer with picture and sound information in the form of pits. *M* mirror coating. *P* protective layer. *G* adhesive layer.

In this article we shall first give a brief description of the LaserVision system in its current form. We shall then examine the manufacture of the discs, considering the reproduction methods in general and the 2p process in particular. We shall then deal with some of the requirements that the video disc has to meet, and with the consequences for the choice of the materials and process control. Finally we shall describe the various materials and the stages in the present production process, indicating some alternatives that are under consideration.

The Philips LaserVision system

The LaserVision disc now being produced at Blackburn consists of two halves, with an external diameter of 300 mm and a central hole of diameter 35 mm, bonded 'back-to-back' with an adhesive. The two disc halves consist of a transparent substrate coated on one side with a primer layer and a lacquer layer in which the picture and sound information is recorded; see *fig. 1*. Since the information is read from the disc in a reflection process, a thin highly reflective metal 'mirror' is applied to this 'information' layer. This in turn is coated with a protective layer to prevent chemical and mechanical damage, which might occur during later stages in the manufacture. The protective layer also serves as a base for the adhesive that bonds the two discs together to form the complete video disc. The information is recorded between the radii of 55 and 145 mm, as a spiral track of small pits 0.12 μm deep and 0.4 μm wide; see *fig. 2*. The length of the pits and the minimum distance between them vary between about 0.5 and 2.0 μm, and the pitch of the spiral is 1.6 μm to 2.0 μm.

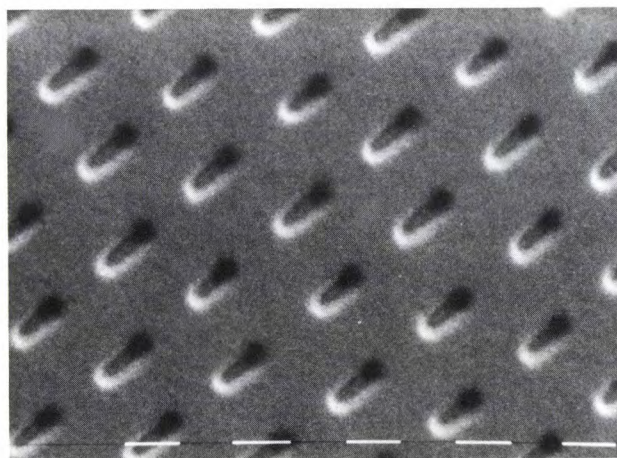


Fig. 2. Scanning-electron microscope (SEM) picture of part of a LaserVision video disc; the dashes correspond to a length of 1 μm. The information track consists of small pits of constant width and depth, with variable length and spacing.

The information is read optically with monochromatic light at a wavelength of 632.8 nm from a helium-neon laser [4]. A moveable lens guided by a control mechanism [5] keeps the laser light focused through the substrate on to the reflecting surface carrying the information. The intensity of the reflected light is modulated by the pattern of pits. The reflected beam is intercepted by a light-sensitive diode that converts the intensity variations into electrical signals [6]; see *fig. 3*. These in turn are translated electronically into picture and sound signals appropriate to a standard television receiver [7].

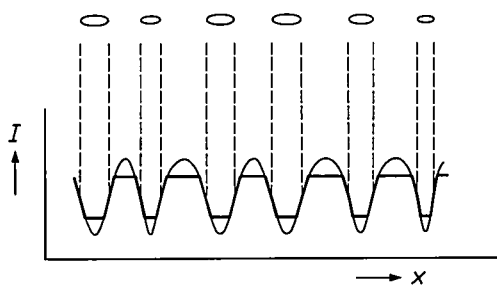


Fig. 3. Diagram showing the relation between the pattern of pits and the LaserVision signal on read-out. The laser light, modulated by the pits, produces a current I in the photodiode of the player. This current depends on the position x along the information track. The trapezoidal waveform is produced by symmetrical amplitude limiting.

Read-out in the LaserVision system can take place at constant angular velocity or at constant linear velocity. In the first case the disc rotates at a constant speed of 25 rev/s, corresponding to the frame frequency of the European television systems (25 Hz), so that each turn of the spiral track contains exactly one frame. As explained earlier in this journal [5], this makes it possible to 'manipulate' television pictures: by moving at the right moment from one turn to another, stationary, speeded-up, slow-motion and reverse-motion pictures can be obtained. The playing time, determined by the speed of rotation, the pitch of the track and the available area on the disc, is more than half an hour on each side. About twice the playing time can be obtained by reading at constant linear speed from a disc in which the information per turn increases with the radius of the turn. The speed of rotation is then gradually reduced as the disc is read

from inside to outside. This method is to be preferred when a long playing time is more important than the ability to manipulate pictures. Both types of disc are currently available in the LaserVision system; the player automatically adapts to either.

Manufacture of LaserVision discs

The manufacture of the video disc begins with a 'master', a glass disc coated with a layer of photoresist. Picture and sound information from a 'master' video tape is written into this layer by means of modulated laser light of peak power about 200 mW, so that it is recorded in the form of a pattern of tiny pits; each side of the disc contains about 25 000 million of these pits. As is usual in the production of gramophone records, a metal copy of the vulnerable master is made by electrochemical methods; this is the 'father' mould. After this negative copy has been detached, the master disc is no longer usable. Further electrochemical copying of the father mould produces the 'mother' moulds, which are positive copies of the master. Copying a mother mould then gives the negative copies that are used as the production mould for mass production. This 'family' process enables many production moulds to be made from a single master disc.

We shall not consider these stages further, but will confine our attention to the following stages, in which the video discs are manufactured from a production mould. We shall now review the methods that can be used for replicating the pattern of pits in the mould.

Reproduction methods

For manufacturing video discs an obvious course was to use the standard method for producing audio long-play records. In this method a preheated plastic blank is placed on the information-carrying mould under a press, the press is then closed and heated to a temperature at which the plastic softens. The plastic is then pressed over the mould under high pressure. The plastic fills up the available spaces and therefore takes up the information contained in the mould. After cooling, the disc is removed from the press. An advantage of this method is that a complete side of the disc, with the correct outside diameter and including the centre hole, is obtained in a single operation. The high pressure and temperature, however, can cause some deformation of the mould, so that local non-circularity of the tracks is introduced. There may also be residual stresses in the plastic, producing undesirable birefringence or a warped disc.

Another familiar method is injection moulding. In this method granulated plastic is melted and injected

[1] Philips tech. Rev. 33, 178-193, 1973.

[2] Appl. Optics 17, 1993-2036, 1978.

[3] F. W. de Vrijer, Philips tech. Rev. 27, 33, 1966.

[4] See for example K. Compaan and P. Kramer, Philips tech. Rev. 33, 178, 1973.

[5] P. J. M. Janssen and P. E. Day, Philips tech. Rev. 33, 190, 1973.

[6] G. Bouwhuis and P. Burgstede, Philips tech. Rev. 33, 186, 1973.

[7] W. van den Bussche, A. H. Hoogendijk and J. H. Wessels, Philips tech. Rev. 33, 181, 1973.

into the mould cavity. After cooling of the mould cavity a hard plastic disc is obtained, containing a negative copy of the information in the mould. This method is used, for example, in the mass production of small gramophone records (singles). The advantages and disadvantages are comparable with those of the pressing process.

In view of the special requirements that the video discs have to meet, a number of alternative reproduction methods were investigated, both at Philips and elsewhere^[8]. In one of them, called the 'printing method', a heated mould is used to impress the information at high pressure and temperature into the surface of a plastic disc at room temperature. The heating softens the surface of the disc so that it can take up the information contained in the mould; the information is then 'frozen in' by fast cooling. In principle only the surface of the disc is affected in this process. However, since it is difficult to obtain disc blanks that are sufficiently flat, the discs must be pre-heated so that complete plastic deformation is possible. In addition there are the disadvantages mentioned earlier of high pressure and temperature, and extra operations are necessary for making the centre hole and obtaining the correct outside diameter.

Another method, using the 2p process referred to in the introduction, was found to have many attractive aspects. These led to the introduction of this process for current production.

Application of the 2p process

In the 2p process, monomers of acrylates (esters of acrylic acid) are polymerized by exposure to ultraviolet light^[9]. The polymerization is started by the addition of a photo-initiator. *Fig. 4* is a diagram of phases in the 2p process. A few millilitres of the 2p lacquer — one or more monomers with a suitable photo-initiator — are applied to the centre of a mould. A transparent plastic disc, the substrate, which is slightly deformed to make it convex, is placed on the mould. The substrate is then pressed flat again by applying a slight excess pressure, causing the lacquer to flow out from the centre to the outside edge, covering the surface of the mould. The space between the mould and the substrate is thus filled by a thin layer (10-30 μm) of liquid lacquer, so that the information contained in the mould is completely transferred to the lacquer. The lacquer is then exposed to ultraviolet light at a wavelength of about 350 nm through the substrate, causing the lacquer to polymerize and become hard. The composition of the 2p lacquer is selected such that it does not adhere to the mould after curing but does adhere to the substrate; a primer layer may be used if necessary^[9]. After the

exposure, the video disc — the substrate with the cured lacquer — is removed from the mould. The mould can be used again immediately for making the next disc.

This process can be used with substrates that already have a centre hole. Further mechanical operations are necessary to give the discs the correct external diameter. Since low-viscosity lacquers are used, the 2p process does not require high pressures, and it can take place at room temperature. Compared with the more conventional thermosetting process, the 2p process has the advantage that no time need be lost in heating the substrate and lacquer and cooling them again. The thin lacquer layer can first be applied between mould and substrate while it is still liquid. The curing does not start until the exposure, and then it is almost instantaneous. An additional advantage of working at room temperature is that there is no thermal shrinkage after curing, and consequently very little stress in the material.

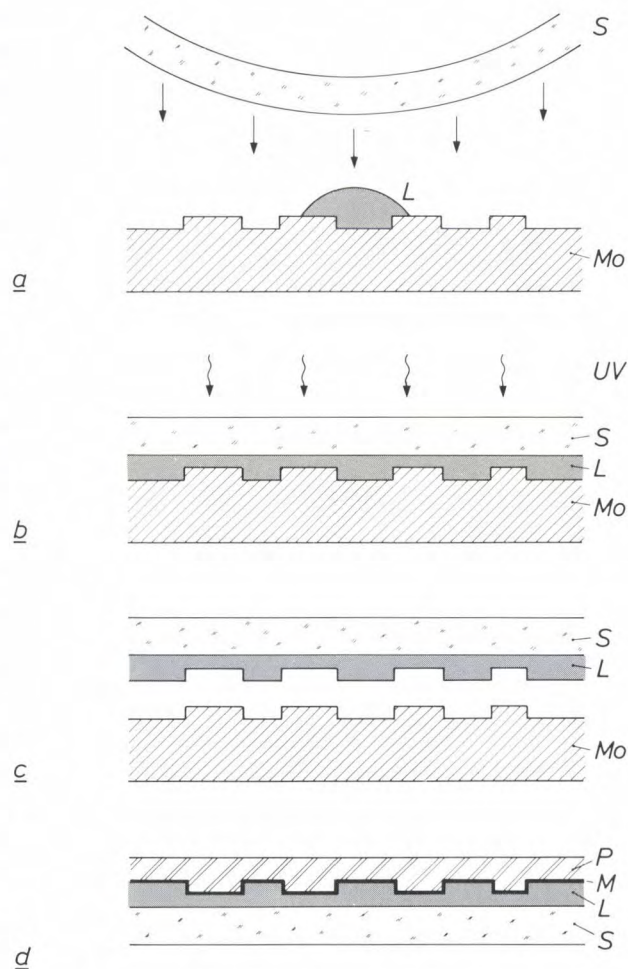


Fig. 4. Diagram showing four phases of the 2p process. *a*) The liquid lacquer *L* is spread over the mould *Mo* by a slightly deformed substrate *S*. The primer layer on the substrate and the centre hole are not shown. *b*) Exposure to ultraviolet light to cure the lacquer. *c*) Substrate with lacquer separated from mould. *d*) Lacquer coated with mirror *M* and protective layer *P*.

Requirements to be met by the LaserVision disc

A video disc on a LaserVision player, operating in accordance with the specifications, must be capable of delivering a good television signal to the television receiver. The specifications of the disc and the player are therefore matched to one another. This means that the disc must satisfy a number of requirements relating to optics, geometry and stability. We shall now take a closer look at these requirements.

Optical requirements

For the stored information to be read properly there must be at least a 75% optical reflection of the laser beam. The absorption and scattering losses that occur as the beam goes through the substrate twice depend strongly on the type of material. The same applies to the losses that arise on reflection from the metal mirror coating. Because a high reflectance is required, very few combinations of materials are suitable for the substrate, 2p lacquer and mirror coating.

Another optical requirement is connected with the use of linearly polarized laser light in the optical system of the player, for separating the incident from the reflected beam. Any change in polarization when the beam returns through the substrate makes this separation less effective, causing a drop in the efficiency of signal detection and an undesired feedback of the modulated light to the laser. This means that the birefringence in the substrate must be no greater than 20°, leading to a further restriction in the choice of materials and in the method of making the substrate.

Geometrical requirements

The specifications laid down for the player and the disc are such as to ensure that any disc can be played well on any player, and also that there is an optimum compromise between the ease of manufacture and production costs of the player on the one hand and of the disc on the other.

The thickness of the substrate was set at 1.25 mm. Thinner substrates are more difficult to make, and moreover with a thickness of 1.25 mm any scratches on the outer surface of the disc are so far outside the depth of focus of the lens of the player that they are no longer troublesome. The tolerance in the average thickness is fairly large, about 100 μm. The specification for the flatness and parallelism of the two surfaces, however, is very tight. Any departure from flatness in the rapidly rotating disc causes vertical movements in the information track. The focus-control system of the player must be able to follow these movements. Translated into the specifications for the focus-control system, the compromise mentioned above corresponds to:

- a maximum vertical deviation of 1050 μm at a frequency $f \leq 1.1$ kHz;
- a maximum vertical velocity of 0.18 m/s at $f \leq 1.1$ kHz;
- a maximum vertical acceleration of 100 m/s² at $f \leq 1.1$ kHz;
- a maximum vertical deviation of 2 μm at $f > 1.1$ kHz.

It follows that the curve of the permissible vertical deviation x plotted against the frequency f consists of four segments, which are shown in blue in fig. 5. The shaded region indicates the excursion that the focus-control system must be able to follow, as a function of frequency.

To calculate the consequences of this for the disc, we assume that there is a simple harmonic disturbance of wavelength λ along the track on the disc surface, with the amplitude in the axial direction. The frequency at which this disturbance is 'seen' by the focus-control system during read-out depends not only on λ

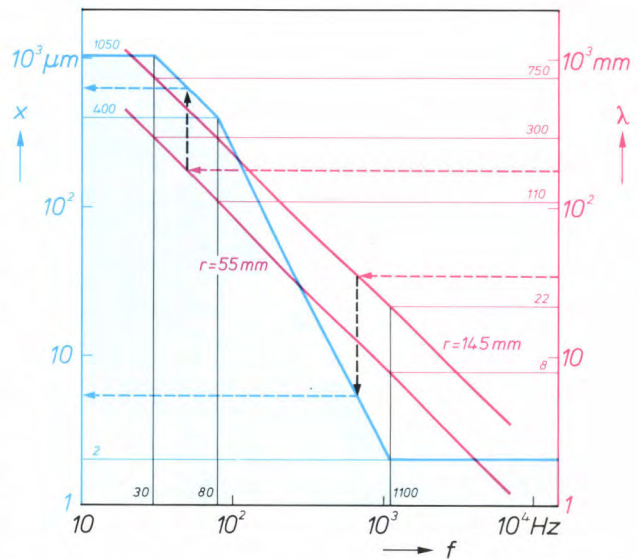


Fig. 5. Blue: Permitted vertical amplitude x as a function of the frequency f , as specified for the focus-control system of the LaserVision player. The curve consists of four segments, with a different specification for each segment (see text). The shaded area indicates the excursions that the control system must be able to follow. Red: Wavelength λ of an axial disturbance plotted against the frequency f at which the control system 'sees' this disturbance while the disc is being played. The two lines relate to a constant speed of 25 rev/s at the minimum read-out radius (55 mm) and at the maximum read-out radius (14.5 mm). Combining the red and blue curves indicates the maximum permitted vertical amplitude of a disturbance at wavelength λ and read-out radius r . The dashed lines give two examples of the results of such a combination. The values of f , x and λ corresponding to the changes in the slope of the blue curve are shown along the axes (see Table I).

[8] J. S. Winslow, IEEE Trans. CE-22, 318, 1976.
 D. G. Howe, H. T. Thomas and J. J. Wrobel, Photogr. Sci. Engng 23, 370, 1979.
 [9] J. G. Kloosterboer, G. J. M. Lippits and H. C. Meinders, this issue, p. 298.

but also on the speed of rotation and the read-out radius of the disc. Fig. 5 indicates in red the relation between λ and f at a constant disc speed of 25 rev/s for the minimum and maximum read-out radii (55 and 145 mm). By combining the red and the blue curves in fig. 5 we can read off the permissible vertical amplitude of an assumed sinusoidal disturbance on the disc at a particular read-out radius and a particular wavelength of the disturbance. The principal results are collected in *Table I*.

The specifications are also matched for the radial control system of the player and the radial deviations of the track on the disc. The specifications for the player are:

- a maximum radial deviation of 80 μm at $f \leq 2.2$ kHz.
- a maximum radial acceleration of 20 m/s^2 at $f \leq 2.2$ kHz.
- a maximum radial deviation of 0.1 μm at $f > 2.2$ kHz.

Fig. 6 indicates in blue and red the relation between the permissible radial deviation y and f and the relation between the wavelength λ of the radial disturbance and f . The consequences for the disc in relation to the permissible amplitude of the radial disturbance are shown in *Table II*.

The specifications for the radial deviations are fairly tight considering the many sources of deviations. These can arise even during the process of making the production moulds. There may be bumps and indentations in the disc, and these can act as a prism and cause radial deflection of the laser light. This is seen by the player as non-circularity of the track. In operation, imbalance in the disc can cause shocks and vibrations in the player.

The tolerance on the outside diameter of the disc is 1 mm. An additional requirement, however, is that the imbalance force during playback should be less than 1.5 N. This imbalance is caused by local thickness variations or more usually by eccentricity. Excessive imbalance can often be corrected by deliberately trimming the outside edge eccentrically.

Stability requirements

The disc must meet its specification immediately after manufacture and during its use by the consumer. The conditions during use are simulated as closely as possible in ‘climate tests’. Among the requirements are that use at a temperature of about 40 °C should have no adverse effect on the optical and geometrical characteristics of the disc. The disc must also remain unaffected by the climatic conditions likely to be encountered during transport and storage. Finally, the disc should have an acceptable useful life. This is

Table I. Permitted vertical amplitude x of an axial disturbance at the minimum and maximum read-out radius r on a LaserVision video disc played at a constant speed of 25 rev/s (see fig. 5).

λ (mm)		x (μm)
at $r = 55$ mm	at $r = 145$ mm	
> 300	> 750	1050
$300 \rightarrow 110$	$750 \rightarrow 300$	$1050 \rightarrow 400$
$110 \rightarrow 8$	$300 \rightarrow 22$	$400 \rightarrow 2$
< 8	< 22	2

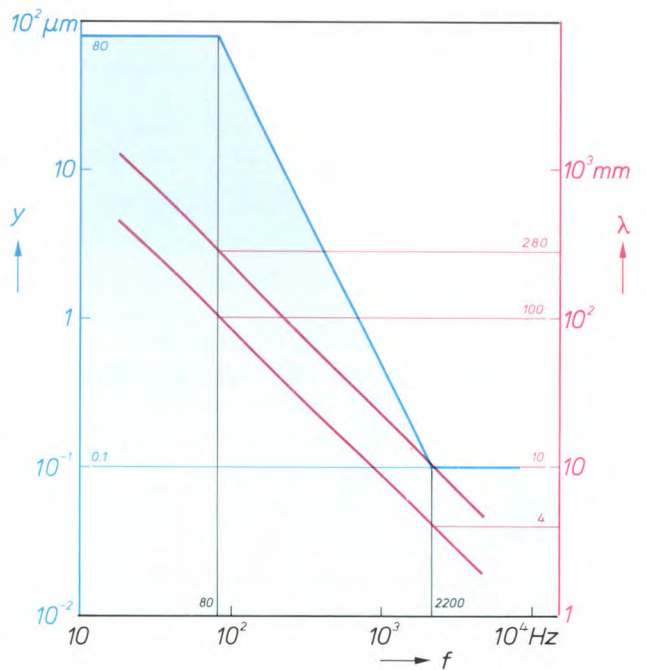


Fig. 6. *Blue:* Permitted radial amplitude y as a function of the frequency f , as specified for the radial control system. The curve here consists of three segments; the shaded area indicates the excursions that the radial control system must be able to follow. *Red:* As in fig. 5, but now for a radial disturbance. The values of f , y and λ corresponding to the changes in the slope of the blue curve are shown along the axes (see *Table II*).

Table II. As in *Table I*, but now for the permitted amplitude y of a radial disturbance (see fig. 6).

λ (mm)		y (μm)
at $r = 55$ mm	at $r = 145$ mm	
> 100	> 280	80
$100 \rightarrow 4$	$280 \rightarrow 10$	$80 \rightarrow 0.1$
< 4	< 10	0.1

tested, e.g. in cyclic humidity tests in which ageing is accelerated by subjecting the disc to frequent alternations of temperature (from 20 to 45 °C, and vice versa) and relative humidity (from 50 to 96%, and vice versa).

Choice of materials and process stages

The requirements mentioned above are clearly rather strict for a mass-production process, and obviously restrict the choice of the materials that can be used. This restriction applies not only to the separate components of the disc but also to the complete assembly. This is related to the complicated structure of the disc (fig. 1) and to the considerable differences in the chemical properties of the materials. Changing one of the materials or one of the stages in the process therefore generally requires an appropriate change in one or more of the other materials or process stages. It is consequently not easy to find the optimum production process. It will therefore come as no surprise that we are also investigating and developing other variations of the production process and testing them in pilot production. We shall now describe the choice of materials and the stages in the current production, together with a number of promising variations.

or even small scratches on the outer surface are usually of little significance, since the relatively large substrate thickness brings this surface well away from the focus of the laser beam. The lens of the player has a very large numerical aperture (0.4), so that the laser beam has a relatively large diameter (about 1 mm) when it enters the substrate.

In spite of this large measure of freedom, however, the choice of substrate material remains limited. For economic reasons it is necessary to use plastics, and only three of these will meet the optical and mechanical requirements: polymethyl methacrylate (PMMA), polyvinyl chloride (PVC) and polycarbonate (PC). *Table III* summarizes some of the advantages and disadvantages of these three materials for use as the substrate in the 2p process. The information given is based on our experience with test samples obtained from production and development departments of outside suppliers.

Table III. Properties of polymethyl methacrylate (PMMA), polyvinyl chloride (PVC) and polycarbonate (PC) for possible application as substrate material in the 2p process for the manufacture of LaserVision discs.

Property	PMMA	PVC	PC
absorption at 630 nm	about 5%	about 15%	about 5%
birefringence	acceptable	marginally acceptable, depending on preparation	problematic
glass temperature	about 110 °C	70 °C	140 °C
dimensional stability	sensitive to moisture	vulnerable at high temperature	good
method of manufacture	extrusion, casting, injection moulding	calendering + pressing	injection moulding, difficult because of high glass temperature
price	reasonable high reasonable	attractive	high
compatibility with 2p process	good, but adhesion layer necessary	good, but mirror coating can be attacked by additives and decomposition products	good, but tendency to stress cracking in contact with 2p lacquer

Choice of substrate

An important difference between the 2p process and the mechanical (pressing) reproduction methods is that the substrate in the 2p process serves purely and simply as an inert carrier of the 2p information layer. It would therefore seem that there is considerable freedom in the choice of substrate and in the method of making it. In general terms this is true. The disc blanks can be made from cast, extruded or rolled plastic sheet, by turning or punching. On the information side of the substrate a limited micro-roughness and thickness variation are permissible because of the smoothing action of the 2p lacquer. Some roughness

Polycarbonate is not used at present because it is too expensive for the mass production of such large discs and poses problems of birefringence. With PVC as substrate material, problems are mainly encountered with the mirror coating. The optical attenuation of the read-out laser beam in the PVC substrate is fairly high: about 15%. To meet the requirement of at least 75% of total reflection, the mirror coating must then have an inherent reflectance of at least 90%, which is difficult to achieve in practice, because of factors such as degradation during life. When they are combined with a PVC substrate, aluminium and silver coatings can be attacked by substances emanating

from the PVC, even with 30 μm of 2p lacquer between the substrate and the metal coating. As Table III shows, PMMA has good optical properties: the absorption at 630 nm is only about 5%, and the birefringence is generally well within the required limits. Because of these properties we have decided for the present to use PMMA as substrate material for the mass-production process.

Compared with the two other substrate materials, PMMA has the disadvantage of being less compatible with the 2p process, because a cured 2p lacquer does not adhere well to a PMMA surface. The substrate surface is therefore given a pretreatment, such as the application of a thin adhesion layer of an acrylic resin.

Another problem with PMMA relates to the water absorption of the material. Unlike PVC, for example, PMMA can absorb a relatively large amount of water vapour (a few per cent by weight), causing the material to swell. Since this is a reversible process, the water content of a PMMA disc is affected by the water content of the ambient air. During the uptake or release of water there are gradients in the water concentration in a PMMA disc, and the associated mechanical stresses can lead to warping of the disc. Fortunately these deformations are largely eliminated when the two disc halves are bonded together to form a symmetrical structure. Nevertheless, careful control of the humidity conditions remains strictly necessary in order to keep the component discs flat during the production process.

Choice of the 2p lacquer

The picture quality of a LaserVision disc depends to a large extent on the properties of the 2p lacquer. The read-out process takes place at the interface between the lacquer layer and the mirror coating, and the depths and lengths of the pits, the slopes of the pit walls, the micro-roughness of this interface and micro-defects in its proximity all play a part in this process. The lacquer must take up the mould information as faithfully as possible, and retain it during the polymerization and afterwards. The cured lacquer layer must also be capable of withstanding the thermal and mechanical stresses during the application of the mirror coating, and in the finished disc it must possess good chemical and mechanical stability.

In addition there are requirements connected with the mass production of the video discs. The lacquer, for example, must have a sufficiently low viscosity to allow a thin layer to be spread out sufficiently fast over the mould surface. It must also polymerize rapidly on exposure to ultraviolet light and the cured lacquer must not adhere to the mould surface. Other con-



Fig. 7. General view of the 2p machine and related equipment. The machine on the right applies the 2p lacquer in manufacturing the LaserVision discs. The equipment on the left includes a post-exposure unit and a mechanical device for transferring the discs with cured lacquer to a holder.

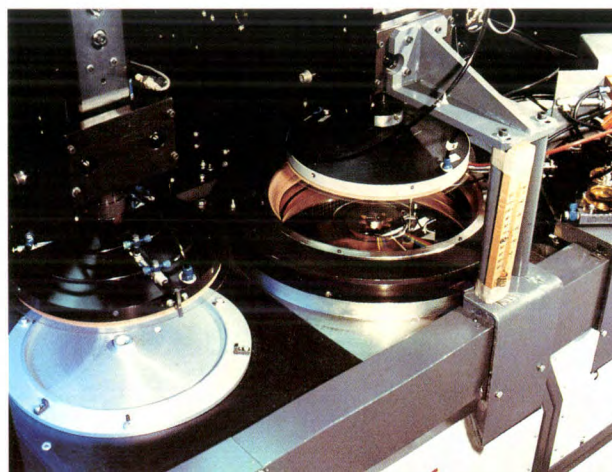


Fig. 8. Photograph showing the disc holder (on the right) for picking up the substrate and spreading out the 2p lacquer. A ring of liquid 2p lacquer can be seen on the mould. The disc holder presses the substrate against the mould, producing a uniform information-carrying layer of lacquer. On the left are the claw device and disc holder that separate the substrate with the cured lacquer from the mould.

ditions relate to matters such as availability, reproducibility and toxicity in the liquid state.

The requirements and conditions can be satisfied reasonably well by lacquers that contain one or more acrylates^[9]. With a suitable photo-initiator, these lacquers polymerize very rapidly during the ultraviolet exposure. A number of these acrylate layers give a reproduction accuracy that is more than sufficient for disc manufacture. They can also be sufficiently hard, so that the lacquer layer is not deformed on application of the mirror coating. To achieve the required hardness it is usually necessary to have a polymer whose chains are chemically cross-linked by side branches. A suitable polymer can be obtained from a mixture of a diacrylate and a triacrylate. By choosing

the correct mix and adding a thinner, the viscosity in the liquid state can be kept sufficiently low. The basic material can also be a single monomer such as hexanediol diacrylate [9].

Application of the 2p layer

Fig. 7 shows the machine designed for applying the 2p layer. After a substrate (with primer layer) has been placed in position, a measured quantity of liquid 2p lacquer is applied and spread out over a nickel mould. First, the area containing the information is exposed to ultraviolet light. Then the superfluous, still liquid part at the unexposed outer edge is removed. The remaining lacquer is then exposed to the ultraviolet light again, to complete the curing, especially at the outer edge, and the substrate with cured lacquer (the 'disc') is removed from the mould.

Some details of the operation of the 2p machine can be explained with reference to fig. 8. The substrate is picked up by a 'disc holder', which applies a holding vacuum along the outer edge of the disc and a small excess pressure in the inner area, causing the disc to sag slightly in the middle. A ring of the liquid 2p lacquer is applied around the pin at the centre of the mould, the centring pin. The disc holder is then lowered, and after alignment by means of the centring pin the substrate is pressed flat against the mould. As a result the lacquer is squeezed out to form a circular uniform layer 10 to 30 μm thick. The layer thickness depends on the excess pressure in the disc holder, the rate of spread and the viscosity of the lacquer. After the layer has been spread, the holder is withdrawn and the disc on the mould is transported to an exposure rack consisting of ultraviolet-emitting fluorescent lamps. The two-stage hardening process takes place here in less than 10 seconds. The disc and the mould are then transported to the separation position, where a claw device and a second disc holder separate the disc from the mould. Outside the 2p machine all the discs are given a final ultraviolet exposure directly on to the 2p lacquer to complete the curing process. The mirror coating can then be applied.

Depending on factors such as the viscosity and the curing behaviour of the lacquer, the total cycle time of the 2p machine is at present 30 to 40 seconds. A considerably shorter cycle time seems feasible, particularly since it may be possible to reduce the curing time to about 1 second.

Depositing the mirror coating

The mirror coating should have a sufficiently high reflectance at the read-out wavelength, since at least 75% of the incident laser light should be reflected. As we have seen, the optical attenuation in a PMMA sub-

strate is about 5%, so that the inherent reflectance of the mirror coating must be at least 80%. For this reason, and in connection with material costs, stability and reactivity, the choice of material is limited to the metals aluminium, silver and copper. These have theoretical reflectances of 85, 99 and 95%, respectively, at the read-out wavelength. These values also apply generally to thin layers with a thickness of only a few tens of nanometres.

Until now little attention has been paid to copper as a possible mirror material, mainly because of the attractive appearance of the disc with an aluminium or silver layer. At the start of disc production it was decided to use vacuum evaporation (vapour deposition) as the method of applying the metal coating. An evaporated silver layer, however, has poor adhesion to the 2p surface, so that a finished disc could fracture at this interface. Although this problem can be solved by chemical treatment of the 2p surface, extra stages in the production process would be necessary. Aluminium, on the other hand, adheres well to the 2p surface, and this material has therefore been chosen for the time being.

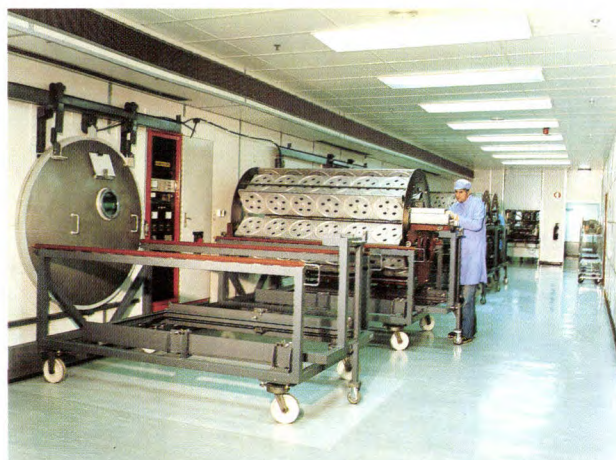


Fig. 9. Photograph of the evaporation area (above) and a carousel holder (below) for the deposition of aluminium on LaserVision discs.

Batches of about 120 discs are aluminized on a rotating carousel holder (*fig. 9*) by evaporation from simple aluminium-coated resistance elements. In this deposition process a problem arises because of the use of PMMA as substrate material. Owing to its water-vapour absorption, PMMA in vacuum desorbs a great deal of water vapour for some considerable time. Evaporation of aluminium at an insufficiently low pressure leads to a lower reflectance, owing to the reaction between water vapour and the growing layer. It was therefore necessary to design special evaporators with very large liquid-nitrogen-cooled cryopump surfaces, that would enable a pressure of 10^{-3} Pa to be reached easily in a short time.

Good mirror coatings can in fact be made by evaporation, but the batch-loading method is not easily integrated into the production process, which is otherwise continuous. The large carousel holders have to be loaded and unloaded, and for each batch a relatively long time is required for pumping down the evaporation chambers and readmission of air after deposition. Labour-intensive maintenance is also necessary to ensure that dust in the chambers, stirred up during pumping and readmission of air, does not give rise to pinholes in the mirrors. For these reasons two alternative metallization techniques are now in development, which may lead to more efficient production.

Alternative metallization techniques

A suitable metallic mirror can also be obtained by 'electroless' silver deposition, a wet chemical process in which the entire surface is coated with a thin layer of silver. A process of this type has long been used for producing glass mirrors. In the electroless silver deposition on a cured 2p lacquer the surface of the lacquer is simultaneously sprayed with an ammoniacal aqueous solution of a silver salt and a solution of a suitable reducing agent. When the two solutions mix on the surface metallic silver is formed. As with an evaporated silver coating, the adhesion to the 2p surface is poor. The surface is therefore given a pretreatment with an aqueous solution of an organic compound that modifies the surface of the 2p lacquer. After further treatment with an inorganic agent containing tin to increase adhesion, the deposition of the silver is started.

Since all the stages in this method of metallization consist of liquid-spraying operations, it was possible to develop a continuous production machine in which the discs undergo all the required operations sequentially on a conveyor belt. A prototype of this machine is illustrated diagrammatically in *fig. 10*, and a photograph is shown in *fig. 11*.

This metallization method takes on average much less time per disc than evaporation. Continuous rinsing of the surface with liquids almost completely prevents the occurrence of pinholes in the silver coating, as required for the LaserVision discs. *Fig. 12* shows a transmission-electron-microscope (TEM) photograph of a silver coating. For this photograph the adhesion treatment was omitted, so that the metal coating could easily be separated from the 2p lacquer and mounted as a separate film in the transmission electron microscope. The photograph covers an area somewhat lar-

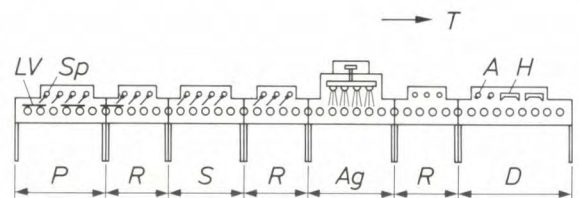


Fig. 10. Schematic prototype of machine for electroless silver deposition. *LV* LaserVision disc on a conveyor belt for transport *T*. *P* pretreatment. *R* rinsing. *S* sensitization. *Ag* silver deposition. *D* drying. *Sp* spray nozzle. *A* air blade. *H* heater element.



Fig. 11. Photograph of the machine for electroless silver deposition. This machine incorporates two metallization stations.

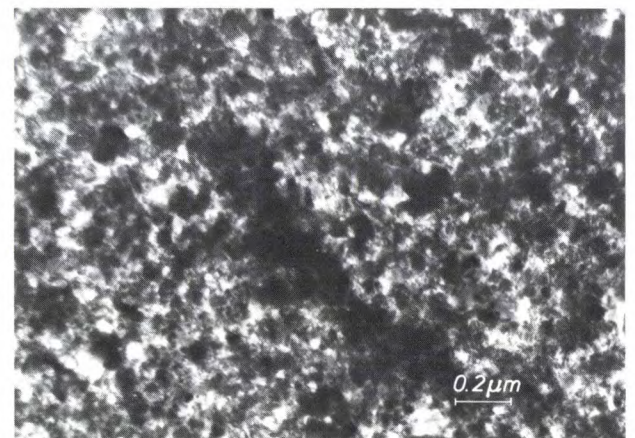


Fig. 12. Transmission-electron-microscope (TEM) picture of a layer obtained by electroless silver deposition on a cured 2p lacquer with a LaserVision pattern. The walls of the pits are also uniformly coated with silver.

ger than one pit, and shows the small crystallites that form the coating. It can be seen that the coverage is uniform, even on the walls of the pits.

A second alternative metallization technique is the method of 'magnetron sputtering' from a cold solid target by means of bombardment with ions (usually argon ions)^[10]. At an argon pressure of about 10^{-1} Pa and a sufficiently high direct voltage, a discharge is formed between a target, acting as cathode, and an anode. With a magnetic field, maintained by permanent magnets behind the cathode, a concentrated discharge plasma is produced immediately above the target surface. Ar^+ ions at an energy of 300 to 500 eV can be extracted from the plasma, and these ions bombard the target surface and sputter its material. The disc to be coated is situated directly opposite the target and outside the plasma region.

Fig. 13 is a diagram of the continuous sputtering machine specially designed for mass production. The design makes allowance for the high water desorption of PMMA. A photograph of the machine is given in fig. 14. The production rate is determined mainly by the sputtering time, and not by the time taken for loading and unloading, so that the metallization is again much more efficient than with evaporation. It also allows more freedom in the choice of the mirror material. Binary and ternary alloys can be applied, some of which have desirable properties for Laser-Vision applications, including high stability and low internal stresses. Another advantage is the generally better adhesion to the cured 2p lacquer. It is even possible to apply silver coatings that adhere well to the 2p lacquer when certain gases are added to the argon plasma.

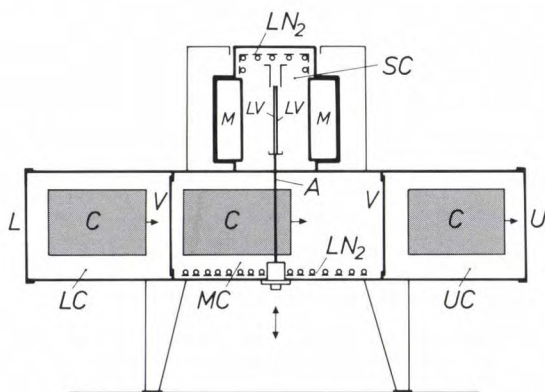


Fig. 13. Magnetron sputtering machine for mass production. *L* loading door. *LC* loading chamber. *C* cassette for discs. *V* valve. *MC* magazine chamber. *LN₂* cryopump surface cooled by liquid nitrogen for pumping water vapour. *M* magnetron. *LV* LaserVision disc. *UC* unloading chamber. *A* lifting arm. *SC* sputtering chamber. *U* unloading door. The vacuum pumps for the four chambers are not shown.

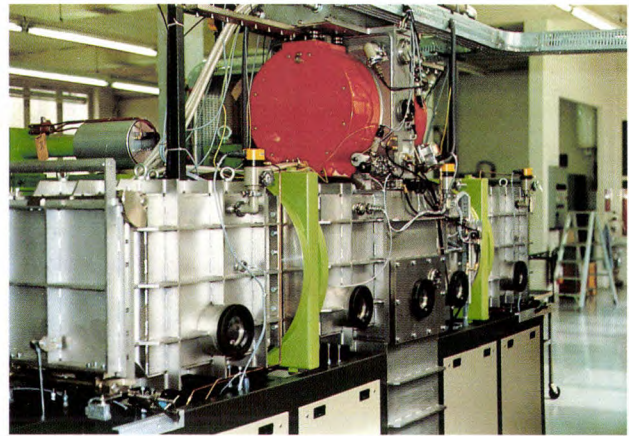


Fig. 14. Photograph of magnetron sputtering machine with the chambers indicated in fig. 13. The magazine chamber is in the centre, between the two green valve housings; the sputtering chamber is above it, with the red enclosures for the two magnetrons.

Finishing operations

After the metallization the sides of the disc are ready for use. At this stage the electro-optical characteristics and the signal quality of the separate sides can be determined, if necessary. First the plastic protective layer is sprayed on to the mirror coating. This layer is also the base for the adhesive used for bonding the two discs back-to-back to form the final double-sided video disc. Both sides are then sprayed with an acrylic-based dispersion contact adhesive, which is left to dry until a rough 'orange-peel' texture forms. This microstructure is necessary to allow the escape of air trapped during assembly, which might introduce local deformations and hence unacceptably large vertical accelerations during playback. The danger of entrapping air can also be avoided by assembling at low pressure.

After a pressing operation over the entire surface, the inside and outside edges of the double disc are again firmly pressed together to ensure a reliable seal, mainly to prevent the penetration of moisture, which could affect the back of the mirror coating. If necessary the disc is balanced by turning the outer edge eccentrically. Finally the labels are applied to both sides, and the disc is inspected, tested and packed.

Summary. The mass production of the video discs for the Philips LaserVision system, now on the market in Europe, is carried out at Blackburn, England. The discs are manufactured in a photopolymerization process: a liquid polymerizable lacquer is cured on a mould by exposure to ultraviolet light, so that the picture and sound information of the mould is imprinted in the lacquer. A complete video disc is formed from two component discs bonded together; each side consists of a transparent substrate, a primer layer, a cured lacquer layer containing the information, a metal mirror coating and a protective layer. The materials chosen for these layers and the methods used in applying them permit video discs to be produced that comfortably satisfy the exacting requirements of the LaserVision system.

^[10] See for example J. J. Scheer and J. Visser, Philips tech. Rev. 39, 246, 1980.

Photopolymerizable lacquers for LaserVision video discs

J. G. Kloosterboer, G. J. M. Lippits and H. C. Meinders

Introduction

As described in the previous article [1], the manufacture of LaserVision discs by the photopolymerization (2p) process places some very exacting requirements on the lacquer in which the information is stored. Some of these requirements, such as a high curing rate, low viscosity and good adhesion to the substrate, relate to the application in the 2p process, while others are connected with the specifications of the discs, e.g. the dimensional stability, durability and lack of odour. In addition, production technicalities such as availability, purity, constancy of composition and price are also important in the ultimate choice of lacquer.

A lacquer that meets all the requirements is not easy to find, largely because some of the requirements are more or less conflicting. For example, it must be easy to release a cured lacquer from a metal mould (sometimes called a 'stamper'), whereas the lacquer must adhere well to the metallic mirror coating applied later. There are two ways of achieving an acceptable compromise. One way is to look for a mixture of different monomers, each providing one or more of the desired properties. This entails an investigation of the dependence of the process and product properties on the composition. Another way is to make a careful study of the polymerization process and the relation between molecular structure and chemical behaviour, and in this way to try and find a monomer that combines as many desirable properties as possible.

In the investigation described in this article both approaches were adopted, the first mostly at the Centre for Metallurgical Chemistry and Lacquers (CMCL) of the Philips Plastics and Metalware Factories (PMF), the second mostly at Philips Research Laboratories. Both approaches have led to the manu-

facture of playable discs of good quality. It is not yet appropriate to draw any conclusion as to which lacquer is the best: process development is still in full progress and a small change in the 2p process may in turn necessitate the use of a different lacquer composition. Nevertheless, the information obtained already enables us to anticipate new developments.

In this article we shall first give a short general description of photopolymerization. We shall consider the mechanism, the choice of substances to be used and the methods of monitoring the polymerization reaction. We shall then discuss the search for suitable lacquers made from mixtures of monomers, and finally the search for suitable single-component lacquers.

Photopolymerization

In polymerization a large number of small molecules, monomers, are linked together to form chains of macromolecules [2]. The linking can take place by condensation or by addition. In condensation the formation of the macromolecules is generally accompanied by the formation of one or more by-products, e.g. water. During the reaction, dimers are first formed, then trimers, tetramers, and so on, which can also react with one another. It is only towards the end of the reaction that this step-by-step process leads to the formation of very long polymer chains of high molecular weight. Addition polymerization, on the other hand, proceeds in the form of a chain reaction, in which a chain, once started, keeps on growing through the addition of one monomer after another. No by-products are formed in this reaction. To start the chain growing, molecules with an unpaired electron (radicals) can be used. These can arise from the dissociation of an activating substance, called the initiator. When this dissociation is produced by the action of visible or ultraviolet light, the process is

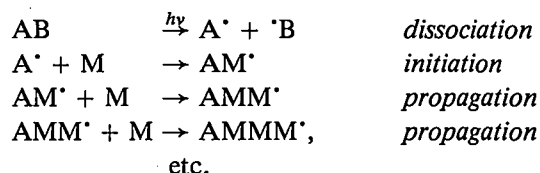
Dr J. G. Kloosterboer and Drs G. J. M. Lippits are with Philips Research Laboratories, Eindhoven; Dr H. C. Meinders is with the Philips Plastics and Metalware Factories (PMF), Eindhoven.

referred to as photopolymerization^[3]; when the dissociation is produced by heating the process is called thermal polymerization.

For replication processes, addition reactions are clearly preferable to condensation reactions, since no by-products are formed and, if a photosensitive initiator is used, the reactions are easy to control externally, since the polymerization does not start until light is absorbed by the initiator. The addition reactions controlled by light are also extremely efficient: the absorption of a single photon by the initiator can trigger the growth of a chain of 1000 to 10000 monomer links. Another important advantage is that the solutions of photo-initiators in monomers used for photopolymerization are stable in the dark, unlike solutions of thermal initiators in monomers. The photosensitive solutions can therefore be used without any danger of spontaneous reaction and explosion during transport and storage.

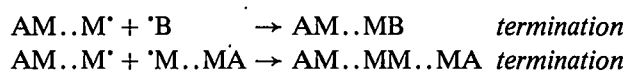
The mechanism

The photopolymerization of a monomer M, started by the initiator AB that supplies the radicals A' and 'B can be represented schematically as follows:



Once the radical AM' has been formed in the initiation process, the reaction proceeds by itself. If the original material contains different monomers, these can be incorporated in the same chain; this is referred to as 'copolymerization'.

In theory a single radical A' should be sufficient to make all the monomer molecules react together to form an enormously long single chain. This does not happen, because reactive radical end-groups are deactivated by termination, for example when two radicals meet:

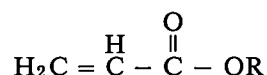


The two unpaired electrons thus form a covalent bond and a non-reactive polymer molecule is produced. Moreover, the growth can also be inhibited or even completely suppressed. This happens, for instance, in the presence of dissolved oxygen, because this reacts with the radicals and the resulting peroxy radicals are not highly reactive. After the oxygen and other inhibitors that may be present have been converted, the rate-determining step of the polymerization process is the initiation. Propagation takes place almost of its

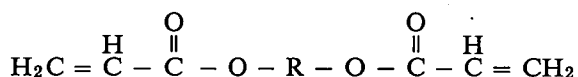
own accord: the process has a low activation energy and a relatively large amount of heat is released (up to 800 J/g). After initiation, a chain can therefore grow very fast until the point of termination is reached. The length a chain can grow to is determined by the probability of termination. This probability depends on the concentration of radicals and hence on the intensity of the light. A high intensity gives rise to the formation of relatively short chains.

Choice of monomers and photo-initiators

For a monomer, or class of monomers, to be used in the 2p process it is necessary to satisfy requirements with regard to rate of polymerization, separation from the mould, adhesion of the mirror coating, dimensional stability, etc. If the process is to be efficient, the polymerization rate must be high. Fast polymerization reactions are particularly easy to obtain with acrylates, which are esters of acrylic acid of the general formula:



where R is an arbitrary group that does not significantly affect the mechanism of the polymerization. This is a justification for insisting on the speed requirements: without too great a sacrifice of polymerization rate it is possible, with a suitable choice of R, to meet the other requirements of the 2p process as well. Thus, a lacquer can acquire high dimensional stability when R has an acrylate group on each side:



Molecules with such a structure will introduce cross-links between different chains. This cross-linking leads to the formation of a three-dimensional network.

The photo-initiator decomposes into radicals when it is irradiated by light in a particular absorption band of the molecule. For efficient initiation the location of this band should match the emission spectrum of the light source as closely as possible. Both should also be situated in the wavelength range above 300 nm, in

[1] H. C. Haverkorn van Rijsewijk, P. E. J. Legierse and G. E. Thomas, this issue, p. 287.

[2] P. J. Flory, Principles of polymer chemistry, Cornell University Press, Ithaca, New York, 1953.
G. Odian, Principles of Polymerization, 2nd edn, Wiley, New York 1981.

[3] In a narrower sense this term is used only when a photon is necessary for every link of a monomer. Chain reactions that are only initiated by light are then referred to as light-induced reactions. In this article, however, we use the term photopolymerization in the broad sense, that is to say when light is necessary in one way or another for the polymerization process.

which the acrylate monomers themselves have extremely low absorption. On the other hand, to produce a colourless product the initiator should absorb the minimum amount of visible light (400-700 nm). Some photosensitive aromatic ketones have an absorption maximum at about 350 nm and therefore meet these requirements. These compounds are somewhat sensitive to the blue-violet part of the visible spectrum, but in a 'yellow room' the lacquers in which they are dissolved can safely be handled without undesired polymerization. An example of a suitable initiator is 1,1-dimethoxy-1-phenylacetophenone (DMPA). On irradiation with light at about 350 nm one of the products is the benzoyl radical (*fig. 1*), which is highly reactive and an efficient initiator for the polymerization of acrylates.

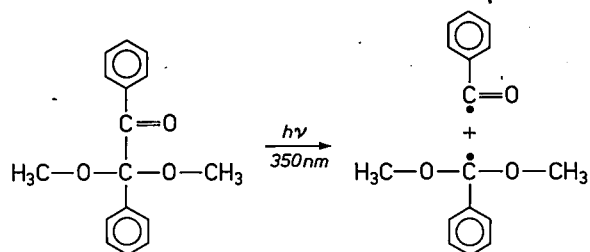


Fig. 1. Decomposition of the photo-initiator 1,1-dimethoxy-1-phenylacetophenone (DMPA) into a benzoyl radical and a dimethoxybenzyl radical by irradiation with ultraviolet light at a wavelength of 350 nm.

Monitoring the reaction

The progress of a polymerization reaction can be monitored by a number of methods: see *Table I*. Classical methods, such as measuring the volume shrinkage (dilatometry), determining the increase in viscosity and weighing the deposited polymer provide information about a macroscopic property of the material. This is also true for the more recent method of gel-permeation chromatography, in which particles of different molecular weight can be separated and the molecular-weight distribution determined. Most of these methods are not very suitable for monitoring a reaction through to a high degree of conversion, especially when there is considerable cross-linking. In such a case a method such as microcalorimetry, in which the heat of reaction is determined as a function of time, is more suitable. With the spectroscopic techniques mentioned in *Table I* it is possible to make measurements on a molecular scale. The concentration of unreacted C=C bonds can be determined by measuring the ultraviolet or infrared absorption and the Raman emission (inelastic light scattering) due to the monomer and polymer molecules. With this

method the progress of a polymerization reaction can also be monitored right through, even when there is cross-linking.

In our investigation Raman spectroscopy and microcalorimetry have been very useful. For fast routine measurements of reactivity a modern version of dilatometry has been particularly appropriate. A thin layer of liquid lacquer is introduced between two glass plates and then irradiated. By measuring the reduction in thickness of this sandwich as a function of time a measure of the reaction rate is obtained. In this way reactions with shrinkage half-value times of less than one second can be monitored. As an example we shall now briefly discuss the method of monitoring a reaction by means of microcalorimetry.

During the polymerization of acrylates a relatively large amount of heat is generated. When one type of reaction (in our case propagation) is dominant, the heat production is proportional to the polymerization rate, i.e. the number of monomers converted in unit time. This also applies to mixtures of acrylates, since all acrylate molecules give virtually the same heat of reaction per acrylate group. The polymerization rate v is proportional to the concentration of radicals and monomers. In the steady state the radical concentration $[R^*] = \sum_i [AM_i]$ can be assumed to be constant: in unit time the dissociation of the initiator produces about as many radicals as disappear by termination. The reaction rate is then proportional to the monomer concentration $[M]$ and will decrease exponentially with time: $v = -d[M]/dt \propto [R^*][M]$.

In a dilute solution of a monomer in an inert solvent this relation is indeed found. At higher monomer concentrations, however, a complication occurs, because the viscosity shows a marked increase during the reaction. This seriously inhibits the diffusion of large macroradicals, while the diffusion of the much smaller monomer molecules is hardly affected. Termination, generally a reaction between two large

Table I. Some methods for monitoring a polymerization reaction.

Method	Measurement of:
Dilatometry	Volume shrinkage
Viscosimetry	Increase of viscosity
Gravimetry	Weight of polymer after separation from monomer by precipitation
Gel-permeation chromatography	Molecular weight distribution of formed polymer
Microcalorimetry	Heat of reaction
Refractometry	Increase in refractive index
C=C spectroscopy	UV absorption (layers to 1 μ m) IR absorption (layers to 10 μ m) Raman emission (inelastic light scattering)

macroradicals, is therefore much more strongly suppressed than propagation, which requires a collision between a large macroradical and a small monomer molecule. In this way, for each chain initiated there will be many more monomer molecules converted, and for a constant initiation rate the reaction rate will therefore increase steeply. This effect is called auto-acceleration and is also known as the Trommsdorff effect [2].

To illustrate this effect *fig. 2* gives a plot of the relative reaction rate against time for the formation of polymethyl methacrylate, also known as Plexiglas or Perspex. After about 100 minutes, with a conversion of about 20%, a very marked auto-acceleration peak occurs: Acrylates give such a peak in the rate right at the start of the reaction. Some acrylates, e.g. ethyl acrylate, even have two different acceleration peaks (*fig. 2*); the explanation of this effect does not come within the scope of this article [4].

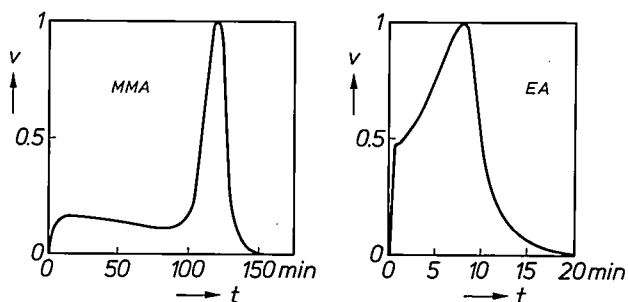


Fig. 2. Reaction rate v (in arbitrary units with 1 as maximum) as function of the time t , for photopolymerization of methyl methacrylate (MMA) and ethyl acrylate (EA) at 20 °C.

Copolymerization

Very many polymers used in practice are copolymers, which result from polymerization of different monomers to form a single type of chain or a single network. When a mixture of monomers is used and its composition is taken as a variable, considerable freedom of choice is obtained for certain properties of the polymer, making it possible to tailor a product to a particular application. This is of great interest when a large number of requirements have to be met at the same time, as in the video disc.

The reaction of a mixture of monomers is often complicated. In copolymerization two different monomers M_1 and M_2 are not incorporated in the network at random if they do not have exactly the same chemical reactivity to the radicals M_1^{\cdot} and M_2^{\cdot} . In the extreme case where M_1^{\cdot} and M_2^{\cdot} have a distinct preference for M_1 and M_2 , respectively, 'block' copolymers are formed in which there are small regions of poly- M_1 and poly- M_2 . On the other hand, if M_1^{\cdot} has a distinct preference for M_2 and M_2^{\cdot} for M_1 , then 'alternating'

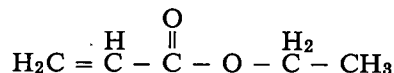
copolymers are formed. The reactivity ratios are known only for a few monomer pairs, and are difficult to determine in cross-linking systems. When the monomers react at different rates, the composition of the as yet uncured part of the mixture will also change during the reaction, as will also the composition of the polymer formed later. This requires accurate optimization of the mixture. To obtain a system that could be studied more easily, we also investigated the possibility of making good video discs with a single monomer, by means of 'homopolymerization'. During the polymerization of such a single-component lacquer there are of course no differences in reactivity and the composition of the polymer formed at different times is constant.

The search for a suitable mixture

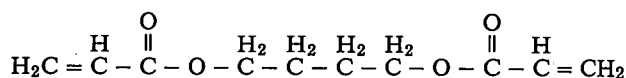
Because of the many possible variations of the R group there are many acrylates and hence very many possible different mixtures. The actual number, however, is considerably reduced when account is taken of the requirements for dimensional stability, shrinkage, adhesion and hardness. Our second selection was based on the required viscosity and rate of polymerization. The possible variations in R that we considered related to factors such as the length, and hence the flexibility of the molecule, and the number of C=C bonds and polar groups.

Selection based on dimensional stability, shrinkage, adhesion and hardness

One of the quantities that affects the dimensional stability of an acrylic polymer is the functionality f of the monomer. The value of f is equal to the number of molecules with which a monomer molecule can form direct bonds. A monomer with low functionality ($f = 2$), is ethyl acrylate:

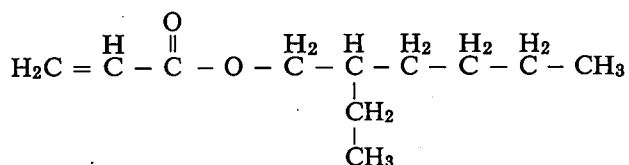


This monomer gives a rather unstable, rubbery polymer with long flexible chains and few cross-links. A monomer with a higher functionality gives more cross-links on polymerization, resulting in increased dimensional stability. An example of this is 1,4-butanediol diacrylate (BDDA), a monomer that can be considered to be built up from two molecules of ethyl acrylate:



[4] H. J. L. Bressers and J. G. Kloosterboer, *Polymer Bull.* 2, 201, 1980.
J. G. Kloosterboer and H. J. L. Bressers, *Polymer Bull.* 2, 205, 1980.

BDDA not only has a higher functionality ($f = 4$), but because the molecule is relatively small, it also has a large number of C=C bonds per unit volume. This results in a dense network, and hence in a stiff, not very flexible material. In addition, there is strong polymerization shrinkage, amounting to about 15% of the volume. A lacquer layer of thickness 50 μm , for example, will cause warp in a 1 mm plastic substrate. If the substrate is made of glass, and of the same thickness, the lacquer layer will crack. The flexibility can be substantially increased by adding a bulky 'auxiliary' monomer with low shrinkage, e.g. 2-ethylhexyl acrylate (EHA):



Besides the concentration of C=C bonds, it is also possible to consider the fraction of saturated hydrocarbon groups k , which is a measure of the groups not taking part in the reaction as a percentage of the total volume of the lacquer.

With the aid of f and k the available monomers can be mapped as in *fig. 3a*. On the right-hand side of the diagram are the soft, paraffin-like monomers with $\text{R} = \text{CH}_3(\text{CH}_2)_n -$, and on the left are BDDA, etc. Above are a few triacrylates ($f = 6$) and one tetraacrylate ($f = 8$), which, because of the low flexibility of their polymers, are even less useful in the pure form than BDDA. Tri- and tetra-acrylates with long chains of bulky groups do exist, but they are not available in a sufficiently pure state and are therefore not included in the figure.

The values of f and k of a few mixtures, and the monomers from which they are built up, are given in *fig. 3b*. This figure also indicates whether the cured lacquer can easily be released from the mould and whether there is good adhesion to the substrate. If the disc is to be released easily the cured product must contain no strongly polar groups such as $-\text{OH}$,

$-\text{SH}$, and $-\overset{\text{O}}{\parallel}{\text{C}}-\text{OH}$, because these groups can cause strong adhesion to certain metals. They can also produce corrosion of the mirror coating that will be applied later. Some mixtures, especially those of monomers with a long aliphatic chain, cannot be used because they do not adhere well even to a pretreated substrate.

A further constraint arises from the need to take account of the required hardness of the polymer; see *fig. 3c*. When a mirror coating is deposited by evap-

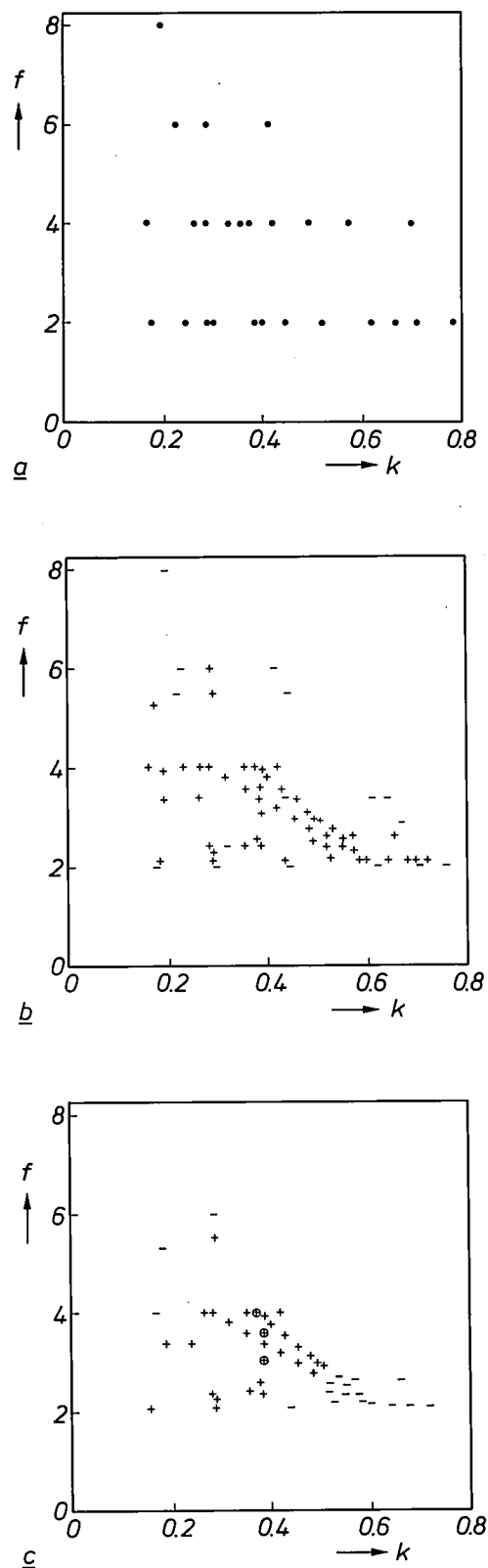


Fig. 3. Saturated hydrocarbon fraction k and functionality f of commercially available acrylate monomers (*a*), and for a number of mixtures composed from them (*b* and *c*). In *b* there is an indication of whether, after curing, good adhesion is obtained at the substrate and whether the layer separates readily from the mould; + yes, - no. In *c* there is a further selection, depending on whether the material remains bright on metallization; + yes, - no. The indications for two mixtures and a single-component layer that will give good video discs are ringed.

oration or sputtering, stresses occur in the metal layer, and on a soft lacquer these could cause shear and deformation. This shows up as a wave pattern on the surface of the lacquer (fig. 4) and produces a

milky discolouration. A good impression of the hardness of a cured lacquer is obtained by measuring the depth of indentation produced by a diamond point under a standard load. In fig. 5 the indentation depth in the polymer is plotted against the BDDA content for mixtures of EHA and BDDA. The hardness increases with the amount of BDDA in the mixtures. In practice it is generally found that no problems are encountered for indentations of less than 3 to 4 μm , and that, with carefully controlled metallization conditions, good discs can also be made with even softer lacquers.

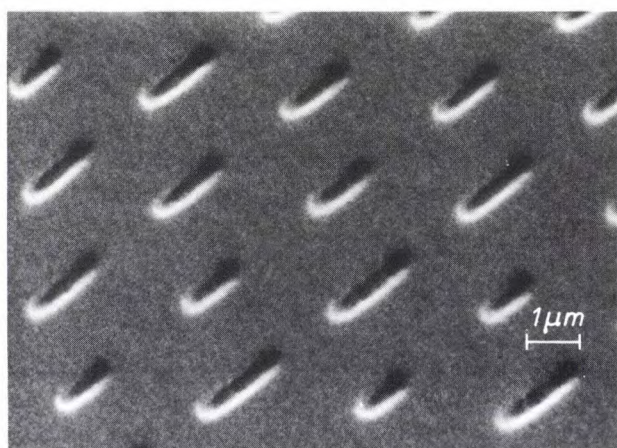
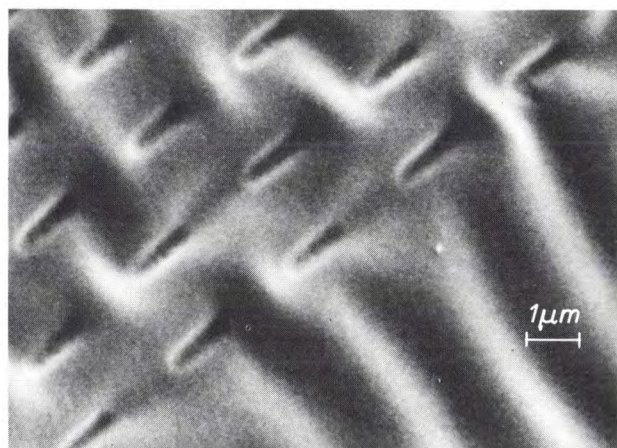
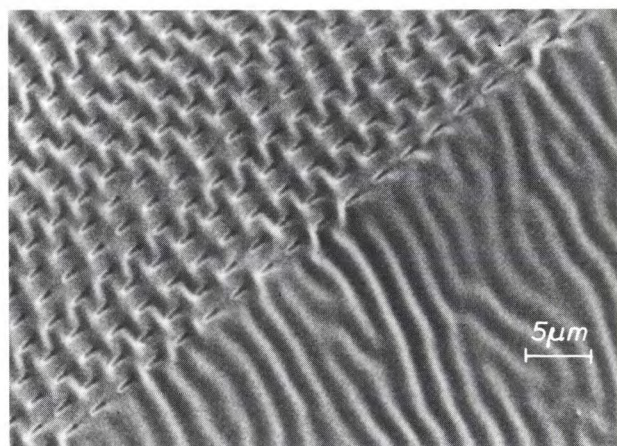


Fig. 4. Scanning-electron-microscope (SEM) pictures of the Laser-Vision pattern after metal deposition on a lacquer of a mixture of 2-ethylhexyl acrylate (EHA) and 1,4-butanediol diacrylate (BDDA). The upper photograph (magnification 2500 \times) and the photograph in the centre (magnification 10000 \times) were obtained with the composition 80% EHA, 20% BDDA. After metallization the soft lacquer appears as a wave pattern. The lower photograph (magnification 10000 \times) was obtained with the composition 20% EHA, 80% BDDA. Since the lacquer substrate is sufficiently hard, there is no wave pattern.

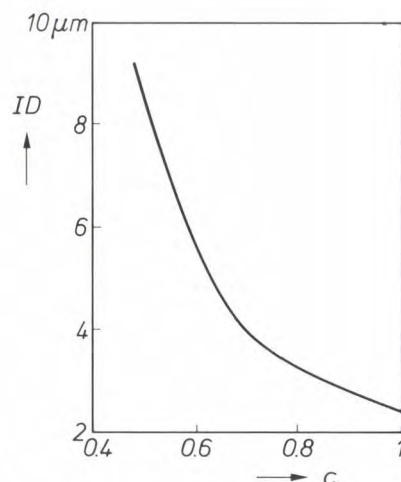


Fig. 5. Indentation depth ID of a standard diamond point in cured mixtures of EHA and BDDA, as a function of the BDDA content c . When c is increased the indentation depth decreases as a result of an increase in the hardness of the cured lacquer.

Selection based on viscosity and reactivity

For the 2p process the viscosity of the still-liquid layer should preferably be between 5 and 15 mPa.s. If the viscosity is too low, the lacquer does not spread out uniformly. If the viscosity is higher than 15 mPa.s, air bubbles can easily be enclosed at high production rates. Furthermore, with the method employed, it is difficult to make polymer layers that are uniformly thick from viscous lacquers.

To achieve a high production rate the reactivity of the monomer must be as high as possible. Table II gives the viscosity and the relative reactivity of a number of potentially useful monomers, and also the indentation depth of a diamond point in the polymer.

Because of their high viscosity, triacrylates cannot be used as the principal component of a mixture. On the other hand there are various diacrylates that do have a suitable viscosity. Of these, tripropylene glycol diacrylate (TPGDA) has the greatest reactivity. Some mixtures of TPGDA and monomers with one C=C bond appear to be even more reactive, as do mixtures

Table II. Properties of some single-component lacquers ^[a].

Monomer	Abbreviation	Viscosity at 23 °C (mPa.s)	Relative reactivity ^[b]	Indentation depth ^[c] (μm)
<i>diacrylates:</i>				
diethylene glycol diacrylate	DEGDA	7.8	1.06	1.9
polyethylene glycol diacrylate	PEGDA	24.3	1.00	10
dipropylene glycol diacrylate		8.6	2.10	2.3
tripropylene glycol diacrylate	TPGDA	12.3	3.00	3.4
hexanediol diacrylate	HDDA	6.7	2.98	2.8
3-methylpentanediol diacrylate		6.6	1.33	3.3
Ebecryl 180 (UCB)		25.8	1.54	2.4
4160 (Diamond Shamrock)		17.0	3.07	3.5
<i>triacrylates:</i>				
trimethylolpropane triacrylate	TMPTA	107	3.60	1.5
4094 (Diamond Shamrock)		172	3.44	2.3

^[a] 96% monomer + 4% photo-initiator, irradiated with TL05 lamp.

^[b] Determined by shrinkage measurements. The value for PEGDA is taken equal to 1.00.

^[c] For a diamond point in the cured lacquer. A small penetration depth corresponds to a high hardness of the lacquer.

of TPGDA and some triacrylates, e.g. trimethylolpropane triacrylate (TMPTA); see *Table III*.

The reactivity of a mixture of TPGDA and TMPTA can be further increased by the addition of N-vinylpyrrolidone (NVP); see *Table IV*. In this case there is obviously a synergetic effect, because pure NVP polymerizes so slowly that the progress of the reaction cannot be monitored by the conventional shrinkage measurement. The cause of the accelerating action of NVP is not yet entirely clear. It could be that the effect is related to the ring structure of NVP, since N-methylpyrrolidone also causes an accelerated reaction. Dimethyl acetamide, which also has a tertiary nitrogen atom next to a carbonyl group, but no ring structure, has the opposite effect of reducing the reactivity, probably through dilution of the lacquer. Of course, the two other substances also function as diluents, but evidently the accelerating effect predominates.

Some results

Various mixtures of TPGDA, TMPTA and NVP meet all the requirements of the 2p process and are suitable for use in the mass production of video discs. After curing, the lacquers do not give off toxic vapours ^[6] and they are clear enough to permit good optical read-out. The components are also available in a fairly constant composition. Properties of two mixtures are presented in *Table V*, which also includes the properties of a suitable single-component lacquer for comparison.

A minor problem that arises during the photolysis of DMPA is that coloured photoproducts can be formed. These give a yellowish tint to thick polymer layers. This yellowing is to some extent reversible: it decreases slightly in the dark and returns again in the light. In some mixtures the discolouration can also be seen in thin layers like those used for video discs. Although the playing of a disc is in no way affected by this discolouration (the discs are read out by red light), it looks slightly less attractive. The yellow discolouration can be avoided by using a different photo-initiator, but unfortunately this photo-initiator is less efficient.

The investigation of mixtures showed that the desired combination of properties can be obtained by using an approximate mix of components. Under the usual conditions, the mixtures mentioned here polymerize faster than the individual components investigated previously (*Tables II and V*).

The search for the right single-component lacquer

From the acrylates available we first made a selection based on the possible degrees of conversion and cross-linking, which determine the dimensional stability and hardness of the polymer. For a further selection we took the polymerization rate as the criterion. We continued to pay attention to aspects such as viscosity, availability, purity and price of the monomers, of course.

Selection based on degrees of conversion and cross-linking

Monoacrylates are not suitable for a single-component lacquer, because they have only one C=C bond per monomer molecule and therefore give far too little cross-linking. Polymers of monoacrylates are generally therefore soft and rubbery.

During the polymerization of di-, tri- and tetraacrylates not all the C=C bonds are equally reactive. If one of the C=C bonds of a monomer molecule is included in a chain, the mobility of the other is greatly reduced. When more and more monomer molecules become attached to the network by one 'arm', the chance of all the C=C bonds being converted decreases rapidly. Nevertheless, an unconverted C=C bond remains a centre of high chemical reactivity. A high degree of conversion is therefore essential if a stable product is to be obtained.

To define the region more precisely, we have correlated the degree of conversion of various acrylate monomers with the achievable degree of cross-linking, the concentration (mol/l) of the cross-links. Besides its dependence on the degree of conversion, the degree of cross-linking also depends on the molar volume of the monomer and on the details of the reaction mechanism, such as ring and ladder formation. To a first

Table III. Properties of some two-component mixtures with TPGDA as main component [a].

Second component	Viscosity at 23 °C (mPa.s)	Relative reactivity [b]	Indentation depth [c] (μm)
vinyl acetate	4.5	3.33	3.7
2-ethylhexyl acrylate (EHA)	7.1	2.14	>10
butoxyethyl acrylate	7.6	3.33	—
<i>t</i> -butyl cyclohexyl acrylate	11.7	1.28	7.3
phenylpropyl acrylate	10.0	1.50	9.3
ethoxyphenyl acrylate	16.2	3.07	>10
trimethylolpropane triacrylate (TMPTA) [d]	13.0	3.24	3.2
VPS - 2051 (Degussa) [e]	13.5	3.60	2.8
VPS - 2052 (Degussa) [e]	13.5	2.86	2.7

[a] 96% monomer mixture + 4% photo-initiator, exposed with TL05 lamp. The mixture contains 80% of TPGDA and 20% of the second component, except in the last three cases, where an 80/20 mixture would have too high a viscosity.

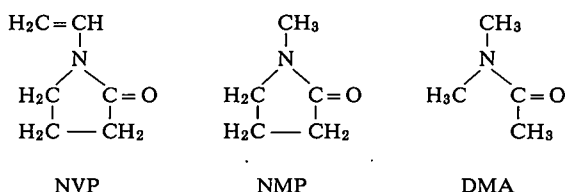
[b] See Table II, under [b].

[c] See Table II, under [c].

[d] 15% instead of 20%.

[e] 17.5% instead of 20%.

Table IV. Properties of a TPGDA/TMPTA mixture [a] after addition of N-vinylpyrrolidone (NVP), N-methylpyrrolidone (NMP) and dimethyl acetamide (DMA) with the following structures:



Addition [b]	Viscosity at 23 °C (mPa.s)	Relative reactivity [c]
—	13.0	3.24
NVP	7.8	6.68
NMP	10.7	4.29
DMA	6.9	2.70

[a] 96% monomer mixture + 4% photo-initiator, exposed with TL05 lamp. Without addition the mixture contains 85% of TPGDA and 15% of TMPTA.

[b] 25% of the TPGDA is replaced.

[c] See Table II, under [b].

Table V. Properties of lacquers that will give good video discs.

Lacquer	Composition	Viscosity at 23 °C (mPa.s)	Relative reactivity [a]	Indentation depth [b] (μm)
A	57% TPGDA 29% NVP 10% TMPTA 4% DMPA	7.8	6.68	1.9
B	61% TPGDA 17.5% NVP 17.5% TMPTA 4% DMPA	12.3	5.99	1.9
C	96% HDDA 4% DMPA	6.7	2.98	2.8

[a] See Table II, under [b].

[b] See Table II, under [c].

approximation, the degree of cross-linking N for a completely converted monomer is equal to

$$N = 1000 (f - 2) \rho / 2M,$$

where f is the functionality, ρ the density (kg/l) and M the molecular weight (g/mol) of the monomer. For the monomers of interest to us the maximum value of N is mainly determined by M : whereas ρ is always between 1.1 and 1.2 kg/l, M can vary by a factor of 3 to 4. Since ρ is usually not known exactly, we shall henceforth consider not N but N/ρ (mol/kg) as a function of the degree of conversion.

The approach to the maximum value of N/ρ depends on the reactivity of the 'free' C=C bonds (r_1) and on that of the 'dangling' C=C bonds (r_2). In fig. 6 the value of N/ρ for a diacrylate is plotted as a function of the degree of conversion x for three special cases. When r_1 is very much greater than r_2 , all the

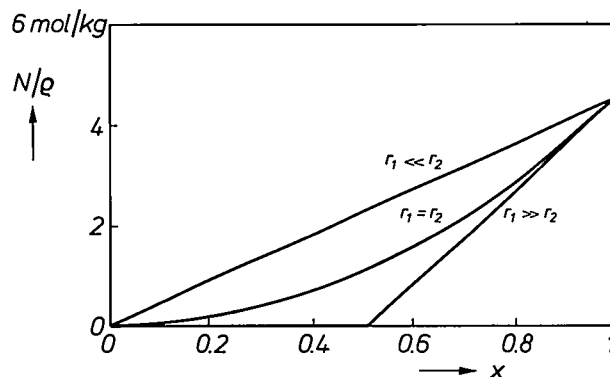


Fig. 6. Theoretical ratio of the degree of cross-linking N to the density ρ , plotted against the degree of conversion x of a diacrylate, for different relative reactivities of the 'free' C=C bonds (r_1) and the 'dangling' C=C bonds (r_2).

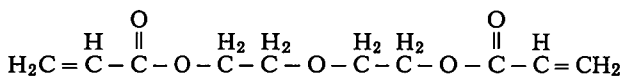
monomer molecules first react with one C=C bond and only then do the 'dangling' C=C bonds enter into the picture so that cross-links are formed. If r_1 , on the other hand, is very much smaller than r_2 , then each molecule that has reacted with one C=C bond will immediately react again with its other C=C bond. If r_1 and r_2 are of equal magnitude, the quantity N/ρ will increase as the square of x . Because of the reduced mobility after the reaction of the first C=C bond, we would expect that in practice the curve of N/ρ against x would lie between the lower two curves.

The actual reactivity ratio r_1/r_2 is not only unknown but may also depend on x . Unfortunately we are not yet able to monitor the development of the

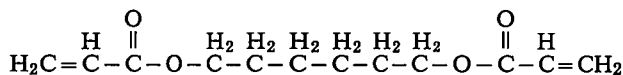
[5] All acrylate monomers are toxic. Differences in toxicity are primarily determined by differences in volatility. During the manufacture of the video discs the vapours are of course exhausted, so that the production personnel are not exposed to harmful concentrations.

cross-linking with the reaction time, nor is it yet possible to give any reasonably well-founded statistical description of the process [6]. Nevertheless, to obtain a simple comparison of various monomers, we always assume that r_1 is very much greater than r_2 , so that the lower curve is followed. For the tri- and tetraacrylates we assume that the reactivities of the third and fourth double bond are approximately equal to r_2 and are thus very much smaller than r_1 . Fig. 7 shows the variation of N/ρ with x for two di(meth)acrylates (one with a large molar volume and the other with a small molar volume), two tri(meth)acrylates (again with large and small molar volumes) and a tetraacrylate. The curves indicate that with tri- and tetraacrylates it should in theory be possible to obtain a substantially higher value of N/ρ than with diacrylates. This is not so in practice, however. We see this, for example, if we indicate on the curves the points based on the measurements of x taken from the literature [7]. Both at 30 °C and at 50 °C a higher value of N/ρ can be obtained with diacrylates, because x can become much higher. For the tri- and tetraacrylates the value of N/ρ does not become higher unless the temperature is gradually raised from 30 °C to 250 °C during the reaction, but obviously this condition cannot be applied to the manufacture of video discs.

The conclusion to be drawn is that it is preferable to start with a diacrylate. To define the choice more closely, fig. 8 gives a plot of N/ρ against x for a number of diacrylates. For clarity, only the theoretical values of N/ρ associated with the measured values of x are shown [7]. At the highest temperature diethylene glycol diacrylate (DEGDA) is superior. This compound has a long flexible connection between the acrylate groups:



The flexibility is mainly attributable to the $\overset{\text{H}_2}{\text{C}}-\overset{\text{H}_2}{\text{C}}-\text{O}-\overset{\text{H}_2}{\text{C}}-\overset{\text{H}_2}{\text{C}}$ fragment. This is because the energy barrier for rotation about O-C bonds is much smaller than for rotation about C-C bonds, where the hydrogen atoms tend to obstruct one another. A good second choice is 1,6-hexanediol diacrylate (HDDA), which has a somewhat longer but slightly less flexible connection:



At 30 °C HDDA has the highest degree of cross-linking. Another candidate indicated in fig. 8 is tetraethylene glycol diacrylate (TEGDA), which has the highest degree of conversion at 30 °C.

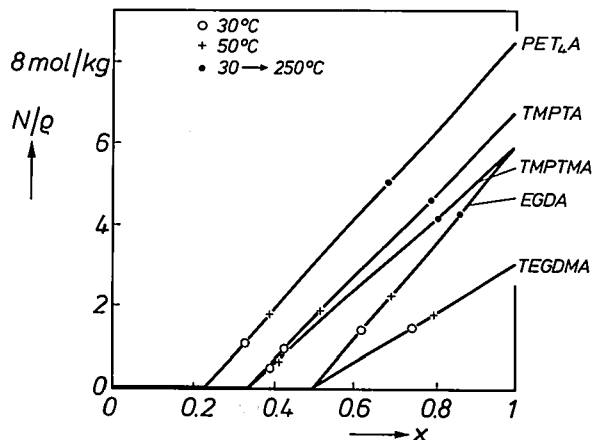


Fig. 7. N/ρ plotted against x for some di-, tri- and tetra(meth)acrylates, i.e. tetraethylene glycol dimethacrylate (TEGDMA), ethylene glycol diacrylate (EGDA), trimethylol propane trimethacrylate (TMPTMA), trimethylol propane triacrylate (TMPTA) and pentaerythritol tetraacrylate (PET₄A). The points on the curves relate to measurements of x reported in the literature at 30 °C, 50 °C and for a thermal scan from 30 to 250 °C.

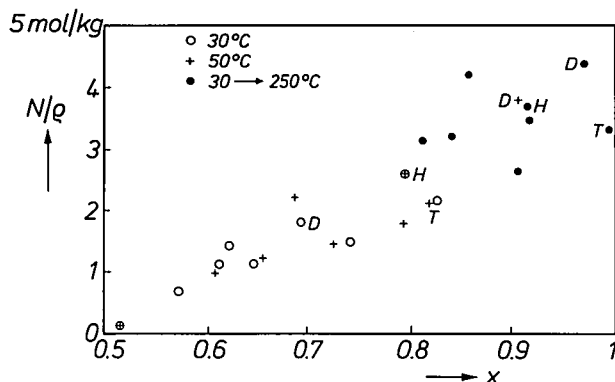


Fig. 8. N/ρ plotted against x for a number of diacrylates. The various points were obtained in the same way as those in fig. 7. The points for diethylene glycol diacrylate (DEGDA), hexanediol diacrylate (HDDA) and tetraethylene glycol diacrylate (TEGDA) are denoted by D , H and T , respectively.

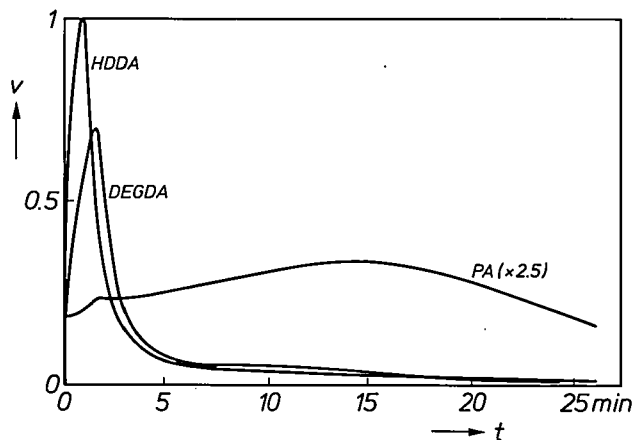


Fig. 9. Conversion rate v (in relative units) as a function of time t , for the polymerization of DEGDA, HDDA and n -propyl acrylate (PA).

Selection based on polymerization rate

As a consequence of the auto-acceleration the polymerization rate cannot be expressed by means of one or more rate constants, because both the propagation-rate constant and the termination-rate constant depend closely on the degree of conversion. These constants can only be determined in a dilute solution at a low degree of conversion, so that the polymer molecules formed do not 'see' each other. For our purposes such determinations are unfortunately not very useful. The reaction rates of diacrylate monomers can best be compared by determining the maximum rates. This comparison, however, is rather difficult because the reaction rate is highly sensitive to the presence of impurities, or inhibitors deliberately added to prevent spontaneous polymerization. We therefore first thoroughly purified the diacrylates selected for our experiments. To make careful measurements and comparisons of the maximum rates, we let the reactions take place much more slowly than in the actual manufacture of the discs. Fig. 9 shows the

DEGDA, HDDA and TEGDA are unfortunately not readily comparable with those in fig. 8). In addition we investigated the purity of the five monomers. The C=C bond content was determined by a complete bromination and the result was compared with the theoretical content at a purity of 100%.

Table VI gives the results of the various measurements. If we require the maximum conversion rate to be as high as possible, with the purity closely approaching 100%, then the only remaining candidates

Table VI. Degree of conversion x , maximum conversion rate v_m at 20 °C and purity z of some diacrylate monomers.

Monomer	x	v_m (%/s)	z (%)
DEGDA	0.71	0.43	94
HDDA	0.70	0.56	100
TPGDA	0.75	0.53	94
TrEGDA	0.76	0.23	99
TEGDA	0.82	0.57	98

Table VII. Maximum conversion rate v_m and the time t_m taken to reach it, for three HDDA-based lacquers and for the mixtures *A* and *B* mentioned in Table V under different curing conditions. For the polymerization in air the time t_i is given during which there is no reaction owing to the inhibiting effect of oxygen.

Conditions Lacquer	Nitrogen, Low initiation rate		Nitrogen, High initiation rate		Air, High initiation rate		
	v_m (%/s)	t_m (s)	v_m (%/s)	t_m (s)	v_m (%/s)	t_m (s)	t_i (s)
HDDA, No. 1 [a]	0.06	1180	5.4	14	2.7	51	19
HDDA, No. 2 [b]	0.49	50	9.4	5	5.2	25	11
HDDA, No. 3 [c]	0.57	50	10.6	6	6.5	25	11
<i>A</i>			12.4	5	11.3	11	5
<i>B</i>			11.4	5	10.2	13	5

[a] Contains 190 ppm of *p*-benzoquinone.

[b] Contains 15 ppm of *p*-benzoquinone.

[c] Obtained after purification of No. 2; contains less than 5 ppm of *p*-benzoquinone.

rate-time curves of DEGDA and HDDA. The maximum conversion rate of HDDA is clearly higher than that of DEGDA.

For comparison fig. 9 also gives the curve for *n*-propyl acrylate. This compound can be regarded as 'half' of HDDA: all the C=C bonds are located in separate molecules, whereas in HDDA they are arranged in pairs. This structural difference has a considerable effect on the maximum conversion rate and the shape of the curve.

We also measured the maximum conversion rate for three other diacrylates with a useful degree of conversion and cross-linking: TEGDA (fig. 8) and the related compounds triethylene glycol diacrylate (TrEGDA) and tripropylene glycol diacrylate (TPGDA). We also determined the degree of conversion at 20 °C (owing to a difference in measuring conditions the results for

are HDDA and TEGDA. Of the two, TEGDA has the higher degree of conversion. For the moment, however, our preference is for HDDA, not only because of its greater purity but also because of its greater availability. It is widely used, for example, in printing inks, coatings and lacquers, and is therefore widely obtainable in a reasonably pure form.

Some results

Table VII gives the maximum conversion rate v_m for three HDDA-based lacquers under different curing conditions. At a low initiation rate in a nitrogen

[6] R. S. Whitney and W. Burchard, *Makromol. Chemie* **181**, 869, 1980.

[7] J. E. Moore, in: S. S. Labana (ed.), *Chemistry and properties of crosslinked polymers*, Academic Press, New York 1977, p. 535.

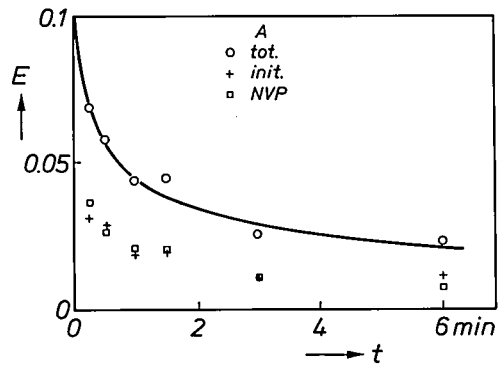
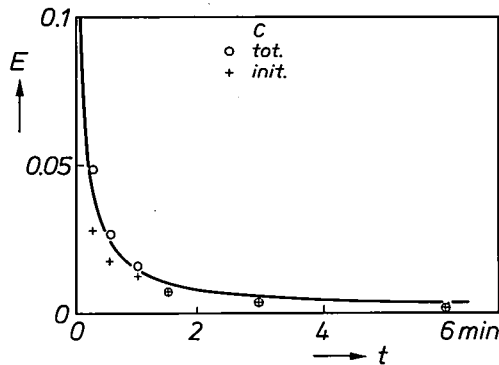
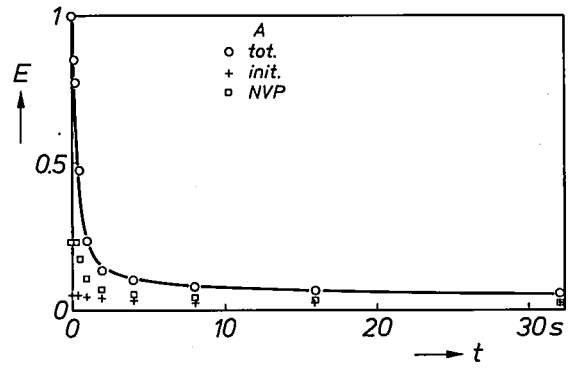
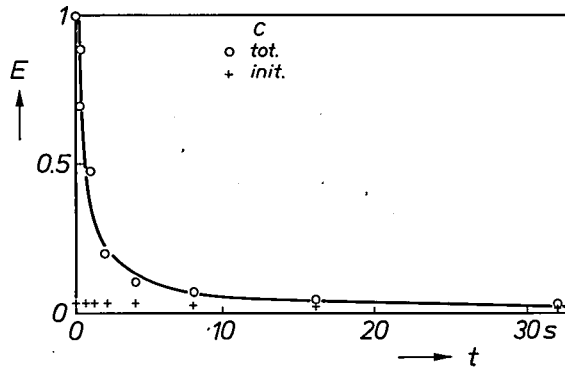


Fig. 10. Extraction curve of the single-component lacquer *C* consisting of HDDA and the initiator DMPA (Table V). The chromatographically determined total quantity of extracted material (E), expressed as the fraction of the initial weight of the lacquer, is plotted as a function of the exposure time t . The proportion of the initiator is indicated separately. Conditions: nitrogen atmosphere, four TL09 lamps, intensity at sample position 1.5 mW/cm^2 .

Fig. 11. Extraction curve (as in fig. 10) of the lacquer *A* consisting of a mixture of TPGDA, TMPTA, NVP and the initiator DMPA (Table V). The proportions of the initiator and NVP are indicated separately. The irradiation intensity at the sample position is 1.2 mW/cm^2 .

Table VIII. Some data on the separation of LaserVision discs from a mould, measured at a pulling rate of 500 mm/min , when mixtures *A* and *B* and the single-component lacquer *C* are used.

Lacquer	Total separation energy (J)	Deformation energy of the substrate (J)	Separation energy of lacquer (J)	Separation time (s)	Maximum pulling force (N)
<i>A</i>	1.27	0.39	0.88	4.2	51
<i>B</i>	0.75	0.27	0.48	3.6	48
<i>C</i>	0.30	0.08	0.22	2.4	31

atmosphere, v_m is closely dependent on the content of the inhibitor *p*-benzoquinone in the lacquer. This compound has a retarding effect, even in the absence of oxygen. If the initiation rate is increased by a factor of 200, the variation in v_m with inhibitor content decreases considerably.

At a relatively low intensity of incident light (0.2 mW/cm^2) these samples do not cure well in air. The photochemical consumption of dissolved oxygen cannot apparently compete with the replenishment of oxygen by diffusion. The video discs are also manufactured in a normal atmosphere, but since the

lacquer layer is enclosed between mould and substrate, oxygen cannot affect the lacquer^[8]. We simulated this situation by exposing the samples and the calorimeter to a flow of air, and then covering the samples with a thin film of mylar, which is transparent to ultraviolet radiation. The polymerization was successful, but there was inhibition, particularly in the sample with the most *p*-benzoquinone.

We also performed the same experiments on the mixtures *A* and *B* listed in Table V. In a nitrogen

[8] Some inhibition does occur at the edge, of course. The incompletely cured edge is removed later.

atmosphere there is little difference compared with the purest HDDA lacquer, but the inhibition on saturation with air is clearly reduced. This is mainly due to the presence of NVP. In separate experiments with mixtures of HDDA and NVP it was found that the influence of NVP depends on the presence of oxygen. In a nitrogen atmosphere a mixture of HDDA and NVP polymerizes more slowly than HDDA, but this order is reversed on saturation with air.

Another subject of investigation related to the small molecules still present in a fully cured lacquer. To obtain sufficient light absorption it is customary to use a surplus of initiator, so that in addition to unconverted monomer some photosensitive initiator remains, possibly accompanied by small quantities of its photoproducts. This may lead not only to the discoloration mentioned above, but also to an undesired continuation of the photochemical processes. Small molecules may also drift slowly to the surface where, if they are hygroscopic or chemically reactive, they may adversely affect the quality of the mirror coating. It therefore seems likely that a dense network offering little scope for transport will give a more durable disc.

To obtain some idea of the quantity and mobility of the small molecules in a polymerized lacquer, we determined the 'extraction' curves by liquid chromatography. The curve for lacquer C with HDDA (Table V) shows that with increasing exposure time the amount of extracted material rapidly decreases, see *fig. 10*. After some time little but initiator is extracted. *Fig. 11* shows that the extract of the NVP-rich mixture A contains (apart from initiator) mainly NVP, which is less rapidly incorporated in the network, since it has one

double bond. After a few minutes some stabilization occurs, and the quantity of extracted material hardly decreases at all with exposure time.

Good video discs can also be made with HDDA. After curing and metallization, the LaserVision pattern is indistinguishable from the pattern shown in the lower photograph in *fig. 4*. The lacquer is easily released from the mould. The energy required for releasing the disc plus lacquer from the mould (the total separation energy) is considerably lower than for the mixtures A and B; see *Table VIII*. This helps to prolong the life of the mould.

Since there is some warping of the disc during its separation from the mould, part of the separation energy must be attributed to the elastic deformation of the disc, so that the actual adhesion energy between lacquer and mould is lower than the total separation energy. The adhesion energies appear to be approximately proportional to the total separation energies.

Finally, we should mention that the video discs made with the mixtures and with HDDA satisfactorily pass all the normal acceptance tests.

Summary. Photopolymerizable lacquers for the manufacture of LaserVision video discs contain acrylate monomers and a photoinitiator that triggers the polymerization of the monomers on exposure to ultraviolet light. Two approaches were adopted to obtain lacquers that meet the many exacting requirements of LaserVision. In one approach a search was made for a mixture of monomers that would collectively provide the desired properties. In the other an effort was made to find a monomer that itself had all the desired properties. Both approaches led to the desired result: lacquers that give good LaserVision discs.

The fining of glass

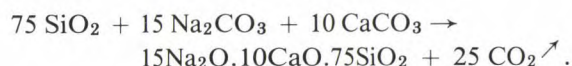
A Raman-spectrometric investigation into the action of arsenic oxides

H. Verweij

At Philips Research Laboratories laser-Raman spectrometry has been used to study the reactions occurring during glass formation. Particular attention has been paid to the action of small quantities of oxides of arsenic, which have traditionally been added to the melt for producing bubble-free glass in an economically acceptable time.

What is fining?

Soda-lime glass has been made for centuries by fusing together sand, soda ash and lime, in a reaction that can be approximated by



Many other familiar types of glass are produced by substituting potassium for the sodium (potash glass) or by using other metals instead of the calcium (to give lead or flint glass). In the glass for television tubes the calcium is largely replaced by barium, to give better X-ray absorption.

The temperature of the batch materials during glass formation is usually increased up to about 1450 °C [1]. At about 900 °C a large amount of CO₂ begins to form in the melt, so that the reacting mixture initially contains a very large number of bubbles (*fig. 1*). The formation of bubbles in the melt may be due not only to the release of CO₂ but also the presence of nitrogen, oxygen, water, etc. During the thermal treatment of the glass some of the bubbles rise to the surface, where they escape from the melt. This bubble formation is a useful feature of the glass-forming process, since the rising movement of the bubbles helps to produce a more homogeneous glass, but it has to be controlled in such a way that the final product is sufficiently free from bubbles.

In the manufacture of glass today, for example glass for television tubes or for optical applications, the specifications relating to the concentration and size of the bubbles remaining in the glass are very strict. This has led to an increased interest in the removal of bubbles from the glass melt, a process traditionally known as 'fining'.

In the fining process the size of the bubbles plays an important part: bubbles of diameter smaller than 10 µm are usually unstable and quickly dissolve again



Fig. 1. Photograph of a cross-section of a glass melt that has been heated to 1450 °C and then rapidly cooled to room temperature. The bubbles to be seen consist mainly of CO₂. To remove such bubbles rapidly, small quantities of a fining agent are added to the ingredients of the melt. The purpose of the research described here was to discover the mechanism of the fining process.

in the melt. Bubbles of diameter greater than 1 mm usually rise quickly to the surface (*Table I*) and therefore present no problems. But bubbles of diameter between these values require special measures for removing them from the melt at an economically acceptable rate.

Table I. The calculated mean time of rise of gas bubbles in a representative silicate glass at the melt temperature, as a function of the diameter of the bubbles.

Diameter (μm)	Rise time
10	1.6 years/m
100	5.8 days/m
1000	1.4 hours/m
10 000	0.8 min/m

The oldest and still most widely used method of fining is to add small quantities of chemicals, called fining agents, to the ingredients for the glass. In the production of window glass and bottle glass the chemicals used for this purpose are sulphates combined with carbon, and in the production of optical glass and glass for television picture tubes, nitrates combined with oxides of antimony are used. A combination of nitrates with oxides of arsenic is also very effective, but arsenic oxides are no longer used in the manufacturing process, largely because of their toxicity.

Most of our knowledge of the action of fining agents has been obtained empirically, and even today little is known with certainty about the mechanism of their operation. Since the production of picture tubes is of such great economic significance to Philips, we decided to undertake a closer study of this mechanism at the Research Laboratories. We began by looking into the action of arsenic oxides, as a preliminary to a more extensive investigation into the action of antimony oxides and the function of the nitrates^[2]. Although arsenic oxides are now little used in the manufacture of glass, the material is highly suitable as a model for a study of the fining action. For the analytical determinations involved in our investigation we have used laser-Raman spectrometry.

Why Raman spectrometry?

We wished to look more closely into the way the reactions take place in the glass melt, and in particular to find out how temperature affects the ratio of the concentrations of trivalent to pentavalent arsenic ions, $[\text{As}^{3+}]/[\text{As}^{5+}]$, in the formation and fining of glasses with an arsenic content. According to an existing hypothesis, this ratio increases with increasing temperature, while oxygen is released. It appeared that

the oxygen released caused the CO_2 bubbles formed earlier in the glass melt to swell and rise to the surface, and that this was the main cause of the fining action of arsenic. Unlike CO_2 bubbles, any remaining oxygen bubbles should dissolve again in the melt on cooling.

The change of $[\text{As}^{3+}]/[\text{As}^{5+}]$ as a function of temperature has usually been determined by 'wet' chemical methods. We wished, however, to obtain additional information that would give us a complete picture of the various reactions that take place during the glass-formation and fining processes. The conditions under which measurements have to be carried out for such an investigation — high temperature or, after sudden cooling to a low temperature, the mixed crystalline/vitreous state of the specimens — greatly restrict the methods of analysis that can be used.

The two methods most suitable for such an investigation are infrared and Raman spectrometry. Both methods provide information about molecular structures (including the valence of the ions that occur in these structures); it does not matter whether the ions are in a crystalline or a glassy environment, and in neither of the two methods is it necessary to destroy any part of the specimens, with the associated loss of information^[3].

In view of the small concentrations of fining agent (less than 1%) used in the fining, Raman spectrometry, which is intrinsically less sensitive, would seem to be less suitable than infrared spectrometry. With the latter method, however, there is generally too much overlap of the bands in the spectra, and since the introduction of lasers in recent years, the sensitivity of Raman spectrometry has also been substantially improved.

Experimental

The analytical measurements required in our experiments were performed with the laser-Raman spectrometer illustrated in *fig. 2*. Our aim in the choice of the instruments and in the performance of the measurements was to give the method of measurement the highest possible sensitivity and accuracy.

[1] A description of the most commonly used methods of *continuous* glass production is given in G. E. Rindone, *Glass Ind.* **38**, 489, 1957, and J. Staněk, *J. non-cryst. Solids* **26**, 158, 1977.

[2] A more extensive description of our research, with all the experimental details, will be found in H. Verweij, *Melting and fining of arsenic-containing silicate glass batches*, Thesis, Eindhoven 1980. This thesis includes H. Verweij, *J. Amer. Ceramic Soc.* **62**, 450, 1979, and H. Verweij, *J. Amer. Ceramic Soc.* **64**, 493, 1981.

[3] In an earlier investigation elsewhere, X-ray diffraction was used for the same purpose, but the information obtained was confined to crystalline compounds. Earlier thermogravimetric determinations and differential thermal analysis had the disadvantage that the data obtained were all indirect. Further details are given in the articles of note [2].

One of the methods of increasing the accuracy is to measure each spectrum a number of times and then take the mean. Since each separate 'scan' takes a considerable time, owing to the low intensity of the signal to be measured, the complete measurement may take hours. For this reason the measurement procedure has been largely automated.

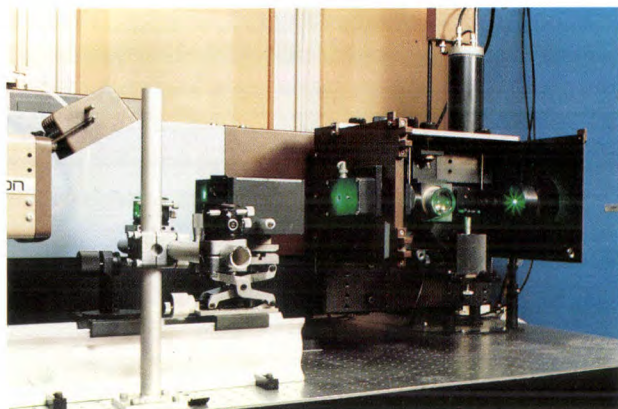


Fig. 2. View of the Raman spectrometer set up at Philips Research Laboratories. On the far left is a laser, which emits green light. On the right can be seen the lens that focuses some of the light scattered from the specimen on to the entrance slit of the monochromator situated behind it. The Raman spectrometer is a 'Ramanor HG2S' made by Jobin Yvon, Long Jumeau, France.

The experimental arrangement is shown in *fig. 3*. The primary beam, generated by a 4W Ar⁺ laser, *La*, which delivers a beam of highly monochromatic light of high energy density at 488.0 or 514.5 nm, is focused by a lens *Le*₁ into the specimen chamber *Sp*. Light scattered at right angles to the direction of the primary beam is then focused on to the entrance slit *En* of the monochromator *Mo*.

To eliminate undesired stray light as far as possible, successive spectral filtering is carried out in the monochromator by means of concave gratings. The control and computing unit *Com* controls the position of the gratings and of the various slits in the monochromator by means of stepping motors. At the exit *Ex* of the monochromator there is a photomultiplier tube *Ph*₁, which converts photons into current pulses; the photomultiplier has a very low dark current. After amplification and pulse-height selection in the amplifier/discriminator system *Am*₁, these current pulses are counted by a 100 MHz counter *Cou*.

To obtain the optimum correction for the drift in the laser intensity during the measurements, the measurement time of the counter (the 'time window') is adapted to the laser intensity in the following way. Part of the primary beam is conducted along a glass fibre to a second photomultiplier tube *Ph*₂. After similar amplification and pulse-height selection, in *Am*₂, the current pulses from this tube are applied to

the counter as external clock pulses. The time window of the counter is set to a fixed number of these external clock pulses, and its 'length' is therefore inversely proportional to the laser intensity — this has the result that the number of signal pulses counted does not vary with the laser intensity.

Another measure designed to increase the accuracy of the measurements is to subject all the measured counter values to a procedure for evaluating their statistical significance. Each time window is divided into twenty equal parts, and the twenty corresponding counter values are examined to see whether they deviate significantly from the mean or not. If they do, they are rejected and replaced by new counter values, and the evaluation procedure is repeated.

Since the spectra are available in digital form, a large number of operations can easily be carried out at a later stage. Examples are the production of difference spectra, obtaining spectra of single components from spectra of mixtures, the accurate determination of separate peaks, various statistical operations on the spectra, and the representation of the spectra in different ways.

The glass composition

The glass specimens to be investigated were prepared from a mixture of 70 mol% SiO₂ and 30 mol% K₂CO₃, with 1 mol% of As₂O₃ added as a fining agent. The composition was made as simple as possible to minimize the number of reaction products occurring during the glass formation, and to avoid phase separations in the melt at the high temperatures at which the glass forms. We used potassium instead

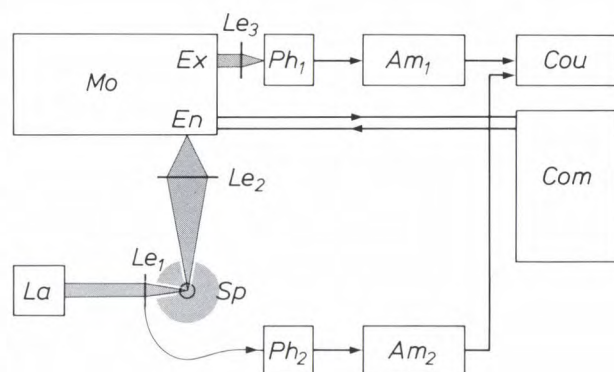


Fig. 3. Diagram of the laser-Raman spectrometer in *fig. 2*. *La* 4W-Ar⁺ laser. *Le*₁, *Le*₂, and *Le*₃ lenses. *Sp* specimen chamber. *En* and *Ex* entrance and exit slits of the monochromator *Mo*. *Ph*₁ and *Ph*₂ photomultiplier tubes. (*Ph*₂ receives light from the primary beam via a glass fibre.) *Am*₁ and *Am*₂ amplifier/discriminator systems. *Cou* 100 MHz counter. (The 'time window' of the counter closes after a fixed number of external clock pulses, which are applied via *Ph*₂ and *Am*₂; in this way the number of signal pulses counted via *Ph*₁ and *Am*₁ is made independent of the laser intensity.) *Com* control and computing unit.

of sodium because the Raman spectra of potassium-silicate glass generally have somewhat sharper peaks and have been studied more extensively than those of sodium-silicate glass. The glass of the composition we used has much the same melting and fining behaviour as glasses in normal use, and both types of glass have the same silicon content.

Results of measurements

The Raman spectra presented in *fig. 4* will now be discussed. These were recorded for powdered specimens with the equipment described above, after the materials had been heated for an hour at a tempera-

ture of 700, 800 or 850 °C and then rapidly cooled to room temperature. Details of the preparation of the specimens and of the rest of the measurement procedure, and also of the way in which the contributions from the various reaction products in the spectrum were identified, have been given elsewhere [2].

After it has been heated to 700 °C, the specimen still only gives the peaks due to SiO₂ (*s*), the carbonate ion in a crystalline lattice (*c*) and the AsO₄³⁻ group, occurring in a crystalline lattice of K₃AsO₄. This can be seen from the lower spectrum in *fig. 4*. At 700 °C the arsenic, which was added in the form of As₂O₃, therefore only occurs in its pentavalent state, presumably owing to oxidation by oxygen in the ambient atmosphere.

After heating to 800 °C the specimen gives the first peaks that indicate the occurrence of a vitreous state. The peaks *mg* are due to non-ordered metasilicate chains (which are represented separately in *fig. 5a*), while the peaks *cg* originate from carbonate ions dissolved in the glass phase. There are also the peaks *d* originating from crystalline disilicate (K₂O.2SiO₂). The upper spectrum shows that the development that starts at 800 °C continues at 850 °C.

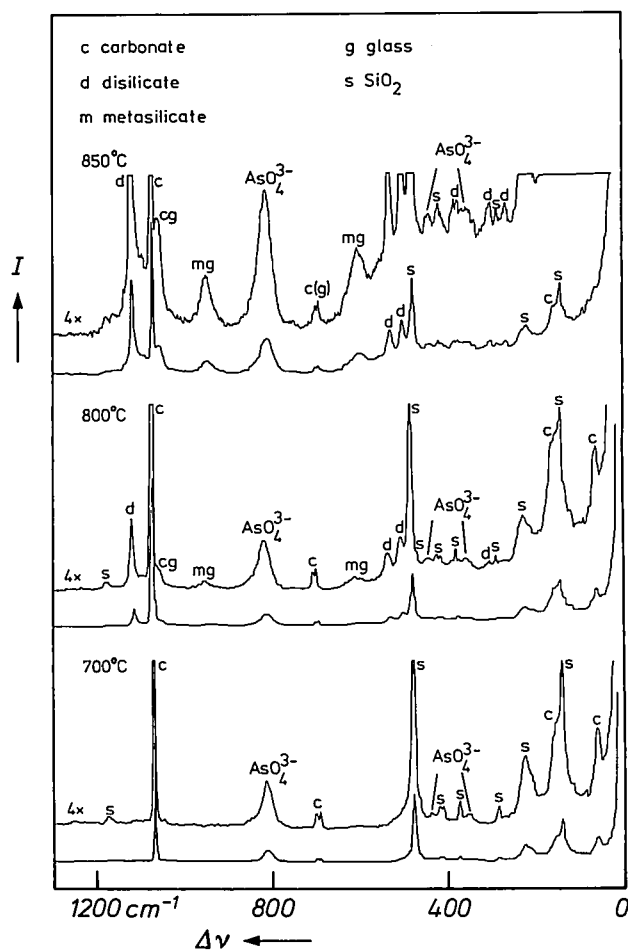


Fig. 4. Raman spectra recorded with the equipment in *figs 2* and *3*. The spectra give the intensity *I* of Raman light scattered at right angles, as a function of the difference in wave number $\Delta\nu$ from that of the primary beam; they were recorded for powdered specimens with an initial composition of 70 mol% of SiO₂, 30 mol% of K₂CO₃ and 1 mol% of As₂O₃. Before the recording the specimens were heated for an hour at 700, 800 or 850 °C and then rapidly cooled to room temperature. The letters and chemical formulae indicate the atomic groups associated with the various parts of the spectrum. *s* SiO₂. *c* carbonate ion in ordered lattice of K₂CO₃. *d* crystalline disilicate (K₂O.2SiO₂). *mg* metasilicate glass. *dg* disilicate glass (see also *fig. 5*). *cg* carbonate ion dissolved in the liquid glass phase. It can also be seen that, after heating to 850 °C, the arsenic occurs only in the form of AsO₄³⁻ groups, that is to say solely in the pentavalent state.

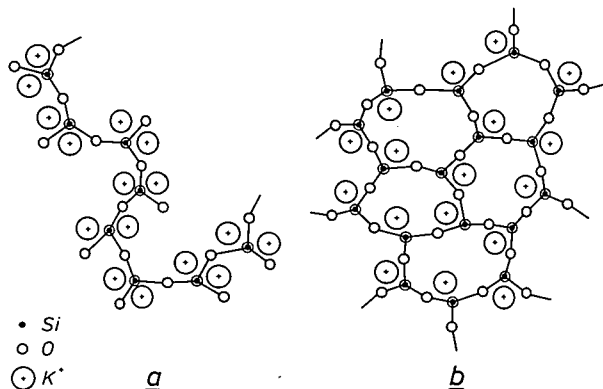


Fig. 5. Schematic 'planar' representation of the structures of metasilicate (*a*) and disilicate glass (*b*). Tetrahedra should be envisaged whose basal plane coincides with the plane of the drawing. The tetrahedra are projected on to this plane. Four oxygen atoms (○) are located at the corners of each tetrahedron, so that there are three oxygen atoms in the basal plane and one at the apex. At the centre of each tetrahedron is a silicon atom (•), which in this schematic projection coincides with the oxygen atom at the apex (⊙). In the linear chains of *a* (metasilicate glass) two oxygen atoms always belong to two tetrahedra. In each tetrahedron only half of the oxygen atoms should therefore be taken into account, and for each tetrahedron we therefore arrive at the formula SiO₃²⁻. Two potassium ions (⊙) should therefore be envisaged as added to each tetrahedron to compensate for the charge. In the more condensed structure of *b* (disilicate glass) three oxygen atoms always belong to two tetrahedra, so that the formula for each tetrahedron is SiO_{2.5}⁻, or, for each pair of tetrahedra, Si₂O₅²⁻. Here only two potassium ions are required for each two tetrahedra to compensate for the charge. By means of such condensations it is possible to compensate for the reduction of the potassium content that occurs in the liquid glass phase during the glass-forming process. Such a compensation is also the result of the condensation of the arsenic-containing AsO₄³⁻ groups into As₂O₇⁴⁻ groups, and also of the reduction of pentavalent arsenic to trivalent arsenic, in the form of AsO₂⁻ groups.

Fig. 6a illustrates schematically the various stages in the progress of these reactions in the specimen. As a result of reaction with K_2CO_3 , layers of crystalline $K_2O \cdot 2SiO_2$ (grey) form around grains of SiO_2 in a medium of liquid metasilicate, whose composition approximates to $K_2O \cdot SiO_2$; carbonate ions are dissolved in this liquid metasilicate.

On heating to above $850^\circ C$ changes occur that appear to be very important in the fining process. Fig. 7 shows Raman spectra for glass specimens that have been heated for an hour at temperatures of 900 , 950 , 1000 or $1100^\circ C$. We see at $900^\circ C$ the first emergence of peaks originating from $As_2O_7^{4-}$ and AsO_2^- groups, the latter thus relating to arsenic in its trivalent state. At this temperature we also see the first peaks originating from disilicate glass (dg), with the composition $K_2O \cdot 2SiO_2$ (fig. 5b). Finally we see that the peak of crystalline carbonate (c) has disappeared at this temperature, as might be expected since the melting point of K_2CO_3 is $891^\circ C$. At $1000^\circ C$ the peak of SiO_2 (s) has also disappeared, and at $1100^\circ C$ so has that of crystalline $K_2O \cdot 2SiO_2$ (d), thus completing the formation of the glassy state.

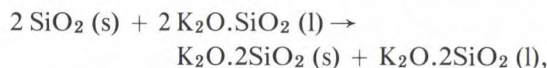
Fig. 6b summarizes a number of reactions that occur in this phase of glass formation: growth of the $K_2O \cdot 2SiO_2$ layer (grey) at the expense of the SiO_2 grain, the formation of disilicate glass in addition to metasilicate glass, the formation of $As_2O_7^{4-}$ and AsO_2^- groups in addition to AsO_4^{3-} groups and the formation of both CO_2 and O_2 bubbles.

Conclusions from the investigation

All these reaction processes in the temperature range up to $1100^\circ C$ considered here may be described as being connected with competition for the potassium ion. As more potassium ions find a place at the outside of the SiO_2 grains forming the compound $K_2O \cdot 2SiO_2$, and hence less positive charge is present for compensation in the glassy phase, so more CO_2 has to escape from the CO_3^{2-} ions in the form of bubbles:



The conversion of metasilicate glass into disilicate glass can also be regarded as a condensation reaction that liberates potassium ions to compensate for the reduced content of potassium in the liquid glass phase:



or, which amounts to the same thing, as a condensation reaction that liberates potassium ions for the solution of SiO_2 :

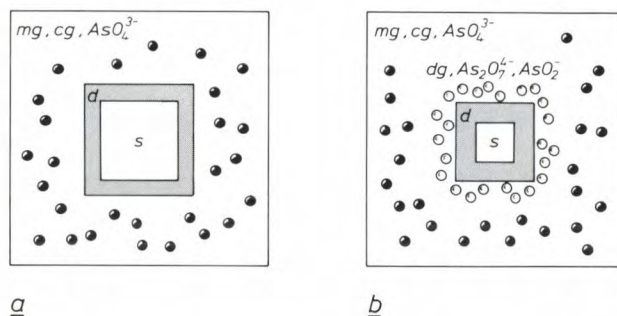


Fig. 6. Schematic representation of the progress of the glass-forming and fining processes. In the first stage (a), which takes place in the temperature range from 700 to $850^\circ C$, we see how the outer layer of a grain of SiO_2 is converted into crystalline $K_2O \cdot 2SiO_2$ (grey). Around this a metasilicate glass forms, in which CO_3^{2-} ions and AsO_4^{3-} groups are dissolved and CO_2 bubbles (\bullet) are generated. With the continued conversion of SiO_2 into $K_2O \cdot 2SiO_2$ at a temperature of 900 - $1100^\circ C$ (b), and with the consequent reduction of the potassium content in the liquid glass phase, there is the formation of disilicate glass, $As_2O_7^{4-}$ and AsO_2^- groups, and in addition to CO_2 bubbles, oxygen bubbles (\circ) are also formed; see also fig. 7. The oxygen formation is the result of the reduction of pentavalent arsenic to trivalent arsenic in the form of AsO_2^- groups.

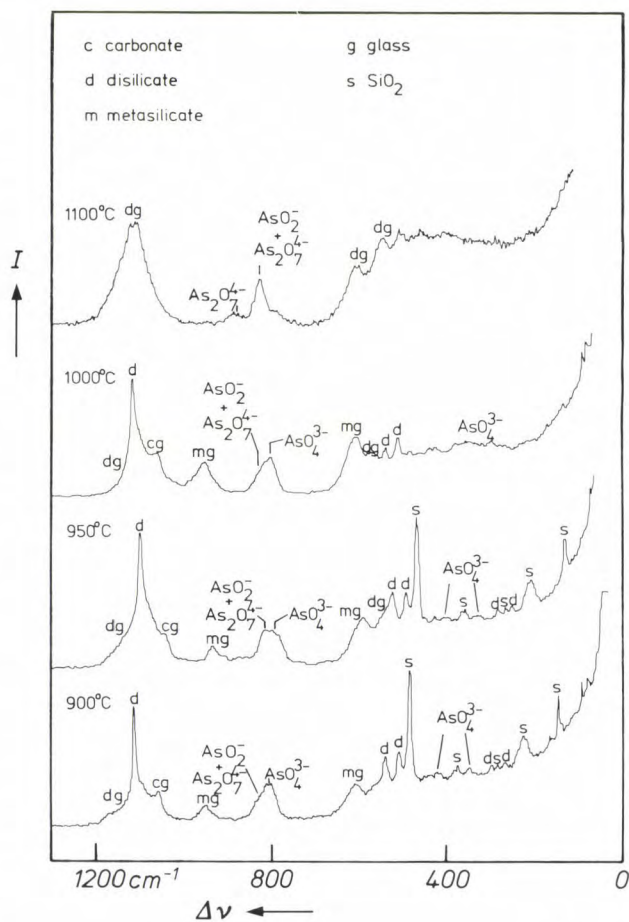
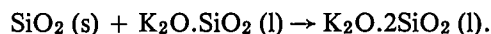
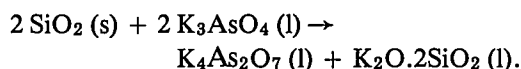


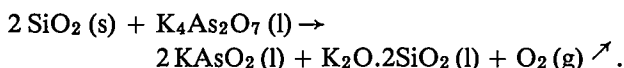
Fig. 7. Raman spectra recorded for specimens heated for an hour at 900 , 950 , 1000 or $1100^\circ C$. Other data as in fig. 4. The occurrence of the different structures in this temperature range corresponds to conversions as represented schematically in fig. 6b. The conclusion to be drawn is that during the glass-formation process the potassium content in the liquid glass phase is reduced, leading to all kinds of condensation reactions and to the reduction of the pentavalent arsenic.



The condensation of the ions containing arsenic can also be related to this:



The reduction of pentavalent arsenic can occur with similar results, accompanied by the formation of O₂ bubbles:



Our investigation thus indicates that changes in the concentration of the potassium ion during the glass-formation process lead to the reduction of arsenic, which is a prerequisite for the release of oxygen, and that this in turn is the most important element of the fining action.

Further research has confirmed that the shift in the As³⁺/As⁵⁺ equilibrium during glass formation is mainly *chemically* induced, owing to the changing content of potassium ions (or, in general, cations), and is not due only to temperature changes, as usually assumed previously.

This chemical induction of the shift in equilibrium does suggest that the oxygen not only forms later than the CO₂, but also forms at a different location, closer to the nucleus of the SiO₂ grains (fig. 6). The formation of O₂, separated from that of CO₂, could provide an explanation for the exceptional effectiveness of the

fining action of As₂O₃: any CO₂ produced would not necessarily appear as bubbles that expand under the influence of oxygen formation and disappear in this way from the melt. The CO₂ would also be displaced by the evolving oxygen. According to this hypothesis, this displacement of CO₂ would happen at the formation stage of the process, when open pores are still present in the reaction mixture, not yet completely liquid, and gas can still escape through these pores.

Although the simple composition of the system investigated permits no definite practical conclusion as yet, it does seem likely that the different light that our investigation casts on the fining mechanism will prove useful in attempts to find the most effective heat-treatment programme for the production of bubble-free glass.

Summary. At Philips Research Laboratories a study has been made of the fining action of arsenic oxides in glass formation. Specimens composed of 30 mol% K₂CO₃, 70 mol% SiO₂ and 1 mol% As₂O₃ were exposed for an hour to temperatures of 700, 800, 850, 900, 950, 1000 or 1100 °C. After cooling the specimens rapidly to room temperature, spectra were recorded with a laser-Raman spectrometer. Above 900 °C there is reduction from pentavalent to trivalent arsenic, accompanied by the release of oxygen. This oxygen has a fining action, either by causing the CO₂ bubbles produced earlier in the melt to expand, or by displacing the CO₂. Any oxygen bubbles still remaining dissolve again in the melt on cooling. The Raman spectra also suggest that the reduction of As⁵⁺ to As³⁺ is chemically induced, by a reduction of the potassium content in the liquid glass phase during the glass-forming process.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

- D. E. Aspnes** (Bell Laboratories, Murray Hill, N.J.) & **J. B. Theeten**: Optical properties of the interface between Si and its thermally grown oxide. *Phys. Rev. Letters* **43**, 1046-1050, 1979 (No. 14). *L*
- R. N. Bates & M. D. Coleman**: Millimetre-wave components made using PCB techniques. *Internecon UK '79*, etc., Brighton 1979, pp. 27-30. *R*
- J. Bloem**: Nucleation of silicon on amorphous and crystalline substrates. *Proc. 7th Int. Conf. on Chemical vapor deposition, Los Angeles 1979* (Electrochem. Soc. Proc. 79-3), pp. 41-58. *E*
- P. W. J. M. Boumans**: Inductively coupled plasma-atomic emission spectroscopy: its present and future position in analytical chemistry. *Fresenius' Z. anal. Chem.* **299**, 337-361, 1979 (No. 5). *E*
- P. W. J. M. Boumans**: ICP: d.c. arc in a new jacket? *Spectrochim. Acta* **35B**, 57-71, 1980 (No. 2). *E*
- J. C. Brice**: The lattice constants of a-quartz. *J. Mat. Sci.* **15**, 161-167, 1980 (No. 1). *R*
- F. J. A. den Broeder & J. van der Borst**: Magnetization reversal in $\text{Fe}_{80}\text{B}_{15}\text{Si}_5$ metallic glass with large uniaxial magnetorestrictive anisotropy. *J. appl. Phys.* **50**, 7116-7121, 1979 (No. 11, Part I). *E*
- A. Broese van Groenou**: Some results on the wear of a bronze-bonded grinding wheel. The science of ceramic machining and surface finishing II, eds B. J. Hockey & R. W. Rice, N.B.S. special publication 562, Washington, D.C., 1979, pp. 147-156. *E*
- A. Broese van Groenou**: Optimization of multistage grinding operations: the choice of consecutive removal rates. The science of ceramic machining and surface finishing II, eds B. J. Hockey & R. W. Rice, N.B.S. special publication 562, Washington, D.C., 1979, pp. 191-200. *E*
- A. Broese van Groenou & R. Brehm**: Experiments on high-speed grinding of ferrites and glasses. The science of ceramic machining and surface finishing II, eds B. J. Hockey & R. W. Rice, N.B.S. special publication 562, Washington, D.C., 1979, pp. 61-74. *E*
- A. Broese van Groenou, N. Maan & J. B. D. Veldkamp**: Single-point scratches as a basis for understanding grinding and lapping. The science of ceramic machining and surface finishing II, eds B. J. Hockey & R. W. Rice, N.B.S. special publication 562, Washington, D.C., 1979, pp. 43-60. *E*
- H. H. Brongersma, G. C. J. van der Ligt & G. Rouweler**: The reaction of bromine and oxygen with a tungsten surface studied by means of low-energy ion scattering. *Philips J. Res.* **36**, 1-14, 1981 (No. 1). *E*
- M. Brouha & J. van der Borst**: The effect of annealing conditions on the magneto-mechanical properties of Fe-B-Si amorphous ribbons. *J. appl. Phys.* **50**, 7594-7596, 1979 (No. 11, Part II). *E*
- K. H. J. Buschow & P. F. de Châtel** (University of Amsterdam); Hydrogen absorption and magnetic properties of intermetallic compounds based on 3d elements. *Pure and appl. Chem.* **52**, 135-146, 1980 (No. 1). *E*
- K. H. J. Buschow & A. G. Dirks**: On the crystallisation behaviour of amorphous alloys of rare-earths and 3d transition metals. *J. Phys. D* **13**, 251-258, 1980 (No. 2). *E*
- K. H. J. Buschow & R. M. van Essen**: Loss of ferromagnetism in YNi_3 after H_2 absorption. *Solid State Comm.* **32**, 1241-1242, 1979 (No. 12). *E*
- S. Colak, B. Singer & E. Stupp**: Lateral DMOS power transistor design. *IEEE Electron Device Letters* **EDL-1**, 51-53, 1980 (No. 4). *N*
- L. E. Cross** (Pennsylvania State Univ.) & **K. H. Härdtl**: Ferroelectrics. *Kirk-Othmer, Encyclopedia of chemical technology, 3rd ed., vol. 10, Wiley, New York 1980*, pp. 1-30. *A*

- P. Delsarte:** A generalization of the Legendre symbol for finite Abelian groups.
Discrete Math. **27**, 187-192, 1979 (No. 2). *B*
- P. Delsarte:** Regular schemes over a finite abelian group.
Geom. Dedicata **8**, 477-490, 1979 (No. 4). *B*
- P. Delsarte, Y. Genin & Y. Kamp:** Two-variable stability criteria.
Proc. 1979 Int. Symp. on Circuits and systems (ISCAS), Tokyo, pp. 495-498. *B*
- P. P. J. van Engelen:** Some ENDOR studies of 3d transition metal ions in semiconductors.
Thesis, Utrecht 1980. *E*
- K. Enke, H. Dimigen & H. Hübsch:** Frictional properties of diamondlike carbon layers.
Appl. Phys. Letters **36**, 291-292, 1980 (No. 4). *H*
- J. van Esdonk:** Verbindingstechnieken voor ultra-hoogvacuumsystemen.
Revue de la Soudure/Lastijdschrift **35**, 123-131, 1979 (No. 3). *E*
- J. A. Geurst:** Zilsel's variational principle derived from Lin's principle in superfluid hydrodynamics of ^4He .
Physics Letters **74A**, 225-228, 1979 (No. 3, 4). *E*
- B. H. A. Goddijn:** Static performance of a hybrid stepping motor with ring coils.
Thesis, Eindhoven 1980. *E*
- W. van Haeringen:** On the choice of tube lengths and diameters in a $\text{He}^3\text{-He}^4$ dilution refrigeration system.
Cryogenics **20**, 153-157, 1980 (No. 3). *E*
- J. Hallais & D. Boccon-Gibod:** Applications des hétérostructures GaAs-(Ga,Al)As aux transistors à effet de champ.
Acta Electronica **23**, 339-345, 1980 (No. 4). *L*
- K. H. Härdtl:** New aspects in ferroelectric ceramics.
Ferroelectrics **24**, 75-80, 1980 (No. 1/4). *A*
- P. Harrop:** Gallium arsenide field effect transistor mixers: theory and applications.
Acta Electronica **23**, 291-297, 1980 (No. 4). *L*
- J. B. Hughes, J. B. Coughlin, R. G. Harbott, T. H. J. van den Hurk* & B. J. van den Bergh* (* Philips Elcoma Division, Eindhoven):** A versatile ECL multiplexer IC for the Gbit/s range.
IEEE J. SC-14, 812-817, 1979 (No. 5). *R*
- F. M. Klaassen, H. J. Wilting & W. C. J. de Groot:** A UHF MOS tetrode with polysilicon gate.
Solid-State Electronics **23**, 23-30, 1980 (No. 1). *E*
- E. Klotz, U. Tiemens & H. Weiss:** X-ray tomography by grid coding.
Appl. Optics **19**, 474-476, 1980 (No. 3). *H*
- G. Kowalski, R. Rieckeheer & W. Wagner:** New means for picture formation in computer tomography.
Optik **55**, 67-86, 1980 (No. 1). *H*
- M. H. Kuhn, H. Ney, R. Geppert & R. Gierloff:** Automatische Zugangskontrolle mit Hilfe der „akustischen Unterschrift“.
VDI-Z **122**, 125-130, 1980 (No. 4). *H*
- F. P. J. Kuijpers & G. F. M. Beenker:** The exact solution of the Stefan problem describing the growth rate of binary III-V compounds for LPE with linear cooling.
J. Crystal Growth **48**, 411-415, 1980 (No. 3). *E*
- J. van Laar, A. Huijser & T. L. van Rooy:** Adsorption of type III and V elements on GaAs (110).
J. Vac. Sci. Technol. **16**, 1164-1167, 1979 (No. 5). *E*
- J. Lohstoh:** The punchthrough device as a passive exponential load in fast static bipolar RAM cells.
IEEE J. SC-14, 840-844, 1979 (No. 5). *E*
- F. Meyer, J. H. J. M. Buster, B. G. Bagley & D. E. Aspnes (Bell Laboratories, Murray Hill, N.J.):** Optical properties of the metallic glass $\text{Pd}_{0.775}\text{Cu}_{0.06}\text{Si}_{0.165}$ over the energy range 0.67 to 5.6 eV.
J. non-cryst. Solids **34**, 441-444, 1979 (No. 3). *E*
- R. F. Milsom:** Three-dimensional variational analysis of small crystal resonators.
Proc. 33rd Annual Symp. on Frequency control 1979, Atlantic City, pp. 263-270. *R*
- A. Mitonneau, A. Mircea, G. M. Martin & D. Pons:** Electron and hole capture cross-sections at deep centers in gallium arsenide.
Rev. Phys. appl. **14**, 853-861, 1979 (No. 10). *L*
- J. Nicolosi & J. Ladell:** X-ray topographic analysis of dislocation line defects in solution grown deuterated triglycine fluoberyllate.
J. Crystal Growth **49**, 120-124, 1980 (No. 1). *N*
- J. M. van Nieuwland & C. Weber:** Eigenmodes in non-rectangular reverberation rooms.
Noise Control Engng **13**, 112-121, 1979 (No. 3). *E*
- V. Pauker:** Amplificateur équilibré large bande hyperfréquences à transistor à effet de champ en GaAs.
Acta Electronica **23**, 299-315, 1980 (No. 4). *L*
- L. J. van der Pauw:** A rigorous set of shell equations derived from the law of conservation of energy.
Philips J. Res. **36**, 31-39, 1981 (No. 1). *E*
- J. G. J. Peelen:** Transparent hot-pressed alumina; I: hot pressing of alumina.
Ceramurgia Int. **5**, 70-75, 1979 (No. 2). *E*
- J. G. J. Peelen:** Transparent hot-pressed alumina; II: transparent versus translucent alumina.
Ceramurgia Int. **5**, 115-119, 1979 (No. 3). *E*
- A. Pirotte:** Fundamental and secondary issues in the design of non-procedural relational languages.
5th Int. Conf. on Very large data bases, Rio de Janeiro 1979, pp. 239-250. *B*
- A. Rabier:** Conception de circuits linéaires assistée par ordinateur.
Acta Electronica **23**, 277-289, 1980 (No. 4). *L*

- H. Rau:** Estimation of the homogeneity range of MoS₂. *J. Phys. Chem. Solids* **41**, 765-767, 1980 (No. 7). *A*
- J. L. Robert, B. Pistoulet, F. M. Roche, P. Girard, J. M. Dusseau, A. Raymond** (all with Université des Sciences et Techniques du Languedoc, Montpellier) & **G. Martin:** Model of conduction in highly compensated semiconductors. Application to semi-insulating GaAs. *Physics of semiconductors, 1978* (14th Int. Conf., Edinburgh; Inst. Phys. Conf. Ser. No. 43), pp. 569-572; 1979. *L*
- J. M. Robertson & H. A. Algra:** Inhomogeneities in bubble films measured by spin wave resonance. *J. appl. Phys.* **50**, 7810-7814, 1979 (No. 11, Part II). *E*
- J. M. Robertson, M. W. van Tol, W. H. Smits & J. P. H. Heynen:** Colourshift of the Ce³⁺ emission in monocrystalline epitaxially grown garnet layers. *Philips J. Res.* **36**, 15-30, 1981 (No. 1). *E*
- C. Schiepers** (Institute for Perception Research, Eindhoven): Response latency and accuracy in visual word recognition. *Perception & Psychophysics* **27**, 71-81, 1980 (No. 1).
- H. J. Schmitt:** Sehen mit Mikrowellen. *Umschau in Wiss. u. Technik* **80**, 82-87, 1980 (No. 3). *H*
- A. Schnell:** Nonlinear charge release of piezoelectric ceramics under uniaxial pressure. *Ferroelectrics* **28**, 351-353, 1980 (No. 1/4). *A*
- P. C. Scholten & D. L. A. Tjaden:** Mutual attraction of superparamagnetic particles. *J. Colloid Interface Sci.* **73**, 254-255, 1980 (No. 1). *E*
- P. J. Severin & H. van Esveld:** On the decomposition of total loss into absorption and scattering loss in compound glass fibres. *Optica Acta* **26**, 1415-1426, 1979 (No. 11). *E*
- B. M. Singer, W. G. Steneck, E. H. Stupp & R. V. Kurczewski** (Magnavox, Mahwah, N.J.): Suppression of pedestal noise in a pyroelectric vidicon. *IEEE Trans.* **ED-27**, 193-198, 1980 (No. 1). *N*
- B. Smets:** Atom formation and dissipation in electrothermal atomization. *Spectrochim. Acta* **35B**, 33-41, 1980 (No. 1). *E*
- J. Snel:** The effect of donors or acceptors on the Si-SiO₂ interface. *Insulating films on semiconductors, 1979* (Conf. Durham; Inst. Phys. Conf. Ser. No. 50), pp. 119-123; 1980. *E*
- J. P. Stagg & M. R. Boudry:** Lateral diffusion of Na⁺ ion at the Si-SiO₂ interface and Na⁺ neutralisation in the presence of chlorine. *Insulating films on semiconductors, 1979* (Conf. Durham; Inst. Phys. Conf. Ser. No. 50), pp. 75-80; 1980. *R*
- D. R. Terrell & U. Killat:** (*N*-vinylcarbazole containing) polymers, II. Sensitization to argon laser for use in single-layer photothermoplastic devices. *Photogr. Sci. Engng.* **24**, 25-31, 1980 (No. 1). *H*
- A. Thayse:** Discrete function expansions in integer powers. *Discrete appl. Math.* **1**, 127-136, 1979 (No. 1, 2). *B*
- A. Thayse:** Programmable and hardwired synthesis of discrete functions, Part I: One level addressing networks. *Philips J. Res.* **36**, 40-73, 1981 (No. 1). *B*
- C. Tsironis:** GaAs dual gate MESFET's and their applications in microwave circuits. *Acta Electronica* **23**, 317-324, 1980 (No. 4). *L*
- C. Tsironis:** 12 GHz receiver with a self-oscillating dual gate MESFET mixer. *Acta Electronica* **23**, 325-329, 1980 (No. 4). *L*
- M. J. Underhill:** Phase lock frequency synthesis for communications. *Symp. on Phase lock loops and applications, Delft 1980*, pp. 62-120. *R*
- M. J. Underhill & R. I. H. Scott:** The effect of the sampling action of phase comparators on frequency synthesizer performance. *Proc. 33rd Annual Symp. on Frequency control 1979, Atlantic City*, pp. 449-457. *R*
- J. A. T. Verhoeven:** Auger surface studies of barium on silicon oxide. *Appl. Surface Sci.* **4**, 242-246, 1980 (No. 2). *E*
- G. Verspui:** CVD of silicon carbide and silicon nitride on tools for electrochemical machining. *Proc. 7th Int. Conf. on Chemical vapor deposition, Los Angeles 1979* (Electrochem. Soc. Proc. 79-3), pp. 463-475. *E*
- J. F. Verwey:** Mobility and trapping of ions in SiO₂. *Insulating films on semiconductors, 1979* (Conf. Durham; Inst. Phys. Conf. Ser. No. 50), pp. 62-74; 1980. *E*
- R. P. Vincent:** Multipath problems in aircraft approach aids — a solution. *IEE Coll. on Modern techniques for combating multipath interference in radio, radar and sonar systems, London 1979*, 4 pp. *R*
- J. Vos:** Design characteristics of an advanced Stirling engine concept. *Proc. 14th Intersoc. Energy Conversion Engng Conf., Boston 1979, Vol. I*, pp. 1191-1196. *E*
- K. R. Whight:** Synthesis and analysis of d.l.t.s. spectra from m.o.s. surface states. *Electronics Letters* **15**, 744-745, 1979 (No. 23). *R*
- R. V. Winkle & C. H. Warner*** (* Kerry Ultrasonics Limited, Hitchin, England): A new method of quality control for ultrasonic wire bonding. *Ultrasonics International 79, Proc. Conf. Graz 1979*, pp. 62-68. *R*
- D. L. Wolters:** The role of water in the oxidation of silicon. *Insulating films on semiconductors, 1979* (Conf. Durham; Inst. Phys. Conf. Ser. No. 50), pp. 18-27; 1980. *E*

An automatic equalizer for echo reduction in Teletext on a single chip

J. O. Voorman, P. J. Snijder, J. S. Vromans and P. J. Barth

The signal standard for broadcast television includes some unused 'space'. In the European 625-line system there are twice 25 lines that are not filled with picture information. These 'field-flyback intervals' contain the frame synchronization pulses and some test signals, but there are also a number of unused lines in which information can be accommodated. Teletext takes advantage of this: digitally coded pictures with text and figures are transmitted simultaneously with the TV programme and can be stored in a memory in the receiver. The contents may range from weather reports and sports results to traffic information and the latest news. Teletext reception can sometimes be upset by strong echoes, even though these may not seriously affect the TV picture if the delay is short. A correction circuit has been designed that equalizes the transmission path, and thus compensates for the echoes. In this circuit the signals are processed by analog rather than digital methods. This reduces the extent and dissipation of the circuit required, and the authors have been able to integrate the circuit — including the capacitors — on a single chip. Their original publication^[1] received the second place in the '1981 Transactions Papers Award' of the IEEE Consumer Electronics Group.

Introduction

Teletext

Teletext is a television broadcasting service — already provided in several West European countries — in which, in addition to the regular television programme, information of general interest is transmitted, such as the latest news, the weather forecast, traffic information, sports results and so on. The Teletext information is arranged in 'pages', each taking up a full screen, and can only be displayed by a receiver with the appropriate circuits. The viewer can select the number of the page he requires from a list of contents, and key it in on a remote control unit, for example.

The Teletext information is transmitted during the field-flyback intervals and is coded in digital form. Echoes with a short delay can sometimes spoil the reception and in this article we present a circuit that

effectively removes the interference caused by such echoes. The system and organization of Teletext differ from one country to another, but this makes little difference to the problems of reliable reception. For convenience, we shall confine ourselves here to the Teletext standard introduced in the United Kingdom. Some details are given in *figs 1 and 2*.

Teletext reproduces the pages by using characters stored in a memory. These include letters and digits and a number of graphic elements from which illustrations, maps etc. can be composed. The display on the screen consists of 24 lines of 40 characters. The transmitted data determines the characters from the memory that will be displayed and their position.

Interference caused by echoes

One bit of the eight-bit words containing the Teletext information is usually a parity bit. This effectively protects the word from distortion, though some im-

Dr Ir J. O. Voorman, Ing. P. J. Snijder and Ing. P. J. Barth are with Philips Research Laboratories, Eindhoven; Ing. J. S. Vromans is with the Philips Video Division, Eindhoven.

portant words are given better protection. The Teletext signal therefore has reasonable protection from noise, interference from other transmitters and most echoes. More often than not, the picture quality of the television programme will have become unacceptable before the Teletext decoding breaks down. An

exception has to be made, however, for echoes with a delay shorter than $1 \mu\text{s}$. These are often strong enough to interfere with the decoding of the Teletext signal, although they may not seriously affect the television picture, since they are more or less concealed by the main signal (*fig. 3*)^{[1][2]}.

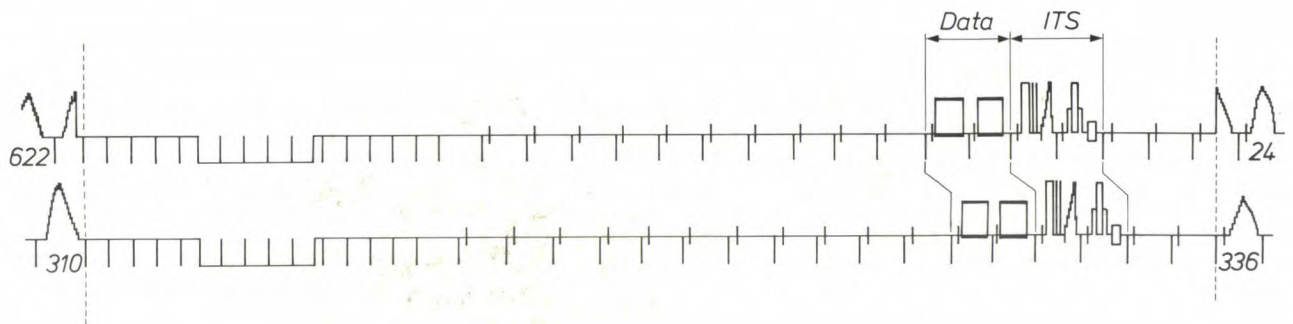


Fig. 1. The signals in the two field-flyback intervals of a television picture. Each line begins with a line-synchronization pulse. Above, from left to right can be seen the last field line, equalizing pulses, the frame-synchronization signal, more equalizing pulses, a number of empty lines, two lines carrying Teletext data (*Data*), two lines carrying test signals (*ITS*), and on the far right the first line for the next field scan. Each field line is preceded by a 'colour burst', a sample of the colour carrier used to synchronize the colour decoding. The numbering of the lines is indicated.

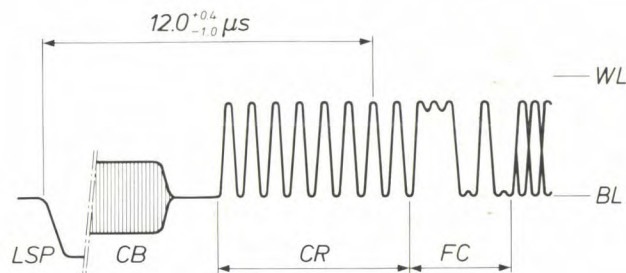


Fig. 2. Beginning of a Teletext line. *LSP* line-synchronization pulse. *CB* colour burst. *CR* clock run-in signal (two words of eight bits '10101010 10101010'). *FC* framing code for word synchronization ('11100100'). *WL* white level, *BL* black level of the TV signal.

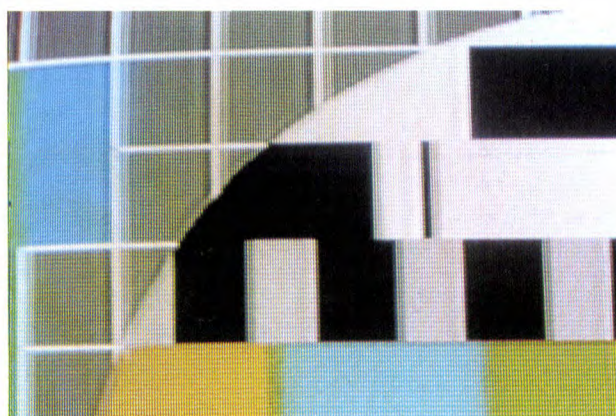


Fig. 3. Visibility in the test picture of an echo with a delay of 300 ns, an amplitude of 45% of the main signal and a phase angle of 0° . Such echoes are not unusual and viewers do not find them particularly annoying.

These echoes originate in the part of the transmission path where the television signal is the modulation on a carrier. They may be caused by multipath reception — for instance by reflections from hills or buildings — or they may be due to mismatched cable terminations in TV distribution systems, misalignment of the antenna or inaccurate tuning of the receiver.

Automatic equalizer

To counteract short-delay echoes we have developed an automatic equalizer^[3]. This is a circuit that eliminates echoes by subtracting a signal of the same waveform. The equalizer operates on the television signal in the baseband, i.e. after demodulation. It adjusts itself automatically to give optimum reduction of the echoes (*fig. 4*). Its operation depends on the binary nature of a correctly transmitted Teletext signal: it has only two levels, with fast transitions between them. Echoes introduce additional levels; these are taken as the starting point for the automatic equalization.

The equalizer takes the form of a transversal filter with variable coefficients; we shall say more about this later. It operates on the television signal in its ordinary analog form; we decided to use analog signal processing because the circuits required for a high information flow and many multiplications are simpler and take less current than digital circuits. We were able to design the equalizer on a single chip, which requires less than 500 mW of power. One special feature of the chip is that it contains integrated capacitors of rela-

tively high capacitance (a total of over two nanofarads; see *fig. 5*).

We shall now consider the waveforms that the echoes can assume after demodulation of the television signal and the use of a transversal filter to eliminate them.

periods, i.e. when the path difference is an integral number of wavelengths (*fig. 6b*). An extra half-wavelength (which is only about 25 cm in the UHF band) gives a negative echo of the waveform shape (*fig. 6c*). A phase shift of 90° or 270° gives the echo a different waveform (*fig. 6d*). It is then known as a quadrature



Fig. 4. Correction of the Teletext information by the equalizer. *Receiver I*: Teletext page degraded by the echo in *fig. 3*. Many graphic elements of the test picture have been replaced by random alphanumeric characters. *Receiver II*: The Teletext page corrected by the equalizer. The oscillograms on the left show the time signals before and after the equalizer; those on the right show the eye patterns.

Principle of echo compensation

Waveform of the echoes

Before demodulation an echo is just a delayed copy of the original signal. In the demodulation the phase of the carrier of the echo relative to that of the original signal is an important quantity. To recover the signal from the vestigial-sideband amplitude modulation used in television, it is necessary to regenerate the original carrier signal in the correct phase. Depending on the phase relative to the original carrier, the echo after demodulation can have various waveforms; see *fig. 6*.

An echo preserves the waveform of the original signal only if the delay is an integral number of carrier

echo. In case *d* the echo also changes the phase of the regenerated carrier signal in the receiver, so that the main signal also becomes distorted on demodulation.

The echo may thus have a different waveform from the original signal. We can consider it as a combination of delayed versions of the original signal, so that,

- [1] R. Klingler, Influence of receivers, decoders and transmission impairments on Teletext signals (UK Teletext and French Antiope), in: G. Cantraine and J. Destiné (eds), *New systems and services in telecommunications* (Int. Conf. Liège 1980), North-Holland, Amsterdam 1981, pp. 37-44.
- [2] Y. Ishigaki, Y. Okada, T. Hashimoto and T. Ishikawa, Television design aspects for better Teletext reception, *IEEE Trans. CE-26*, 622-628, 1980.
- [3] J. O. Voorman, P. J. Snijder, P. J. Barth and J. S. Vromans, A one-chip automatic equalizer for echo reduction in Teletext, *IEEE Trans. CE-27*, 512-529, 1981.

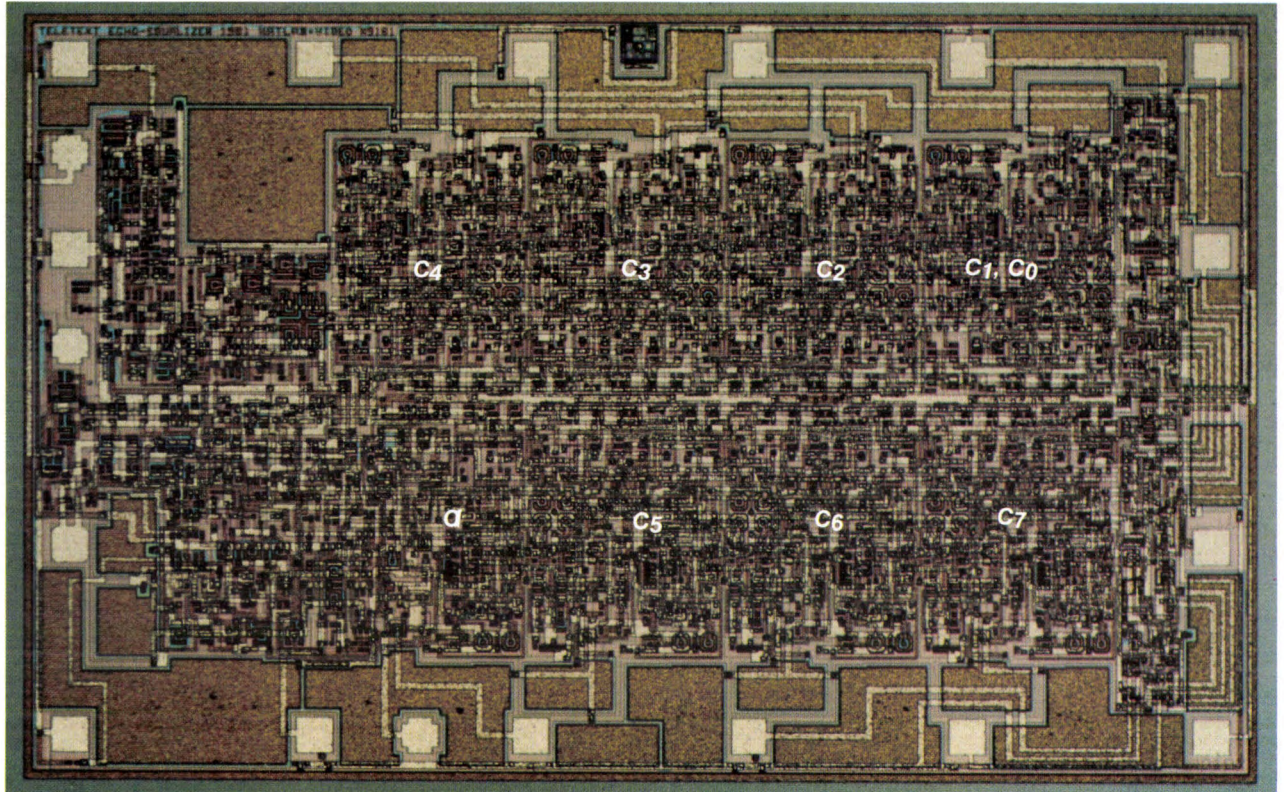


Fig. 5. Photomicrograph of the integrated equalizer (magnification $36\times$). The integrated capacitors (total capacitance more than 2 nF) can be seen along the edges. c_0 determination of d.c. correction. $c_1 \dots c_7$ the seven sections of the filter. a determination of amplitude of auxiliary signal.

to a first approximation, the quadrature echo in fig. 6d may be regarded as a positive echo immediately followed by a negative one.

We must not forget, of course, that the primary aim of the echo reduction is to permit error-free decoding of the Teletext data. The decoding is possible when the oscillogram of the received data signal has open

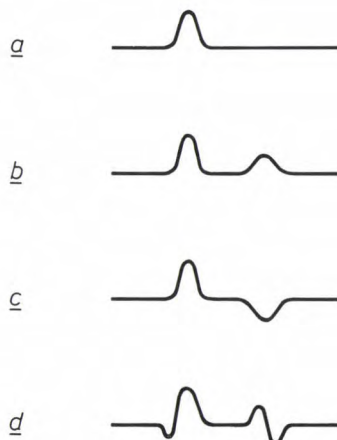


Fig. 6. Waveform of the echo after demodulation. *a*) Original signal. *b*) Echo in phase. *c*) Echo out of phase. *d*) Quadrature echo with a phase angle of 90° ; this produces a pre-echo before the main signal.

'eye patterns' (see fig. 4)^[4]. This does not require the echo reduction to be extremely accurate. The self-adjustment of the equalizer can therefore be made very fast, enabling it to track moving echo patterns, e.g. due to movement of an antenna in the wind.

Echo compensation by means of a transversal filter

The principle of echo compensation with a transversal filter is illustrated in fig. 7. It is based on the elimination of the echo by subtraction of a version of the main signal that has been delayed and given the appropriate amplitude. It can be seen that the elimination of either post- or pre-echoes leaves a higher-order and weaker echo; the more sections the filter has, the weaker is the higher-order echo and the further away it is from the main signal. Complete elimination is not possible with a transversal filter. The operation of the circuits, as represented in fig. 7, depends on the incorporated delays being equal to the time difference between main signal and echo. In practice, however, this time difference is not known, nor is it always the same. We could split up the delay line into as many sections as we expect to have echo delays. The subdivision need be no finer, however, than the value corresponding to the bandwidth of the television signal

(5 MHz): from the sampling theorem this is fully described by sampling at intervals of 100 ns. It is therefore unnecessary to make the delay per section shorter than 100 ns. The transversal filter then has the configuration shown in *fig. 8*. There are only seven sections, yet this is sufficient for the compensation of short echoes and the circuit need only be small. The attenuations in the branches — the ‘tap coefficients’ or ‘tap weights’ of the filter — are continuously adjusted automatically to give optimum cancellation of all echoes within the range determined by the filter length.

represents a convolution. We denote the filter coefficients by c_1, c_2, \dots, c_n and the response of the complete filter to an input signal $x(t)$ is

$$c_1x(t) * h_1(t) + c_2x(t) * h_2(t) + \dots + c_nx(t) * h_n(t).$$

After the addition of a d.c. term c_0 to make the output data symmetrical with respect to zero, the output signal becomes

$$y(t) = c_0 + c_1x(t) * h_1(t) + \dots + c_nx(t) * h_n(t). \quad (1)$$

This expression is easily written as the scalar product of the vectors ^[6]

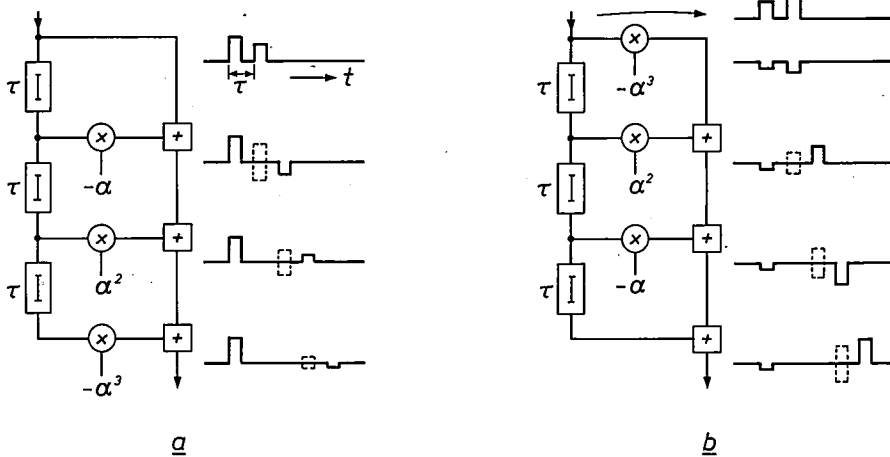


Fig. 7. Principle of echo reduction with a transversal filter. The delay τ is equal to the spacing between main signal and echo. *a)* Post-echo. Main signal and echo are delayed by a time τ , so that the delayed main signal coincides in time of arrival with the original echo. The delayed signal is multiplied by a factor α , giving the delayed main signal the same magnitude as the original echo. The delayed main signal is subtracted from this. At the same time an echo is produced at a spacing 2τ , but this is weaker than the original echo by a factor α . The process is repeated in the following sections of the filter. *b)* Pre-echo. The received signal is multiplied by a factor $-\alpha^3$. A version, delayed by a time τ , is multiplied by a factor α^2 . On summation the non-delayed main signal cancels the delayed pre-echo. This process is repeated in the following sections of the filter.

Automatic control of the filter coefficients

The circuit derives the appropriate filter-coefficient values from an error signal. A required condition is that the control process should converge towards the correct solution — the set of coefficients that minimizes the residual echoes — and that the solution should be a stable end state. We shall now take a closer look at this control process.

Keeping the description general, we assume that we have a filter with n taps. Let the response at tap k to a unit impulse at the input be $h_k(t)$. The response to an input signal $x(t)$ is then $x(t) * h_k(t)$, where the asterisk

$$c = (c_0, c_1, \dots, c_n)$$

and

$$x = (1, x(t) * h_1(t), \dots, x(t) * h_n(t)),$$

hence as

$$y(t) = c \cdot x. \quad (2)$$

If all the coefficients have the correct value, the output signal is as free from echoes as possible; we designate this reference state as

$$y_r(t) = c_r \cdot x. \quad (3)$$

If the reference signal is available for comparison, then control towards the end state can be achieved by determining the error $y - y_r$ and varying each coefficient c_k in the direction that will minimize the square of this error:

[4] F. W. de Vrijer, Modulation, Philips tech. Rev. 36, 305-362, 1976; see page 345.

[6] We do not treat these vectors as single-column matrices (as has been the usual practice in the literature). On multiplication we therefore use the notation for the scalar vector product and there is no transposition from column to row.

$$\frac{d}{dt} c_k = -f \frac{\partial}{\partial c_k} (y - y_r)^2 = -2f(y - y_r) \frac{\partial}{\partial c_k} (y - y_r),$$

where f is a positive loop gain.

Since y_r does not depend on the variable coefficients

$$c_k, \frac{\partial}{\partial c_k} y_r = 0, \text{ and from (1):}$$

$$\frac{\partial}{\partial c_k} y = x * h_k.$$

With the abbreviated notation $x * h_k = x_k$ it follows that

$$\frac{d}{dt} c_k = -2fx_k(y - y_r). \tag{4}$$

From this it can be shown that:

$$\frac{d}{dt} |c - c_r|^2 = -4f(y - y_r)^2. \tag{5}$$

As long as the output signal y differs from the reference signal y_r , the coefficient vector c converges to c_r .

A system in which a receiver generates a reference signal at the beginning of each Teletext line, for comparison with the same received signal, would certainly be feasible. However, as we noted earlier, we decided to adopt a solution based on the intrinsic binary nature of the Teletext signal, to provide us with information about the presence of echoes. This has the advantage that international agreement on a standard reference signal is not required.

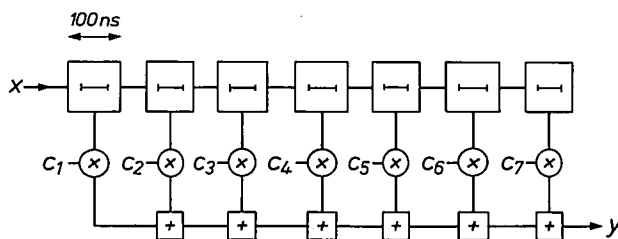


Fig. 8. The transversal filter. x input signal. y output signal. $c_1 \dots c_7$ the tap coefficients of the filter. Each delay element has a separate output for a tap.

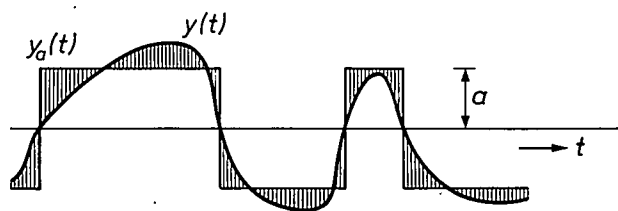


Fig. 9. The auxiliary signal $y_a(t) = a \operatorname{sgn} y(t)$, which is generated as an approximation to the correct Teletext signal. Adaptive control of the filter coefficients causes both signals y and y_a to converge towards the echo-free binary Teletext signal.

Our system does not therefore have a reference signal. Instead, we generate an auxiliary signal,

$$y_a(t) = a \operatorname{sgn} y(t), \tag{6}$$

which we may regard as a better approximation to the Teletext signal than $y(t)$ itself (see fig. 9). We try to make this signal and the output signal $y(t)$ converge by varying the coefficients c_k and the amplitude a . Both signals then tend towards the echo-free waveform.

This is done by starting from the error $y - y_a$, and as before we have

$$\frac{d}{dt} c_k = -f \frac{\partial}{\partial c_k} (y - y_a)^2 = -2fx_k(y - y_a). \tag{7}$$

In equation (7) it is assumed that $\partial y_a / \partial c_k$ may be neglected. This is reasonable, since y_a changes only at the zero-crossings, where the function $(y - y_a)^2$ is continuous.

In vector form this relation becomes:

$$\frac{d}{dt} c = -2fx(y - y_a). \tag{8}$$

An expression comparable with (5) cannot be derived from this vector equation. This would require dc_a/dt to be equal to zero, which is not the case, unlike dc_r/dt (c_a is the vector of the coefficients c_k , where $y = y_a$). Instead we can describe the convergence process as follows:

$$\begin{aligned} \frac{d}{dt} |c|^2 &= 2c \cdot \frac{d}{dt} c \\ &= -4fc \cdot x(y - y_a) \\ &= -4fy(y - y_a) \\ &= -4f|y|(|y| - a). \end{aligned} \tag{9}$$

This gives three possible equilibrium states: $y = 0$, $y = a$ and $y = -a$.

In the state $y = 0$ all the coefficients c_k are zero. To avoid this trivial solution we give one of the tap coefficients a fixed value: $c_m = 1$. If y is close to zero, (9) becomes

$$\frac{d}{dt} |c|^2 \approx +4f|y|a, \tag{10}$$

where a is positive and y can no longer become zero, since $c_m = 1$. It follows that if y is small, $|c|^2$ increases, so that the equilibrium state $y = 0$ is unstable.

There are then two equilibrium states left, $y = +a$ and $y = -a$. The circuit can become locked in one of these two states or alternate between the two. This

third mode is the required one. To prevent latching at a fixed value of $+a$ or $-a$, the amplitude a is set at zero if there have been no more zero-crossings of $y(t)$ (e.g. during half a line period). The d.c. correction c_0 now makes the mean value of the signal $y(t)$ equal to zero, so that zero-crossings occur again and the alternation is resumed.

Finally the coefficients c_k (including c_0) and the amplitude a are obtained by integration. Putting

$$y(t) - y_a(t) = \varepsilon(t) \tag{11}$$

to simplify, we have:

$$c_0 = -2f \int_0^t \varepsilon(\tau) d\tau,$$

$$c_k = -2f \int_0^t \varepsilon(\tau) x_k(\tau) d\tau, \quad k = 1, 2, \dots, n, \quad k \neq m,$$

$$a = +2f \int_0^t \varepsilon(\tau) \operatorname{sgn} y(\tau) d\tau. \tag{12}$$

Fig. 10 shows how the filter in fig. 8 is developed to form a network that performs the required operations. It contains the delay elements D_1, D_2, \dots, D_7 , each with an integrator that stores charge only during

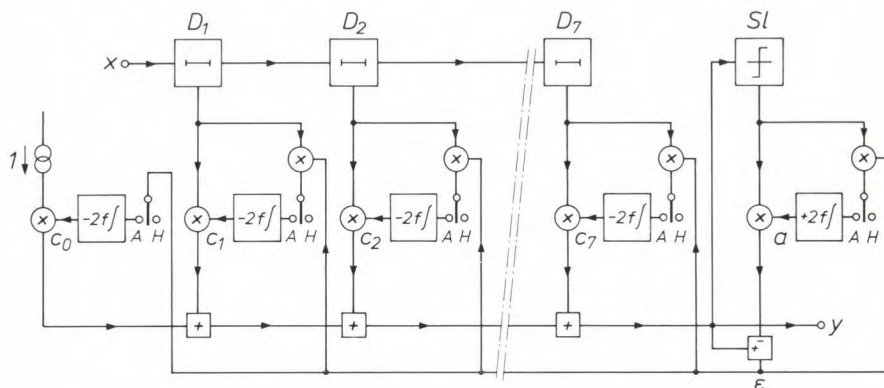


Fig. 10. Block diagram of the equalizer. x input signal. y output signal. $D_1 \dots D_7$ delay elements. A/H adapt/hold switches, which are closed only during reception of the Teletext signals, to permit adaptation of the coefficients $c_0 \dots c_7$ and the amplitude a . Sl slicer.

reception of a Teletext signal, because that is the only time when the switches A/H 'Adapt/Hold' are closed. While the remaining lines are being received the switches are open and the resulting values of the coefficients are preserved.

In the first section the d.c. correction c_0 is generated by integration of ε from eq. (12), and in the final section the amplitude a of the auxiliary signal $y_a(t)$ is generated. This requires the sign of the output signal, $\operatorname{sgn} y$, which is derived from a zero-crossing detector called a 'slicer'. The signal $\operatorname{sgn} y$ is also very suitable as input to a Teletext decoder.

All the tap coefficients of the filter should have a value < 1 except $c_m = 1$. This forces the main signal to take tap m ; the other taps reduce the echoes. This is consistent with the definition of the main signal as the 'largest echo' and improves the adaptation rate of the coefficients.

The electronic circuit

Principle of the delay line and other subcircuits

In principle, the required operations might be performed either digitally or by an analog circuit. In our case, as we saw, there are arguments against a digital approach. The data rate is high (7 Mbits/s), there are many multiplications and the ratio of the smallest time constant (the delay per section) to the largest (in the control loop) is three or four orders of magnitude. A digital circuit that gave the necessary performance would require a relatively large chip area and a high supply current.

We started from a common bipolar process with two layers of interconnection. A particular feature is the inclusion of integrated dielectric capacitors of

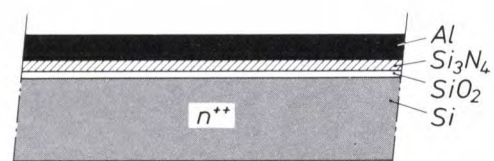


Fig. 11. Cross-section of an integrated capacitor.

relatively high value (see fig. 11). These are necessary because the filter coefficients are determined by the integral of a current. This integral is generated during the Teletext lines. Since the values obtained have to be held until the next Teletext lines arrive, the leakage current has to be extremely small.

Capacitors of this type are also used in the delay line. The analog delay line consists of a cascade of first-order phase shifters.

The transfer function from the filter input to tap k is:

$$H_k(p) = \frac{1}{1 + pRC} \left(\frac{1 - pRC}{1 + pRC} \right)^{k-1}, \quad (13)$$

where $p = \sigma + j\omega$ is the complex frequency, $\omega = 2\pi$ times the frequency in hertz, R is the resistance and C the capacitance in each phase shifter. The transfer function is the Laplace transform of the Laguerre function of order k , and the filter is therefore sometimes called a Laguerre filter [6].

The basic circuit of the phase shifters is given in *fig. 12*. The low-frequency component of the input current flows through the resistor R and is tapped (i_t). The high-frequency component flows through the capacitor C to a 'current mirror', which reverses the current, so that the high-frequency component is subtracted from the output current. This output current is the input current for the next phase shifter.

To make the amplitude ratio of the input and output currents frequency-independent in the practical circuit, correct matching is necessary and compensation for parasitic effects may be required because of the large bandwidth (5 MHz). First-order compensation can be obtained by adjusting the gain of the current mirror.

The principle of a multiplier circuit is given in *fig. 13*. Multiplication depends on the logarithmic relationship between the base-emitter voltage and the collec-

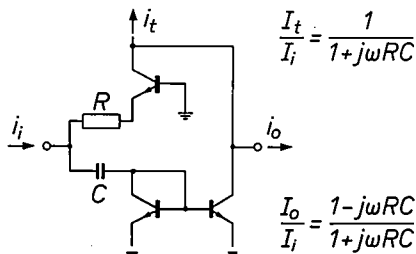


Fig. 12. Principle of the first-order phase shifter or 'Laguerre filter cell'. i_i input current, i_o output current. i_t tap current.

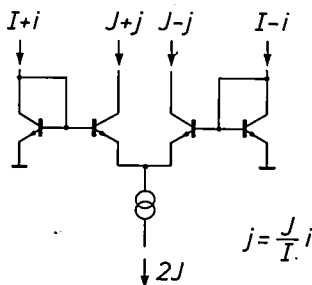


Fig. 13. Principle of the multiplier. Addition of the base-emitter voltages of bipolar transistors corresponds to multiplication of the collector currents, subtraction corresponds to division.

tor current of a bipolar transistor; addition of base-emitter voltages corresponds to multiplication of collector currents and subtraction corresponds to division. In *fig. 13*, where I and J are the supply currents, and i and j are the signal currents, then for identical transistors:

$$\frac{I + i}{J + j} = \frac{I - i}{J - j}, \quad (14)$$

so that

$$j = \frac{J}{I} i. \quad (15)$$

The signal current i can thus be multiplied by a desired factor by giving the supply current J the appropriate value.

The sign function $\text{sgn } y$ is generated in the 'slicer', a zero-crossing detector, shown in essentials in *fig. 14*. The output voltage y of the transversal filter is applied to the bases of the differential amplifier $T_{1,2}$. The cross-connected differential stage $T_{3,4}$ presents a negative resistance, which neutralizes most of the resistance of the base-emitter diodes of $T_{1,2}$. Consequently, when the input signal passes through zero, the differential stage $T_{1,2}$ switches very quickly from one extreme state to the other. The output transistors $T_{5,6}$ deliver a corresponding output current, whose amplitude is determined by the constant-current source a ; the difference current at the output is thus $a \text{sgn } y$.

The complete circuit

The complete circuit of a single section of the adaptive transversal filter is shown in *fig. 15*. *Fig. 15a* gives the 'Laguerre filter section'; the input signal is divided

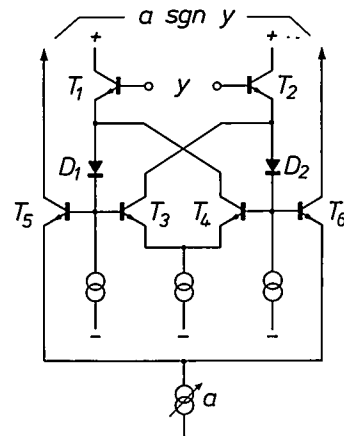


Fig. 14. Principle of the slicer. The emitter diode resistances of the differential stage $T_{1,2}$ are largely neutralized by the negative resistance presented by the cross-connected differential stage $T_{3,4}$, so that when y passes through zero the differential stage $T_{1,2}$ switches very rapidly from one extreme state to the other. $D_{1,2}$ diodes for d.c. level shift. $T_{5,6}$ output transistors.

between the resistor R and the capacitor C . The tapped signal x_k appears between the terminals A and B .

These appear again at two places in fig. 15b. The multiplication by the error signal ε (see (12)) takes place at the lower right; the current $\varepsilon(t)x_k(t)$ charges the capacitor C_{int} . The voltage across this capacitor is the integral of the product of the multiplication; in the multiplier circuit at the lower left the product $c_k x_k$ is produced; this product, expressed in terms of a difference current, goes to the output terminals P and Q , which are summation points for all taps.

In the filter section with a fixed coefficient ($c_m = 1$) the right-hand part of the circuit is not necessary, and

Practical results

The equalizer has been tested in the laboratory and in field trials, in the Netherlands and Switzerland and recently in Australia and Norway. The laboratory tests show that the equalizer is capable of restoring the Teletext signal even when there are exceptionally strong echoes. This is verified by the field trials, which showed that multiple echoes are also reduced effectively and that the Teletext reception only fails when the ordinary television reception has become unacceptable.

Our example of a laboratory test, has already been shown in fig. 4. An artificial echo is added with a variable delay, phase and amplitude. The input and

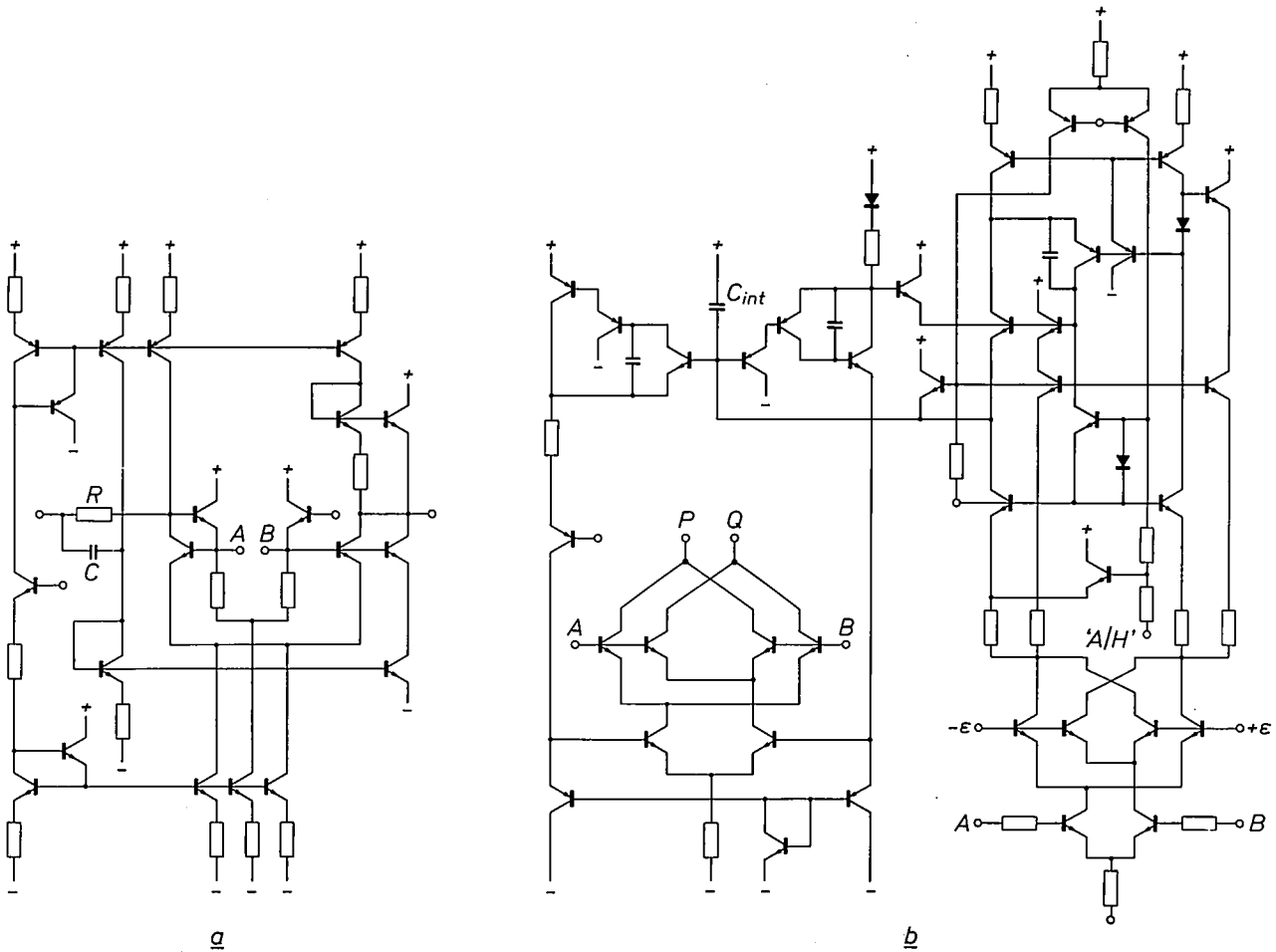


Fig. 15. Complete diagram of a section of the equalizer. a) Delay element (R, C) and tap (A, B). b) Right: Multiplication by the error signal $\varepsilon(t)$. 'A/H' adapt/hold control (see fig. 10). Left: Integration of the product in C_{int} and multiplication of the tap signal by the filter coefficient. P, Q difference-current summation points for all taps.

for calculating the d.c. correction term c_0 the left-hand part is redundant. The locations of the various circuits on the chip are indicated in fig. 5. The first section (c_1) of the transversal filter has a fixed setting; the other six ($c_2 \dots c_7$) are adaptive, as are c_0 and a . The integrated capacitors are located along the edge of the chip.

output signals of the equalizer after addition of an echo with a delay of 300 ns, an amplitude of 45% of the main signal and phase differences of $0^\circ, 90^\circ$ and 180° are given in fig. 16 a, b, c. In each picture the upper trace gives the input signal — part of a Teletext

[6] Y. W. Lee, Statistical theory of communication, Wiley, New York 1960.

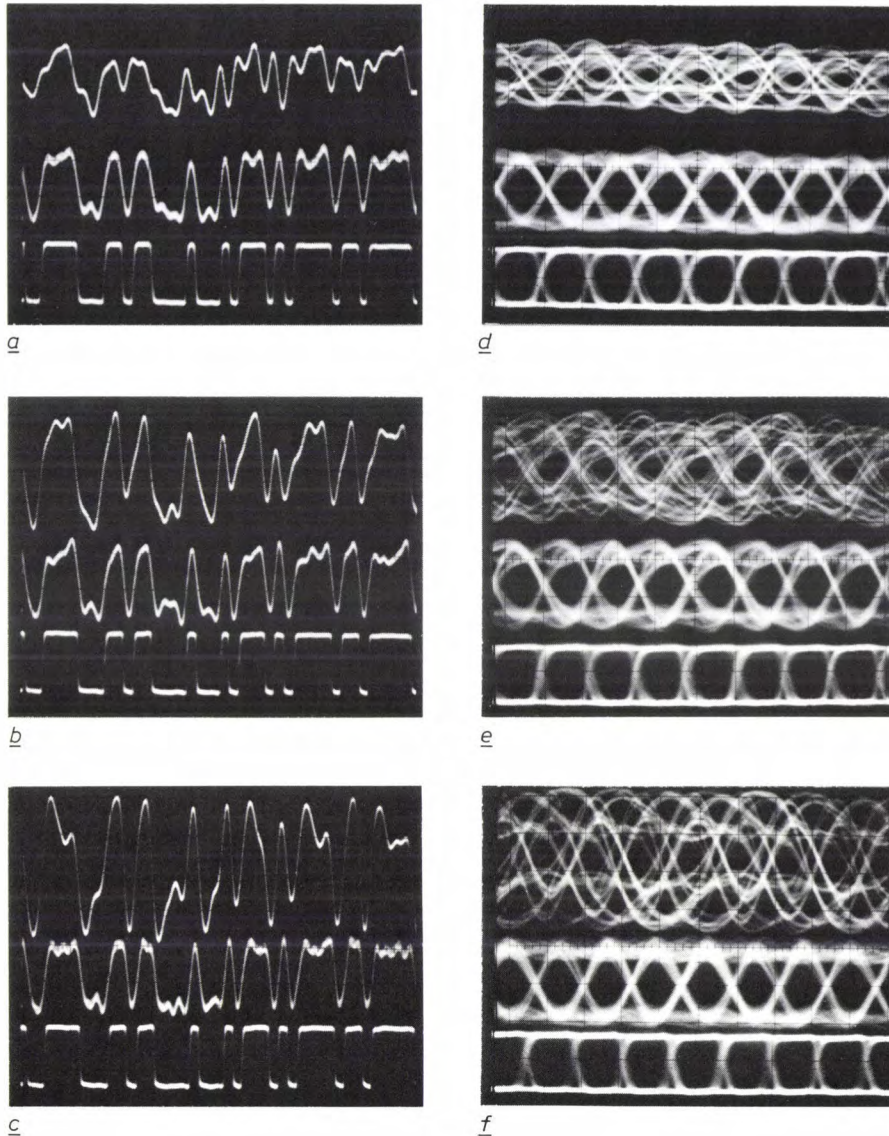


Fig. 16. Input and output signals of the equalizer. *a, b, c*) Waveforms. *d, e, f*) Eye patterns. The traces show, from top to bottom, the input signal, the direct output signal (y) and the output signal from the slicer ($\text{sgn } y$) for an echo with a delay of 300 ns, an amplitude of 45% and a phase angle of 0° (*a, d*), 90° (*b, e*) and 180° (*c, f*). It can be seen that in all cases the equalizer restores the binary nature of the Teletext signal and opens the eye patterns.

data signal chosen at random, with the echo added. It can clearly be seen that plateaus at about half height appear in the signal, making it difficult to decide whether a '1' or a '0' is intended. The centre trace gives the output signal y from the equalizer, showing that the binary nature of the signal has been restored. The lower trace gives the same output signal after it has passed through the slicer, which generates the function $\text{sgn } y$ from the function y and accentuates the binary nature of the signal, which is passed in this form to the Teletext decoder.

By way of further illustration, the eye patterns obtained under the conditions described are shown in fig. 16 *d, e* and *f*. These show that even with such a strong echo as this the conditions can be created for reliable Teletext reception.

These experimental results have been confirmed by the field trials. Tests on equalizers designed with dif-

ferent numbers of filter sections have shown that, even with multiple echoes, seven sections were sufficient to ensure reliable Teletext reception when an acceptable television picture was received.

Summary. The decoding of the digital Teletext signal can be seriously impaired by echoes. In particular, echoes that reach the receiver less than $1 \mu\text{s}$ after the main signal are often strong and occur quite frequently. They are accepted because they do not generally seriously degrade the quality of the television picture. The echoes can be compensated by subtraction of delayed and attenuated versions of the main signal. An equalizer has been developed that performs this function; it consists mainly of a transversal filter of seven sections, in which the filter coefficients are continuously and automatically controlled so as to produce a binary Teletext signal. The equalizer is integrated on a single chip, including a few capacitors with a total capacitance of more than 2 nF. Experiments and field trials have shown that the equalizer can ensure error-free Teletext reception in all conditions in which acceptable television reception is possible.

A linear d.c. motor with permanent magnets

L. Honds and K. H. Meyer

Take two permanent magnets and three iron bars and wind copper wire around one of the bars. This is the idea behind the linear motor designed at Philips Forschungslaboratorium in Aachen and described in the article below. The method of calculating the characteristics of the motor is just as simple. Measurements and calculations with MAGGY — a comprehensive program package for computing magnetic fields — demonstrate that the method of calculation yields the right results. The practical value of this kind of motor is demonstrated by its application in the 12 channel Transokomp 250 multipoint recorder developed by the Philips Science and Industry Division.

Introduction

Linear movements are frequently used in engineering. They are commonly produced by a drive mechanism operated by a rotor-type electric motor, since these are widely available commercially in many versions. However, the conversion mechanism introduces extra friction, takes up space and also increases the moving mass.

For certain applications, it is usual to produce a linear motion directly by means of a linear motor. The electromagnet with a movable armature, for example, is used in different variations in electric bells, in relays, in vibrators and so on. Although linear machines of this type can produce a considerable force, the stroke is limited. Also, the force is not the same for the various positions of the armature.

Other types of linear motors are suitable for traction applications. Linear machines of this type have attracted attention again because of the recent developments associated with passenger transport in densely populated urban areas. Such machines are required to have an 'unlimited' stroke and their force should not vary significantly during the displacement. Although the various electromechanical principles of rotary electrical machines still apply here ^[1], the machines used are usually linear induction motors or linear synchronous motors.

The type of linear motor we shall discuss in this article has a stroke varying between that of the two types mentioned above. In practice the stroke is limited to about half a metre. The force is virtually independent of the displacement. No current has to be sup-

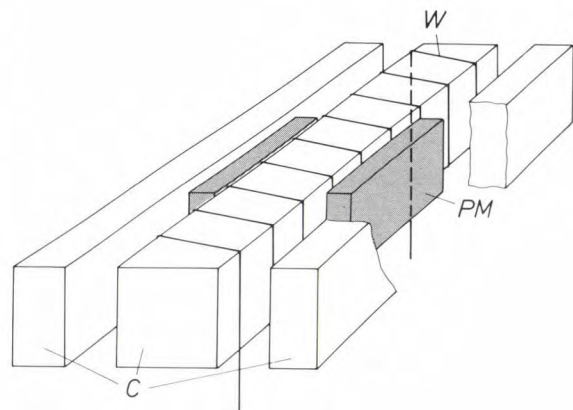


Fig. 1. Schematic diagram of the linear motor. *C* iron cores. *PM* permanent magnets. *W* turns of the coil, wound directly around the inner core, with no former. The two permanent magnets are fixed to a carriage (not shown).

plied to the moving parts and the design is extremely simple. It will therefore come as no surprise that the motor has proved highly successful in the 12 channel Transokomp 250 multipoint recorder developed by the Philips Science and Industry Division.

The configuration of our linear motor is shown in *fig. 1*. The non-moving part consists of three soft-iron cores *C*, which can be made of structural steel that has only undergone a simple anneal. The centre core is wound directly with insulated copper wire *W*; there is no former. The two permanent magnets *PM* form the moving part of the motor and are fixed to a carriage (not shown). The magnets are generally standard products (for example of Philips Ferroxdure, type 330), which require no finishing treatment. The version of

the motor used in the recorder gives a force of 2 N (constant to within 5%) over a stroke of 250 mm. The cross-section of this motor is 67 × 29 mm, and in the stationary state about 7 W is dissipated in the coil at a force of 2 N.

Several variations of the version shown in fig. 1 can be made. It is possible, for example, to ‘halve’ the motor, leaving only one permanent magnet and two iron cores. The magnet can also be stationary, with the coil moving. The current in that case, however, must be supplied by means of flexible wires or carbon brushes. The method of calculation used for these arrangements is almost identical, and we shall therefore confine ourselves to the configuration given in fig. 1.

First we shall explain how the force on the permanent magnets can be calculated. We shall then indicate the value to which the current in the coil has to be limited to prevent saturation of the iron of the cores. Next, we shall give the graphic results of calculations made with the computer program MAGGY. Finally we shall show how the motor design can be optimized for a particular application, and we shall give the results of measurements on the motor used in the Transokomp 250.

Calculation of the force on the permanent magnets

To calculate the force acting on the permanent magnets we start by considering the energy, as in calculating the torque acting on the permanent-magnet rotor of a synchronous electric motor [2]. In our case we then arrive at the following general equation [3] for the resultant force *F* on the permanent magnets:

$$F = I \frac{d\Phi_{CM}}{dp} - \frac{dW_M}{dp} + \frac{1}{2} I^2 \frac{dL_C}{dp}, \quad (1)$$

where *I* is the current in the coil, Φ_{CM} is the flux linkage produced by integrating the flux from the permanent magnets for each winding over the length of the coil, *p* is the coordinate defining the position of the permanent magnets, W_M is the magnetic field energy due to the permanent magnets, and L_C is the inductance of the coil; see fig. 2.

In equation (1) the second term on the right-hand side is equal to zero, because W_M does not depend on the position of the magnets. We assume here that the return field of flux density B_G is uniform and that the stray field of flux density B_{GE} does not depend on the position of the magnets either — provided that $-l < p < l - l_M$, where l_M is the length of the magnets and $2l$ is the total length of the cores. If we assume in addition that the relative permeability of the material of the permanent magnets is equal to 1, then the third

term on the right-hand side of (1) is also equal to zero. Equation (1) can then be expressed more simply as

$$F = I \frac{d\Phi_{CM}}{dp}. \quad (2)$$

Calculating the flux linkage Φ_{CM} and then differentiating with respect to *p* we obtain

$$F = I (B_E + B_G) \frac{l_M w}{l} b, \quad (3)$$

where B_E is the flux density of the effective field, *w* the number of turns of the coil and *b* the dimension of the cores perpendicular to the plane of the drawing; see fig. 2.

To determine Φ_{CM} we calculate the flux from the permanent magnets that is linked by $w dx/2l$ turns of the coil and integrate the result over the ranges $-l \leq x < p$, $p \leq x < p + l_M$ and $p + l_M \leq x < l$. For the first range we have

$$\Phi_{CM1} = \int_{-l}^p \{ \Phi_{GE} + 2(l+x)b B_G \} \frac{w dx}{2l},$$

where Φ_{GE} is the flux from the field of flux density B_{GE} emerging at the head ends of the inner core. The contribution from the second range to the flux linkage is

$$\Phi_{CM2} = \int_p^{p+l_M} \{ \Phi_{GE} + 2(l+p)b B_G - 2(x-p)b B_E \} \frac{w dx}{2l}.$$

For the third range the contribution to the flux linkage is:

$$\Phi_{CM3} = \int_{p+l_M}^l \{ \Phi_{GE} + 2(l+p)b B_G - 2l_M b B_E + 2(x-p-l_M)b B_G \} \frac{w dx}{2l}.$$

(Strictly speaking the three integrals should also contain a contribution from the stray field of flux density B_{GS} at the sides (see fig. 3), but in fact this makes no contribution in equation (3) since B_{GS} is independent of the variable *p*.)

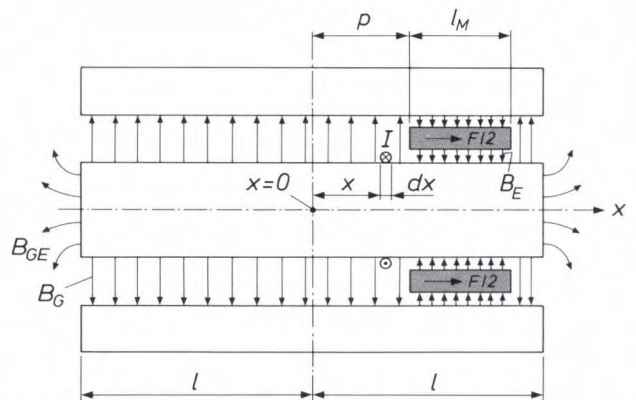


Fig. 2. Quantities for calculating the flux linkage of the coil for the field of the permanent magnets. B_E flux density of the effective field in the air gaps between the permanent magnets and the inner core, B_G flux density of the return field between the cores, B_{GE} flux density of the stray field at the head ends. *x* coordinate in the longitudinal direction with $x = 0$ for the centre of the motor, *p* coordinate defining the position of the permanent magnets. *I* current in the coil; the associated field is not shown. *b* height (not indicated) of the cores perpendicular to the plane of the drawing, l_M length of the magnets, $2l$ length of the cores. *F* resultant force on the permanent magnets.

If we assume that the total length $2l$ of the motor is large with respect to the length l_M of the permanent magnets, we can disregard the magnetic flux density B_G of the return field in comparison with the effective magnetic flux density B_E , and equation (3) simplifies to

$$F = \frac{l_M}{l} w I B_E b. \quad (4)$$

In the following we shall use equation (4) as our basis for the force calculations. Although this means that the theoretical value calculated for the force is slightly too small, the error is approximately compensated by

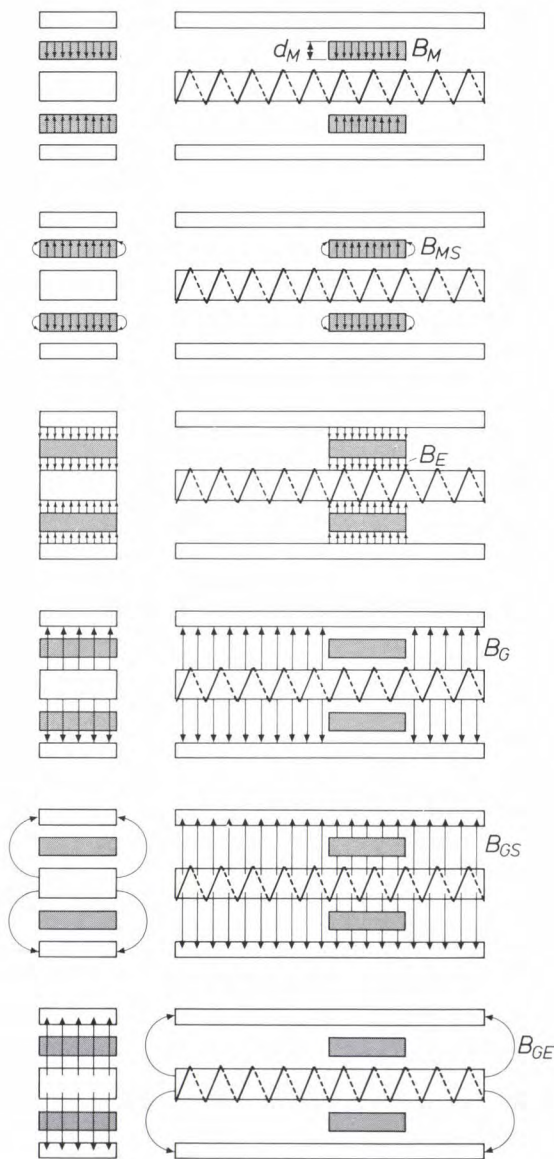


Fig. 3. The magnetic field from the two permanent magnets, split into a number of subfields. These are approximated by uniform or two-dimensional fields and are thus easily calculated. B_M flux density of the field in the permanent magnets. d_M thickness of the magnets. B_{MS} flux density of the stray field around the magnets. B_E flux density of the stray field at the sides of the cores. See also the caption to fig. 2.

the error due to assuming an idealized linear relationship between flux density and field-strength in the materials.

To calculate F from equation (4) it is necessary to determine the effective flux density B_E . This can be done in a simple way by dividing the total magnetic field into a number of subfields that are easier to calculate. Fig. 3 shows how the three-dimensional magnetic field of the permanent magnets can be split into a number of uniform and two-dimensional fields. The field due to the current in the coil has not been included here and we have assumed that the permeability of the iron is infinitely large. The subfields are:

- the field in the permanent magnets, of flux density B_M ;
- the stray field around the magnets, of flux density B_{MS} ;
- the effective field in the air gaps between the magnets and the inner core, of flux density B_E ;
- the return field between the iron cores, of flux density B_G ;
- the stray field at the two sides, of flux density B_{GS} ;
- the stray field at the two head ends, of flux density B_{GE} .

The calculation of the effective field

In the calculations that follow we shall consistently consider one half of the motor, consisting of one permanent magnet and half the coil. Fig. 4 indicates how the magnetic circuit with the coil not energized can be represented by an electrical equivalent circuit with constant-voltage source, resistances and currents.

The permanent magnet is the constant-voltage source in the equivalent circuit and delivers a magnetomotive force F_M . This quantity has the unit A (amperes) and

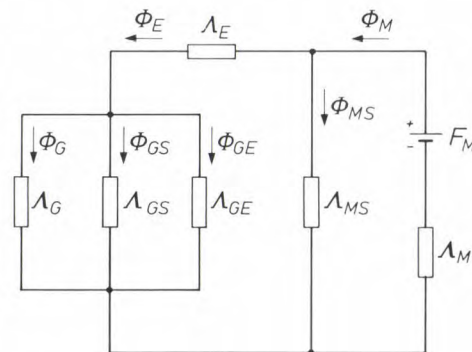


Fig. 4. Electrical equivalent circuit for a half of the magnetic circuit. $\Lambda_M, \Lambda_{MS}, \Lambda_E, \Lambda_G, \Lambda_{GS}$ and Λ_{GE} permeances of the different subfields, as indicated in fig. 3. F_M magnetomotive force, due to one of the permanent magnets. $\Phi_M, \Phi_{MS}, \Phi_E, \Phi_G, \Phi_{GS}$ and Φ_{GE} magnetic fluxes in the subfields of fig. 3.

[2] E. M. H. Kamerbeek, Electric motors, Philips tech. Rev. 33, 215-234, 1973.
 [3] This follows from equation (9) on page 223 of [2].

is the product of the apparent coercive field-strength H_c^* , see fig. 5, and the thickness d_M of the magnet:

$$F_M = d_M H_c^*. \tag{5}$$

The currents in the circuit are represented by the magnetic fluxes Φ in Vs or Wb (webers). The permeances Λ in H (henrys) of the magnetic resistances for a uniform field are equal to

$$\Lambda = \mu_0 \mu_r A/d, \tag{6}$$

where A is the area perpendicular to the direction of the lines of force, d is the length in the direction of the lines of force, μ_0 is the permeability of free space in H/m and μ_r is the relative permeability.

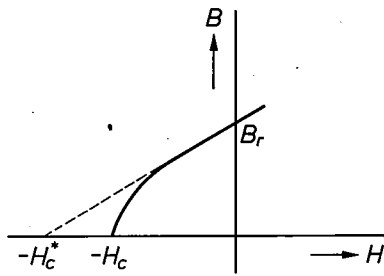


Fig. 5. Part of the magnetization curve for the material of the permanent magnets. B magnetic flux density, H field-strength. B_r remanent flux density, H_c coercive field-strength, H_c^* apparent coercive field-strength.

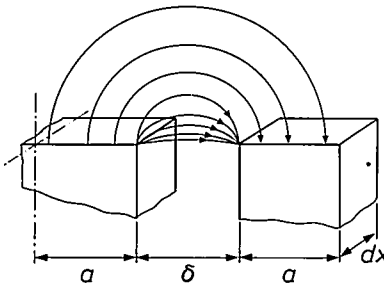


Fig. 6. Idealized geometry of the lines of force. The outer four lines of force are arcs of a circle, the others are parts of an ellipse. Idealizations of this kind can be used for deriving simple expressions for the permeance [4]. In this case, with the dimensions indicated in the figure, the permeance $d\Lambda$ of an element dx is given by

$$d\Lambda = \mu_0 \left[0.26 + \frac{1}{\pi} \ln \left(\frac{\delta + 2a}{\delta} \right) \right] dx,$$

where δ is the width of the air gap and a is the half-width of the inner core and the full width of each outer core.

Λ_E , Λ_G and Λ_M are easily calculated from equation (6). In determining Λ_M the value of B_r/H_c^* must be taken for the product $\mu_0 \mu_r$; see fig. 5. The permeances Λ_{MS} , Λ_{GS} and Λ_{GE} can be calculated from simple expressions that apply for idealized circular or elliptical lines of force [4]; see fig. 6. In this figure and in our calculations the half-width a of the inner core has been taken equal to the width of the outer cores.

[4] H. C. Roters, *Electromagnetic devices*, Wiley, New York 1970.

All the permeances and the magnetomotive force in the circuit of fig. 4 are therefore known. We can now set up a number of equations in which the fluxes are the unknowns and which we can solve for Φ_E :

$$\Phi_E = \Phi_G + \Phi_{GS} + \Phi_{GE} \tag{7}$$

$$\Phi_M = \Phi_E + \Phi_{MS} \tag{8}$$

$$\frac{\Phi_{MS}}{\Lambda_{MS}} = \frac{\Phi_E}{\Lambda_E} + \frac{\Phi_G + \Phi_{GS} + \Phi_{GE}}{\Lambda_G + \Lambda_{GS} + \Lambda_{GE}} \tag{9}$$

$$\frac{\Phi_{GS}}{\Lambda_{MS}} = F_M - \frac{\Phi_M}{\Lambda_M} \tag{10}$$

If we put $\Phi_E = \Lambda_M B_E$, with Λ_M equal to the area of the poles of the permanent magnet and $F_M \Lambda_M = B_r \Lambda_M$ (from (5) and (6)), we have the following expression for the effective flux density:

$$B_E = \frac{B_r}{\left(\frac{1}{\Lambda_G + \Lambda_{GS} + \Lambda_{GE}} + \frac{1}{\Lambda_E} \right) (\Lambda_{MS} + \Lambda_M) + 1} \tag{11}$$

Neglecting the effect of the stray fields, we have $\Lambda_{GS} = \Lambda_{GE} = \Lambda_{MS} = 0$, and eq. (11) reduces to the simpler relation:

$$B_E = \frac{B_r}{\left(\frac{1}{\Lambda_G} + \frac{1}{\Lambda_E} \right) \Lambda_M + 1} \tag{12}$$

The field in the cores

It was assumed above that when the subfields were summed the flux density of the core iron was on the linear part of the magnetization curve. To test the validity of this assumption, the flux density in the cores should be calculated. The field due to the current in the coil must then be taken into account as well. We have already assumed that the inner core has twice the cross-section of each of the outer cores. The flux densities in the three cores are therefore identical.

First of all we calculate the effect of the permanent magnet. In fig. 7a the fields in air are given for the case where the permanent magnet is in the extreme left-hand position, disregarding the stray fields. The flux density B_E of the effective field follows from (12). The flux density in the core reaches its maximum value $B_{C,max}$ at the position indicated and is equal to

$$B_{C,max} = \frac{l_M}{a} B_E. \tag{13}$$

We have also assumed that the fields in the air gaps are uniform. The flux density B_C in the core, after reaching $B_{C,max}$, therefore decreases linearly, as indicated by the solid line in fig. 7b. It is easily seen that the flux density behaves as indicated by the dashed

line if the permanent magnet takes up an intermediate position. The extreme values $B_{C,max}$ and $-B_{C,max}$ are thus reached at the extreme positions of the magnet.

Next we calculate the flux density B_I in the core due to the current I in the coil. It can be shown that

$$B_I = \frac{w I \mu_0}{4a l \delta} (l^2 - x^2), \quad (14)$$

where δ is the width of the air gap; see fig. 8a.

We calculate equation (14) by taking the line integral around the shaded rectangular area given in fig. 8a:

$$H_L \delta = \frac{x}{2l} wI,$$

where H_L is the field-strength in air at the position x . The field-strengths in the iron (with large μ_r), and also in air at the position $x = 0$ (reversal of the field), are set equal to zero in the integral. The magnetic potential $U_x = xwI/2l$ thus produces across the element dx of the air gap of permeance $d\Lambda = b\mu_0 dx/\delta$ a flux $d\Phi_I$ equal to

$$d\Phi_I = \frac{bwI\mu_0}{2l\delta} x dx.$$

Integration over the range from x to l gives an equation for the flux Φ_I in the core:

$$\Phi_I = \frac{bwI\mu_0}{4l\delta} (l^2 - x^2).$$

With $\Phi_I = abB_I$ we then obtain equation (14).

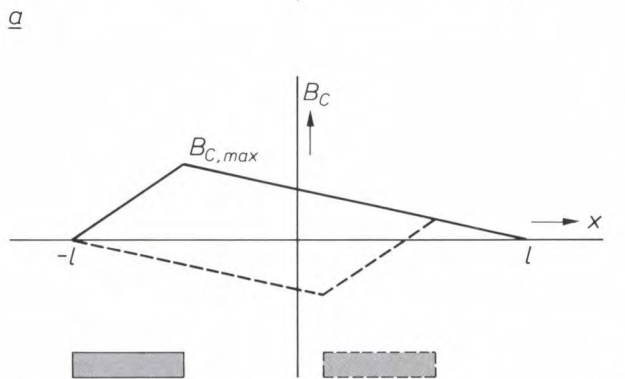
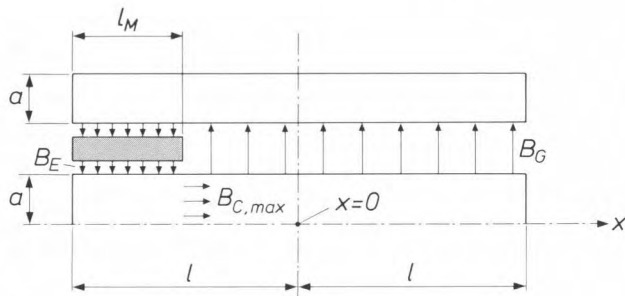


Fig. 7. a) Diagram of a 'motor half', for calculating the maximum flux density $B_{C,max}$ produced in the inner core by one of the permanent magnets. The effects of stray fields are not included. See also the captions to figs 2 and 6. b) Variation of the flux density B_C in the inner core as a function of the coordinate x . The solid line shows the variation when the magnet is in the extreme left-hand position; the dashed line relates to an intermediate position.

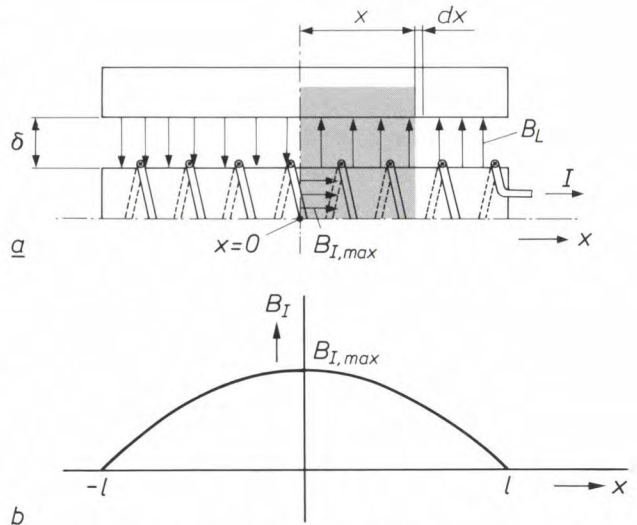


Fig. 8. a) Diagram of a motor half for calculating the maximum flux density $B_{I,max}$ produced by the current I in the coil, in the inner core. Here again, the effects of stray fields are not included. δ width of the air gap between the cores. B_L flux density of the field in air, equal to zero at $x = 0$. A line integral around the area shown shaded gives the magnetic flux in the core. b) Parabolic variation of the flux density B_I in the inner core, as a function of the coordinate x .

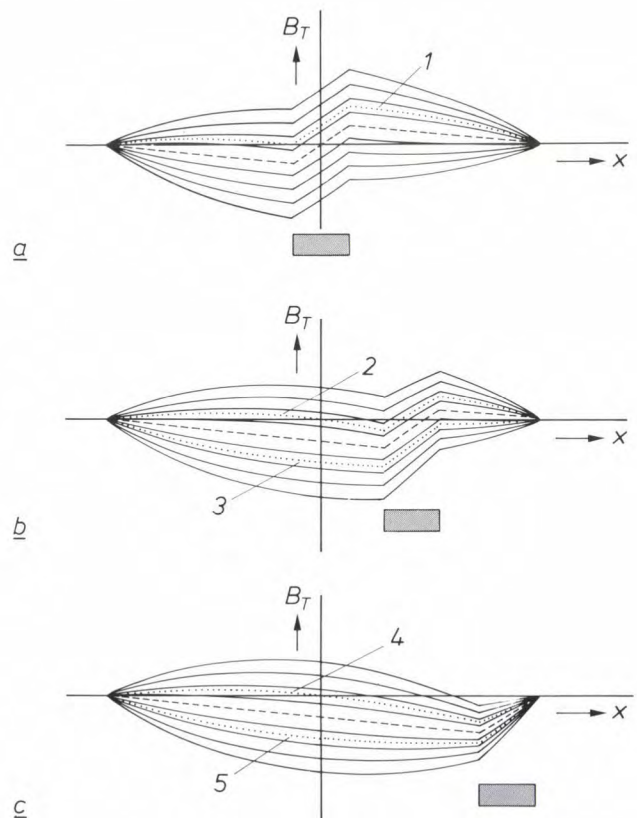


Fig. 9. Variation of the flux density B_T , the sum of the flux densities B_C and B_I of figs 7 and 8, as a function of position x . The dashed lines apply to the case $I = 0$, the solid curves above them relate to values of the current with the direction as indicated in fig. 8a, the solid lines below them relate to currents with the reverse direction. The dotted lines 1 to 5 correspond to figs 10a to e. a) Curves for the central position of the permanent magnet. b) Curves for an intermediate position. c) Curves for the extreme right-hand position of the permanent magnet.

The maximum value of the flux density in the core as given by eq. (14) is

$$B_{1,\max} = \frac{wI l \mu_0}{4a\delta}, \quad (15)$$

so that (14) can be written as

$$B_1 = B_{1,\max} \left(1 - \frac{x^2}{l^2} \right). \quad (16)$$

This quadratic behaviour of the flux density produced in the core by the coil is shown in fig. 8*b*. It follows from (15) that the height of this curve is proportional to the current I . To calculate B_C and B_1 more exactly, we must also consider the effects of the stray fields, as we did in determining B_E in (11); see fig. 3.

below saturation. It will now be clear that it is undesirable to have a motor configuration in which the ends of the cores are magnetically linked. Since the resistance of the magnetic circuit is then much smaller, the flux increases strongly, and the saturation of the iron is therefore reached at much lower values of I . The force F — for the same motor cross-section — consequently becomes smaller.

Calculating the field with MAGGY

To verify the foregoing we calculated the field for a number of cases with the program package MAGGY^[5]. Fig. 10 shows the lines of force calculated with the program for one of the motor halves. In figures 10*a-e* the densities of the lines of force

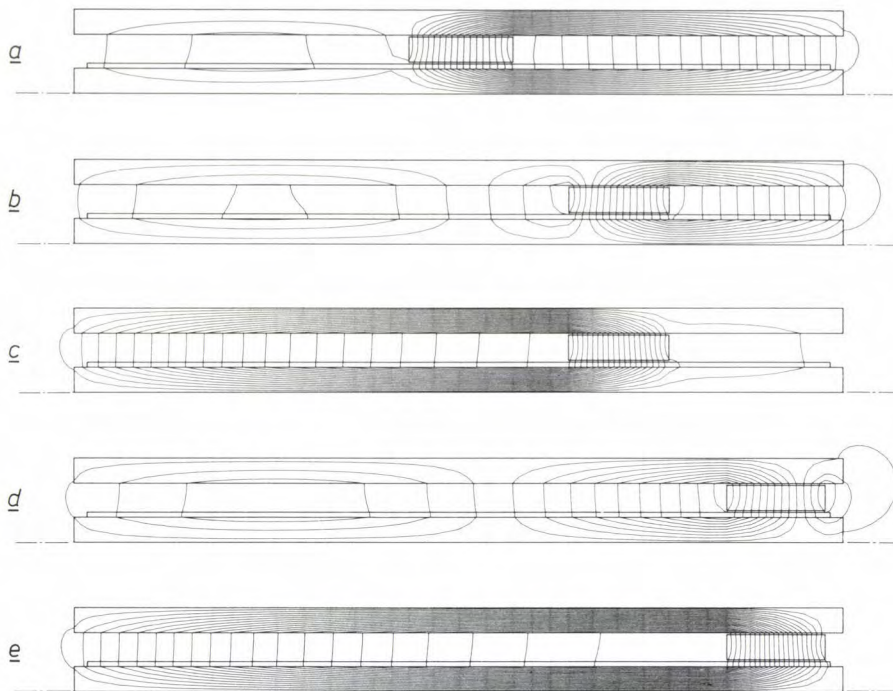


Fig. 10. Some results from the calculation of the field of a motor half with the computer program MAGGY. The densities of the lines of force (corresponding to the value of the magnetic flux density) are comparable for the five figures. *a*) to *e*) The results for different positions of the magnet and for different directions of the current (with the same absolute values), comparable with curves 1 to 5, respectively, in fig. 9.

If we add together the values of B_C and B_1 shown in fig. 7*b* and 8*b*, we obtain the curves given in fig. 9*a-c* for three different positions of the magnet and for different negative and positive values of the current through the coil. The total flux density $B_T = B_C + B_1$ reaches the highest value when B_C and B_1 have the same direction.

With the aid of the values calculated for B_T we can dimension our motor, and determine a maximum value for I , such that the iron of the cores is only just

(which are a measure of the flux density) are comparable with each other. The associated parameters are approximately equal to those of the respective dotted lines 1-5 in fig. 9. It is clear that the maximum flux density occurs in fig. 10*e* (curve 5); this corresponds to the calculations above. In fig. 10*b* and *d* the flux density in the iron changes direction twice. Fig. 10*d* does not completely correspond to curve 4 here because the cores protrude slightly beyond the permanent magnet.

Optimization of the motor design

The specifications of a motor with given dimensions can of course be calculated most accurately with MAGGY, because this program can take the nonlinear properties of the material into account. A MAGGY calculation requires a large computer, however and each 'call' on the program package is expensive. If it

computer. To calculate the dimensions from a given specification we use the following procedure.

A permanent magnet of Ferroxdure is selected and its specifications and those of the motor to be calculated are input to the program. The height b and lengths l and l_M are thus fixed. A first approximation is taken for the thickness of the copper winding and

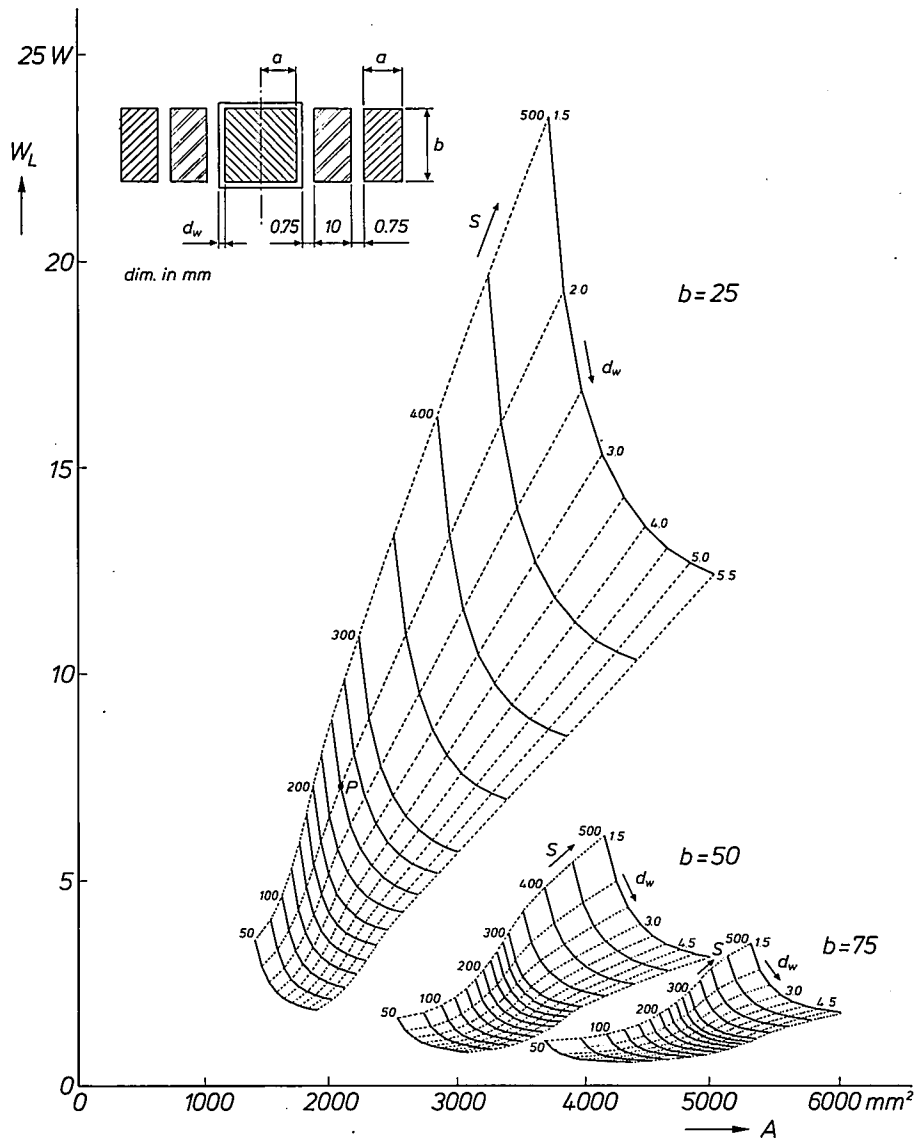


Fig. 11. The dissipation W_L in the coil as a function of the area $A = (4a + 2d_w + 23)b \text{ mm}^2$ of the motor cross-section, for different values of the stroke S , the thickness d_w of the copper turns, and the height b of the permanent magnets and the cores. The curves give the results of calculations with the authors' computer program. The point P corresponds to the dimensions of the motor used in the Transokomp 250 recorder. The force F on the permanent magnets is always 2 N.

is desired to calculate the dimensions for a given specification, a number of iterative steps are necessary and the use of MAGGY then becomes even more expensive.

This is why we decided to develop a faster program, which makes use of the method of stray-field idealization described above, and can be run on a simpler

for the clearance of the magnet between the cores: B_E can then be calculated from equation (12). From equation (4) the product wI can next be determined. Equation (15) gives a first approximation for the

[6] S. J. Polak, A. de Beer, A. Wachters and J. S. van Welij, MAGGY2, a program package for two-dimensional magnetostatic problems, Int. J. num. Meth. in Engng 15, 113-127, 1980.

dimension a , with a value of 1.2 Wb/m^2 (teslas) assumed for the maximum flux density in the core.

In a number of iterative steps — five are generally sufficient — the configuration of the motor is calculated more exactly. Since after the first step the dimensions are known approximately, the stray fields can be taken into account, so that equation (11) can be used. The total computer time required for this program is only a few seconds.

Fig. 11 gives a number of calculated values for a force of 2 N , with the stroke S , the copper thickness d_w of the coil and the height b of the permanent magnet as parameters. The three families of curves relate to permanent magnets 25 mm, 50 mm and 75 mm high; all have a cross-section of $40 \times 10 \text{ mm}$ and are made of Philips Ferroxdure type 330 [6]. The magnet of dimensions $40 \times 25 \times 10 \text{ mm}$ is a standard product and has been used in the motor for the Transokomp 250 recorder mentioned earlier [7]. In the graph the area of the motor cross-section $A = (4a + 2d_w + 23)b \text{ mm}^2$ is plotted along the horizontal axis and the steady-state dissipation W_L in the coil is plotted along the vertical axis. For each family of curves d_w varies from 1.5 to 5.5 mm and S varies from 50 to 500 mm. The point P relates to the dimensions of the Transokomp 250 motor.

With the aid of graphs like those in fig. 11 it is possible to determine the parameters of an optimum motor for a particular application. It is not possible to formulate general rules for this optimization. This is because it is the application that determines which is the most important parameter: the moving mass (the dimensions of the permanent magnet), the external dimensions of the motor (the quantity A), the energy consumption (the quantity W_L) or the cost of materials (mainly determined by d_w).

Provided the dimension d_w remains the same, the actual thickness of the wire in the winding does not matter. Halving the wire thickness has the effect of quadrupling w , but I is then reduced by a factor of four, since the maximum flux density in the core is determined by the product wI . The resistance of the coil then becomes sixteen times larger, so that the dissipation remains unchanged. In fig. 11 no account has been taken of the transfer of heat dissipated in the coil. Once a particular configuration has been selected, it is necessary to determine by calculation or experiment whether the temperature in the coil will become too high.

It can also be seen from fig. 11 that with increasing stroke the dimensions soon become larger. This is more clearly illustrated in fig. 12, where the calculated values of fig. 11 are plotted for $d_w = 2 \text{ mm}$, for larger values of S , with the width $2a$ of the centre core and

the stroke S along the two axes. With a stroke of 1 m the motor becomes inadmissibly large. This also follows from the result that with a stroke of 1 m the total mass becomes 82 kg, 59 kg or 56 kg for a height b of 25, 50 or 75 mm respectively. For $b = 75$, however, the acceleration available is only a third of that for $b = 25$. Although a clear dividing line cannot be drawn, in practice the stroke must be limited to a value of about 500 mm. Better values are obtained if

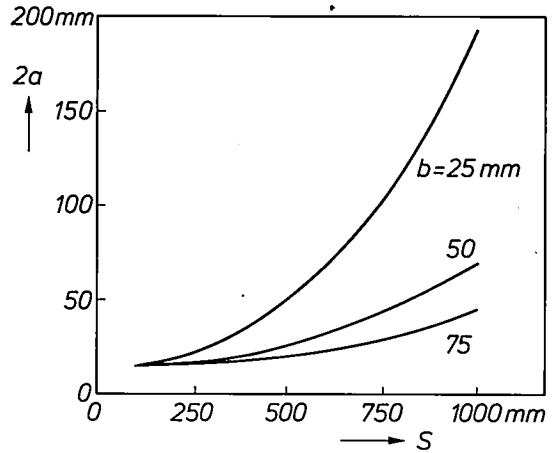


Fig. 12. The width $2a$ of the inner core as a function of the stroke S , for the data in fig. 11, but with larger values of S , with $d_w = 2 \text{ mm}$. It can be seen that for S greater than about 500 mm the motor becomes very large. For $S = 1000 \text{ mm}$ the total mass of the motor is 56 kg (for $b = 75 \text{ mm}$), 59 kg (for $b = 50 \text{ mm}$) or 82 kg (for $b = 25 \text{ mm}$).

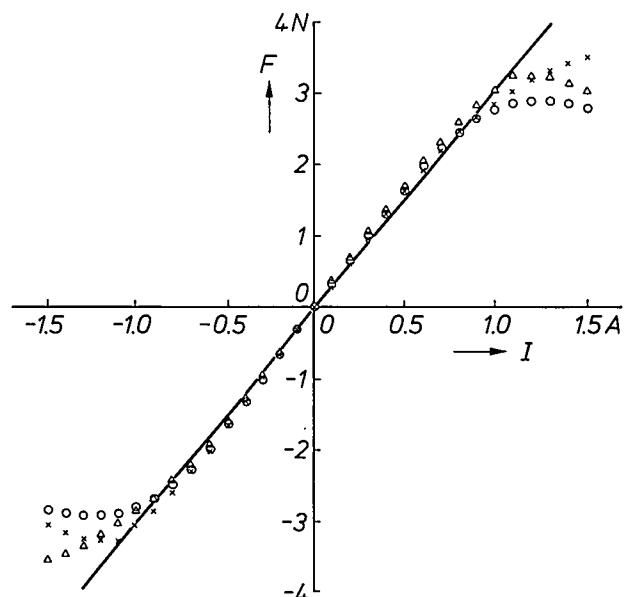


Fig. 13. Measurement of the force F as a function of the current I , for the motor used in the Transokomp 250 recorder. The measured points indicated in three different ways correspond to three different calculated positions of the permanent magnets. The solid line relates to the calculated results. Up to a value of $F = 2.5 \text{ N}$ the calculations and measurements are in good agreement; beyond this value the agreement deteriorates as the iron of the cores becomes saturated.

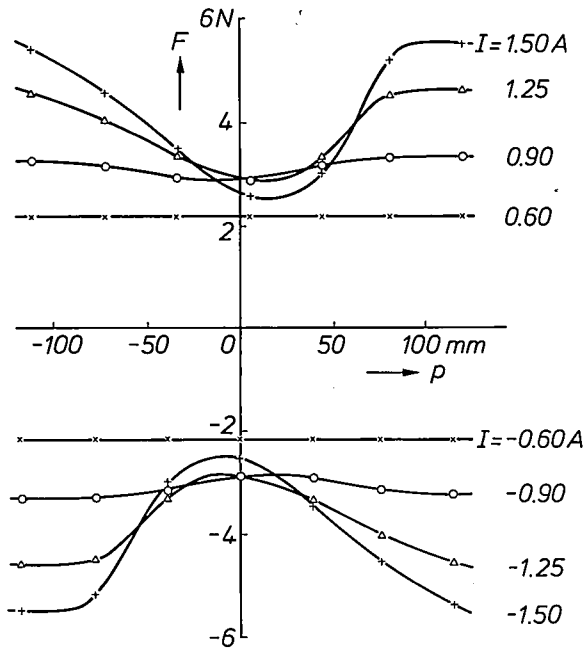


Fig. 14. Measurement of the force F as a function of the position p of the permanent magnets, for different positive and negative values of the current I . F is constant provided the iron of the cores is not saturated. If the absolute value of the current becomes greater than 0.6 A, the force in the central position of the magnets becomes smaller than at the extreme positions. If the current assumes an even higher value, the increase in the current is no longer sufficient to compensate for the decrease in the effective flux density. The force in the central position then becomes smaller than at lower current values.

current is supplied only to the part of the coil situated in the effective field of the permanent magnets. In that case, however, carbon brushes or other aids have to be used, and this makes the arrangement less attractive. In many cases the moving mass can be reduced by causing the coil to move and keeping the components of the magnetic circuit stationary. But flexible supply cables or carbon brushes are then necessary, and the length of the motor is almost doubled.

Verification of the calculations by measurements

To verify the calculations we carried out a number of measurements on the motor we designed for the Transokomp 250 recorder. Fig. 13 gives a plot of the

force F on the permanent magnets as a function of the current I in the coil. The solid line represents the calculated behaviour; the points correspond to the measurements at different positions of the permanent magnet. It can be seen that the calculations and measurements are in good agreement as long as the iron of the cores is not saturated. The specified force of 2 N is easily reached. When the current is increased to make the force exceed a value of 2.5 N, the deviation between the calculated and measured values increases rapidly.

Fig. 14 gives a plot of the measured force F as a function of the position p of the permanent magnets, measured at different currents I . As long as the iron of the cores is not saturated, at a current of 0.6 A or less, the force is virtually the same in all positions. When the current is increased, the force in the central position is less than in the outer positions. When the current is increased still further, the decrease in effective flux density is no longer compensated by the increase in current. The force in the centre position is then less than it is at lower currents. The higher saturation of the cores apparently adversely affects the location of the operating point of the permanent magnets. It is found from the measurements that operating the motor in the linear part of the magnetization curve of the core iron is not only a condition for the validity of the method of calculation, but is also necessary if the motor is to have the optimum characteristics.

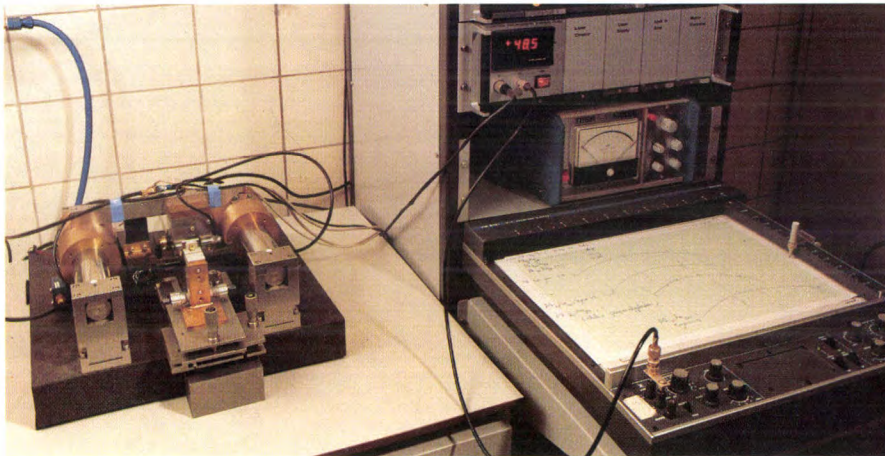
[6] Philips Data handbook, Components and materials, Part 16, 1982.

[7] J. Langlois and H. Foetzki, Linearmotor und Regelung des Nachlaufsystems für einen 12-Kanal-Punktdrucker, ETZ Arch. 3, 77-83, 1981.

Summary. The linear motor described here will provide strokes of about 100 to 500 mm and a force that is virtually constant over the full range of the stroke. The behaviour of the motor is easily calculated by assuming geometrically simple lines of force of the stray fields. The limit that the current in the motor coil must not exceed is determined by the saturation of the core iron. The validity of the calculations has been demonstrated by measurements and by making use of the MAGGY program package — it is however less expensive to use the program devised by the authors. Finally, procedures for optimizing the design of a motor are discussed.

An instrument for measuring the curvature of reflecting surfaces

A. F. Foederer, J. L. M. Hagen and A. G. van Nie



In engineering metrology the shape of a surface is determined by establishing its lines of intersection with planes perpendicular to a reference surface. These measurements may be made, for example, to determine variations in flatness of the surfaces of granite or cast-iron surface plates used in workshops or metrology laboratories. The measurements can also be made on a smaller scale for determining mechanical stresses in coatings applied to test samples [1].

There are various methods that can be used for measuring differences of height on a surface. The most important thing is to find some suitable 'setting criterion' that will give sufficient reproducibility for the measurements perpendicular to the reference plane. Most methods require physical contact with the surface. The force then exerted on the surface affects its shape, however, and may even damage it. Since we wish to determine the radius of curvature of test samples in the form of small strips only about 100 μm thick, which are bowed (or twisted or slightly saddle-shaped) after application of a reflecting coating, a non-contact method is desirable. One method would be to use a microscope with a wide-angle objective

lens, and to take the sharpness of the observed image as the setting criterion. This method does not give sufficient reproducibility, however, and it is not suitable for an automated measurement procedure. In an earlier arrangement we used a gas-laser interferometer system [2]. This is an incremental method, however, and has the disadvantage that a flaw in the surface under test may cause the loss of a measurement unit of half a wavelength. The measurement then has to be started all over again.

By applying some of the technology developed for the Compact Disc and LaserVision players, we have succeeded in designing a practical instrument that satisfactorily meets the requirements: fully automated measurement of surface profiles with radii of curvature between 3 m and 300 m, with no contact with the surface, over a test length of 50 mm. Our instrument uses the AlGaAs semiconductor laser described earlier in this journal [3] and one of the methods used for focusing the optical reading system on to the surface of a LaserVision disc [4]. Other special features of our design are air bearings [5] with a tapered air gap [6] and friction-wheel reduction mechanisms [7].

The title photograph gives a view of the instrument with its XY recorder, which plots the surface profile,

A. F. Foederer, J. L. M. Hagen and Ir A. G. van Nie are with Philips Research Laboratories, Eindhoven.

see fig. 1a. The principle of the instrument is illustrated in fig. 1b. The beam *B* is supported on the shafts *S*₁ and *S*₂ by the air bearings *B*₁ and *B*₂, permitting movement perpendicular to the plane of the drawing (in the *x*-direction). The shafts are attached to a granite

lens *L*₁. The second lens *L*₂ is mounted on the lens carriage *LC*, which can be moved to and fro by means of the rotating shaft *S*₄. This shaft is pressed against *LC* by a spring (not shown). If the surface of the workpiece *W* is below the rear focal plane of *L*₂, the

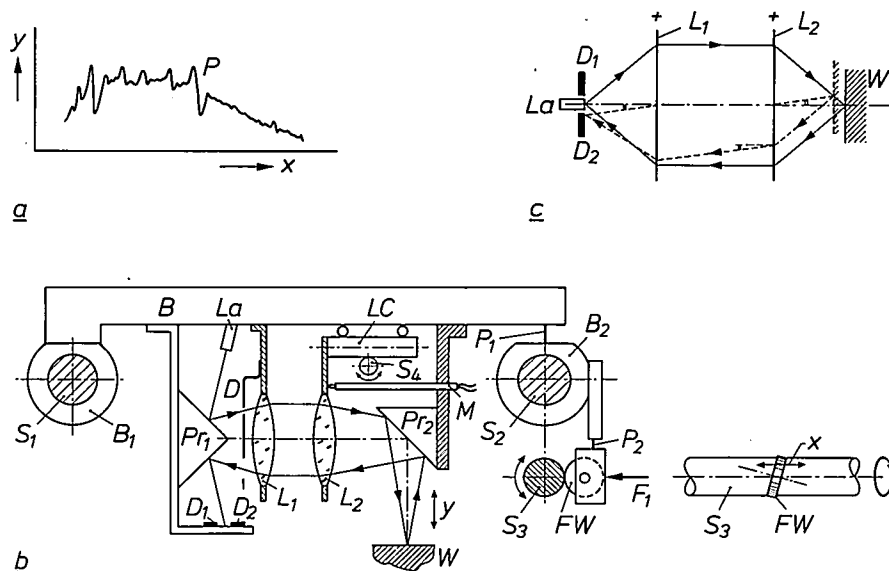


Fig. 1. The measuring instrument, shown schematically. a) The result of a measurement in the form of a surface profile *P* drawn by the XY recorder. b) Cross-section through the instrument, perpendicular to the *x*-axis, with a part of the right-hand side view shown on the right. *B*₁, *B*₂, air bearings for the movement of the beam *B* in the *x*-direction on the shafts *S*₁ and *S*₂. *P*₁ leaf spring; this is necessary since the distance between *B*₁ and *B*₂ and the distance between *S*₁ and *S*₂ may differ slightly. *S*₃ rotating shaft for driving *B* in the *x*-direction by means of the friction wheel *FW*. *P*₂ pivot permitting *FW* to be pressed against *S*₃ with a force *F*₁. *La* semiconductor laser. *D*₁, *D*₂ photodiodes integrated in a dual photocell. *D* diaphragm. *L*₁, *L*₂ lenses. *Pr*₁, *Pr*₂ prisms with reflecting surfaces for beam deflection. *W* workpiece under test. *LC* lens carriage. *S*₄ shaft for driving *LC* in a direction corresponding to the *y* movement. A control system, which takes the output currents of *D*₁ and *D*₂ as the measurement signal, causes *LC* to move so that the reflecting surface of *W* coincides with the rear focal plane of *L*₂. *M* inductive displacement gauge that measures the displacement of *LC*. c) Simplified diagram of the optical system. The dashed lines relate to the situation where the surface *W* is not coincident with the rear focal plane of *L*₂. The reflected beam then strikes one of the photodiodes, in this case *D*₂, and the dual photocell goes 'out of balance'.

baseplate by steel blocks. The beam can be displaced in the *x*-direction by the friction wheel *FW*, which is pressed against the rotating shaft *S*₃, since the axes of rotation of *S*₃ and *FW* are not quite parallel. The displacement of *B* is measured by an inductive displacement sensor (not shown). The output signal from this sensor is applied to the *x* input of the XY recorder.

As we noted earlier, the surface of the workpiece *W* is a reflecting surface. The principle of the optical system for setting in the *y*-direction can be seen from fig. 1b and 1c. The diaphragm *D* (in which the upper hole is smaller than the lower one) passes a narrow beam of the radiation (wavelength 800 nm) from the semiconductor laser *La*. The emitting surface of the laser and the photodiodes *D*₁ and *D*₂ — two photodiodes integrated on one silicon chip to form a dual photocell — are located in the front focal plane of the

reflected beam strikes the photodiode *D*₁; if the workpiece surface is above the rear focal plane, the beam strikes the diode *D*₂. The output currents of *D*₁ and *D*₂ are unequal in both cases. An electronic control system displaces *LC* in the direction that will make the two output signals equal. The displacement of *LC* is

[1] A. G. van Nie, A method of measuring mechanical stresses in passivation layers, Philips tech. Rev. 39, 130-133, 1980.
 [2] A. G. van Nie, A method for the determination of the stress in, and Young's modulus of silicon nitride passivation layers, Solid State Technology 23, No. 1, pp. 81-84, Jan. 1980.
 [3] J. C. J. Finck, H. J. M. van der Laak and J. T. Schrama, Philips tech. Rev. 39, 37, 1980.
 [4] G. Bouwhuis and J. J. M. Braat, Video disk player optics, Appl. Optics 17, 1993-2000, 1978.
 [5] I. C. Tang and W. A. Gross, Analysis and design of externally pressurized gas bearings, ASLE Trans. 5, 261-284, 1962.
 [6] J. L. M. Hagen, Aërostatistische lagers in de werktuigbouw, Mikroniek 21, No. 5, pp. 7-13, 1981 (in Dutch).
 [7] M. J. J. Dona, Constructie-elementen ten behoeve van micro-verplaatsingen, Mikroniek 21, No. 6, pp. 7-10, 1981 (in Dutch).

measured by an inductive displacement gauge *M*. The output signal is applied to the *y* input of the XY recorder. Since the reflecting surface of *W* is located in one of the focal planes of the lens *L*₂, the position of the 'image' of *La* on *D*₁ and *D*₂ does not in theory depend on the angle between *W* and the optical axis. For practical reasons, however, partly connected with the magnitude of the diaphragm apertures, *W* should not deviate too greatly from the perpendicular position. It is therefore not possible to measure radii of curvature smaller than 3 m.

There must be no play in the drive mechanism for the beam *B* and the lens carriage *LC*. Since conventional transmission systems are therefore unsuitable for this application, we decided to use friction-wheel transmission. The shaft *S*₄, which drives *LC*, is connected to a reduction mechanism, shown schematically in *fig. 2a*. This reduces the speed of the motor and its 'apparent backlash' to the desired low values. The friction 'wheels' in this case are three steel balls, denoted by *B*. The two inner tracks along which the balls roll form part of the input shaft *S*₅ and have the same radius *a*; see *fig. 2b*. The axes of rotation *I-I* and *II-II* of the two balls in *fig. 2a* are therefore parallel to *S*₅. The two outer tracks of the balls are located in the stationary part of the reduction mechanism and in a 'cup-shaped' part of *S*₄. If the radii *s* and *r* of the outer tracks were also equal, the output shaft would be stationary. The radii are not exactly equal, however, so that a low transmission ratio $i \approx 1/50$ is obtained. The entire system is prestressed by a force *F*₂ from a diaphragm spring (not shown).

The air bearings for the longitudinal movement have a tapered gap, as mentioned earlier. Aerostatic air bearings have long been known and are mainly used in precision lathes and high-precision measuring instruments. These bearings have two special advantages: their friction is low, and deviations in the shape of the shaft and bearing bore are effectively averaged. In a bearing with a parallel gap the tolerance on the gap height of 10 to 20 μm is no more than a few microns, and this tolerance must be maintained over the full length of the bearing bore. In a bearing with a tapered gap this tight tolerance is only necessary at the two ends of the bore. Bearings with a tapered gap are therefore easier to make, and have the added advantage that they are less sensitive to dirt and damage. *Fig. 3a* gives the dimensions of our air journal bearing. *Fig. 3b* gives a plot of the load *F*_b as a function of the displacement δ of the centre of the bearing, as the result of a number of measurements. It can be seen that with a load of 100 N and an air-supply pressure $p_s = 7$ bars the displacement of the bearing centre is only 3 μm. In our application the load does

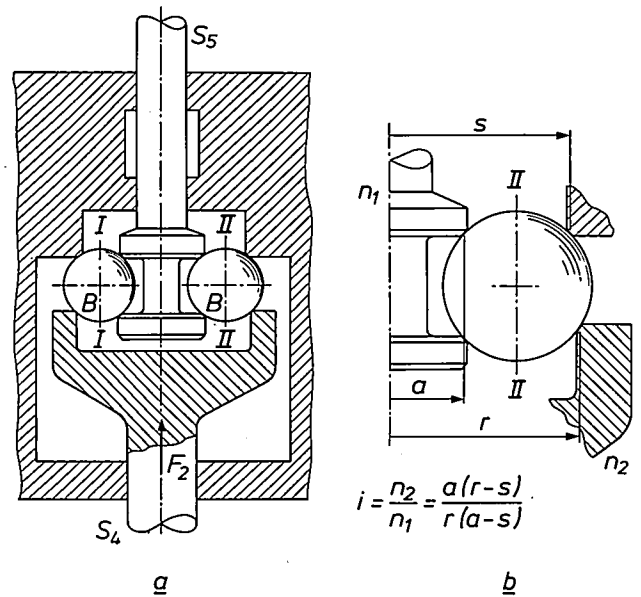


Fig. 2. *a*) The reduction mechanism for driving the lens carriage *LC* (see *fig. 1b*). *S*₅ input shaft (coupled to the motor), *S*₄ output shaft. *B* steel balls (there are three). *F*₂ prestressing force, applied by a diaphragm spring (not shown). *I-I* and *II-II* axes of rotation of the balls. *b*) Calculation of the transmission ratio *i*, the ratio of the speed of revolution *n*₂ for the output shaft and *n*₁ for the input shaft. The radii of the two inside tracks of the balls are both equal to *a*, the axes of rotation *I-I* and *II-II* are therefore parallel to the shaft *S*₅. The two outer tracks have radii *s* (the stationary track) and *r* (the moving track in *S*₄). The transmission ratio $i = n_2/n_1 = a(r-s)/r(a-s)$ can readily be calculated by assuming that the axis *II-II* is fixed and that the ball performs one revolution. If the values of *r* and *s* are close, a low value of *i* is obtained; in our case *i* is about 1/50.

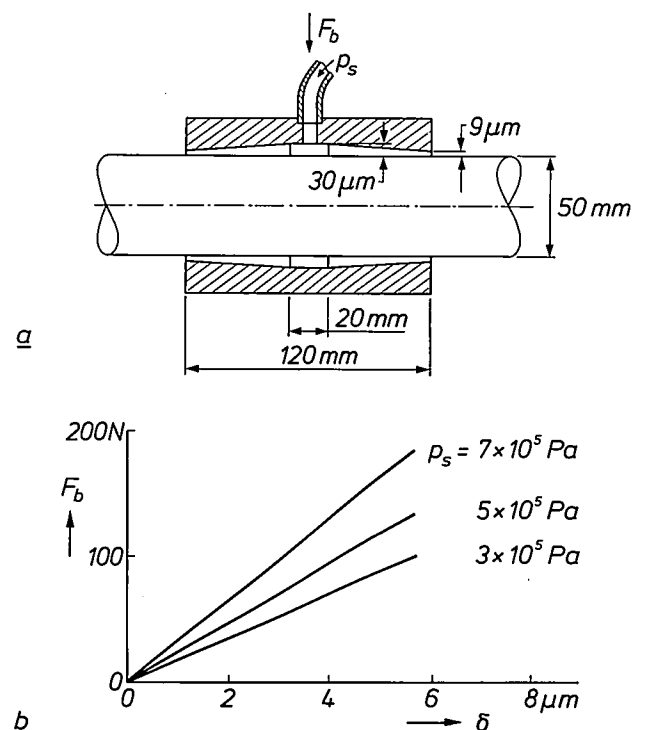


Fig. 3. The aerostatic air journal bearing with tapered gap for supporting beam *B* in *fig. 1b*. *a*) The dimensions of the bearing; *p*_s supply pressure, *F*_b load. *b*) Load *F*_b as a function of the displacement δ of the centre of the bearing, as the result of measurements made at three different supply pressures. The bearing has a high stiffness: with $p_s = 7 \times 10^5$ Pa the displacement at *F*_b = 100 N is only 3 μm.

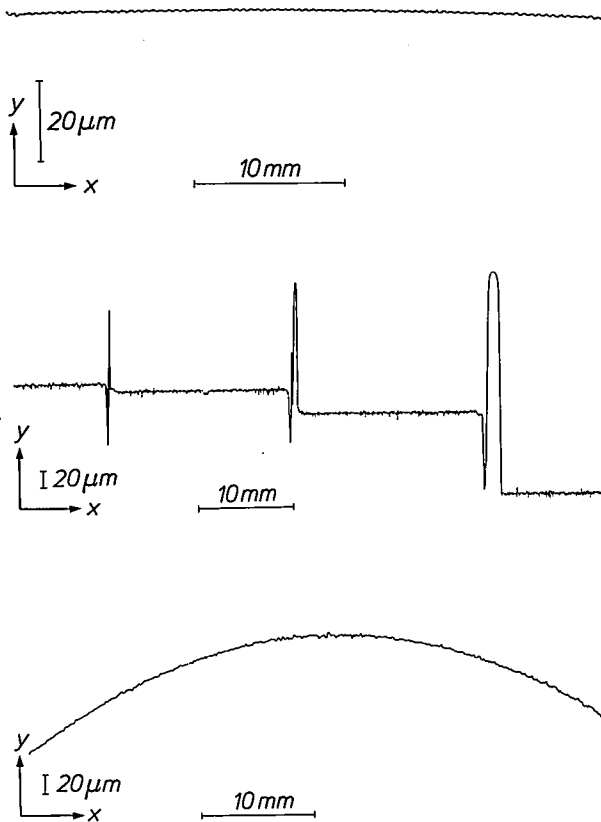


Fig. 4. The results of measurements on test objects. In the interpretation of the traces; the peaks can be ignored since they are due to particles of dust, scratches or sudden irregularities in the surface. When such disturbances are encountered, the reflected beam (see fig. 1b) is deflected too sharply and the instrument overshoots.

a) Measurement of the linearity of the x movement by means of an optically flat mirror. The measured deviation is $1.2 \mu\text{m}$ over a length of 40 mm . The dimensional error of the mirror as determined in the metrology laboratory with a Talysurf 5 instrument is $0.5 \mu\text{m}$ (in the same direction) over 50 mm . The deviation in the linearity of the x movement is thus about $0.7 \mu\text{m}$.

b) Measurement of a workpiece, made on a precision lathe, with step heights of $11 \mu\text{m}$, $31 \mu\text{m}$ and $104 \mu\text{m}$, measured in the metrology laboratory, again with a Talysurf 5. The results of the measurements give step heights of $11 \mu\text{m}$, $32 \mu\text{m}$ and $106 \mu\text{m}$. The maximum deviation is thus about $2 \mu\text{m}$.

c) Measurement of a curved mirror. Height measurements are made at regular distances along the trace ($16\times$). The radius of the circle of best fit was calculated with the aid of a computer program based on the Gauss-Newton method^[2]: the value was 2984 mm . The heights were also measured in the linear metrology laboratory with a Zeiss UMM500 instrument. The same program gives a radius of curvature of 3032 mm . The difference between the two results is 1.6% .

not change during the movement of beam B , so that the effect of the bearing system on the accuracy of the measurement is extremely small.

To test the measuring accuracy of our instrument, we determined the shape of three test objects. Their geometry was also measured in the metrology laboratory at Philips Research Laboratories. Fig. 4a gives the result of the measurement on an optically flat mirror. The deviation measured with our instrument is $1.2 \mu\text{m}$ over a test length of about 40 mm , and is the sum of the $0.5 \mu\text{m}$ dimensional error of the mirror, measured with a Talysurf 5 roughness-measuring instrument in the metrology laboratory, and the deviation of the linearity of the x movement, which is clearly about $0.7 \mu\text{m}$. Fig. 4b shows the result of a measurement on a workpiece with a stepped taper, specially made on the COLATH precision lathe^[8]. The step sizes measured with our instrument and those that were determined with the Talysurf 5 in the metrology laboratory differ by $2 \mu\text{m}$ or less. We estimate the inaccuracy of the 'rise' of surface profiles as measured with our instrument to be less than $\pm 2.5 \mu\text{m}$ per $100 \mu\text{m}$. Fig. 4c shows the result of the measurement on a mirror with about the smallest measurable radius of curvature. Its shape was also determined by height measurements with a three-dimensional Zeiss UMM500 instrument. The results from the two methods of determining the radius of curvature (caption, fig. 4) support our statement that the measurement of a radius of curvature of about 3 m with our instrument has an inaccuracy of less than $\pm 2.5\%$.

As we said at the beginning, the instrument described here has been used for determining the radius of curvature of test strips after applying a metal coating. The instrument can also be used for measuring the movement of the bimetallic strip in starters for fluorescent lamps. Further laboratory applications are expected.

^[8] T. G. Gijssbers, Philips tech. Rev. 39, 229, 1980.

The attachment of leadless electronic components to printed boards

R. J. Klein Wassink and H. J. Vledder

In view of the ever-increasing demands on the quality, compactness and price of electronic circuits, it is to be expected that increasing use will be made of leadless components, that is to say components without wire terminations. These have proved to be suitable not only for hybrid thin-film and thick-film circuits but also for printed circuit boards, which are still more widely used. Compared with components that have wire terminations, leadless components are more suitable for automatic mounting, and offer other advantages such as a higher packing density and greater reliability. At various places work is now under way on the technology required for attaching leadless components to printed boards. It is also the subject of extensive study at the Philips Centre for Manufacturing Technology, and in the article below we present a broad survey of the present state of the art.

The use of leadless electronic components

Electronic components like resistors, capacitors and inductors have for many years been made with wire terminations (leads), to provide the electrical connection to other components of an electronic circuit. Since the fifties it has become common practice to use printed wiring for the interconnection of components. Conductor patterns are formed on an insulating substrate of a phenolic paper laminate or an epoxy resin reinforced with glass fibre, and small holes are made in the substrate. The leads are inserted through the holes from one side of the substrate and soldered to the conductors on the other side^[1]. A substrate with holes and conductor patterns will be referred to here as a 'printed board'; other terms in use are 'printed-circuit board' or 'PC board'. *Fig. 1* shows a number of electronic components with leads. The components differ considerably in shape and dimensions, but one thing they have in common is that their leads must be inserted through holes in the board before the components are finally attached to it.

For some considerable time, manufacturers have been producing many components that do not have wire leads but are provided with short flat-based terminations or with metallized bonding pads; see *fig. 2*.

Ir R. J. Klein Wassink and Dr H. J. Vledder are with the Philips Centre for Manufacturing Technology, Eindhoven.

Examples are 'block-shaped' resistors and multilayer capacitors, cylindrical resistors and diodes, and transistors in an 'SOT' encapsulation^[2]. Until recently such components were used almost exclusively for hybrid circuits. These are produced by means of thin-film or thick-film technology, the conductor patterns and certain components being applied in integrated form to a substrate, often of aluminium oxide (Al_2O_3); the discrete components are then added. The substrate is usually no larger than 50×50 mm and the thickness of the applied film may be of the order of 10 nm in the case of thin-film technology^[3], rising to 0.2 mm for thick-film technology^[4].

In recent times the attachment of leadless components to printed boards of a laminate material (phenolic paper or glass-fibre-reinforced epoxy resin) has become increasingly common, since in many applications such components have distinct advantages over conventional types. Their small dimensions, for example, permit high horizontal and vertical packing densities, and this is a great advantage in equipment such as portable television cameras and car radios. The absence of leads is also advantageous for high-frequency applications. Furthermore, since the components are usually simple in shape, it is easy and inexpensive to mechanize the process of positioning or

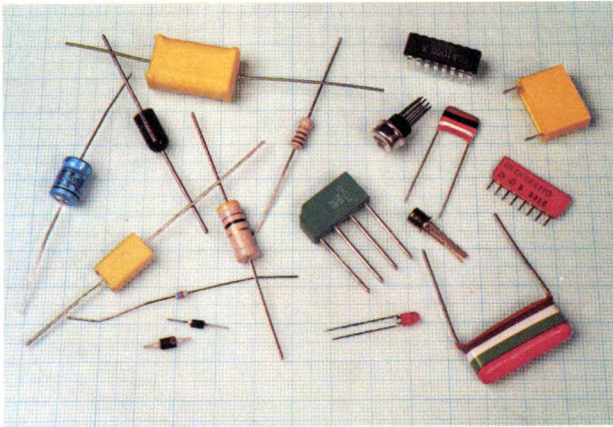


Fig. 1. Electronic components with lead wires for connection to a printed board.

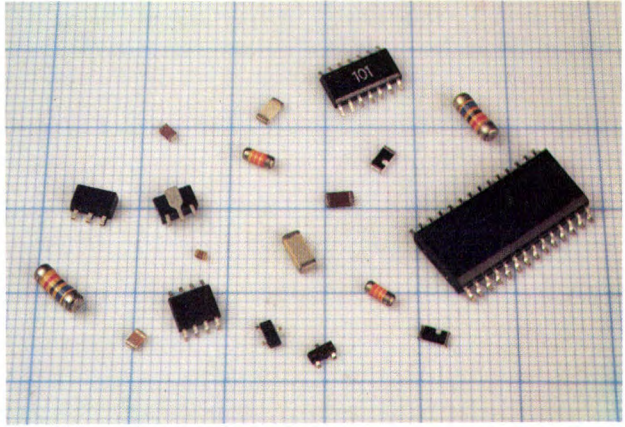


Fig. 2. Leadless electronic components with bonding pads or short terminations.

‘placing’ them on the boards, because no holes in the board are required and no wires have to be inserted through holes.

The exclusive use of leadless components is not yet possible, however. Some components such as inductors and electrolytic capacitors are not yet generally available in leadless versions, and other components, such as power resistors, are not really practical in leadless versions. Printed boards with only leadless components are therefore seldom used. As a rule the boards either have wire-terminated components alone — the conventional situation — or they have a combination of both types, in various possible configurations; see *fig. 3*.

The combination of components of both types imposes a significant constraint on the method of placing and attaching the leadless components. Before considering the attachment technology, we shall first look at some aspects of the placing of components.

Placing of components

Leadless components are very suitable for automatic placement on boards. For large numbers of components to be placed rapidly, they must be packaged in such a way that they can readily be fed to the printed boards. *Fig. 4* shows a machine for placing leadless components contained in tapes wound on reels. The number of components to be placed varies from one to many hundreds per board. The boards range in size from 40 by 100 mm to 100 by 250 mm (in some cases up to 200 by 300 mm). The required positional tolerance is of the same order as the tolerance observed in the application of conductor patterns, about 0.15 mm, so that the total tolerance is about 0.3 mm.

Various pick-and-place mechanisms have been developed that bring the components rapidly and ac-

curately to the appropriate positions on the board. The method used in the machine shown in *fig. 4* is illustrated schematically in *fig. 5*. A number of vacuum pipettes, which can all be rotated and moved in the horizontal and vertical directions, pick up the components from recesses cut into tapes that are fed in automatically.

There are three possible methods of finally attaching the leadless components to the printed boards. The components can be placed on small droplets of an electrically conducting adhesive applied to the conductors beforehand, and the adhesive can then be ‘cured’ (thermally hardened, thermosetting). The adhesive is generally an epoxy resin, filled with metals such as gold and silver. This method is used successfully in thick-film technology, but not for the mass production of boards with leadless components. The components can also be placed in solder paste on the conductors, and the paste then melted, e.g. by heating

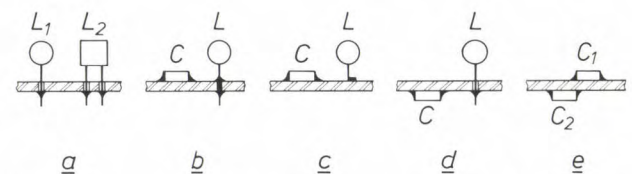


Fig. 3. Different configurations on a supporting panel. *a*) Conventional situation: components L_1 and L_2 with lead wires through the holes in the board. *b*) Combination of leadless component C and conventional component L , on the same side of the board. L is connected with the other side through a metallized hole. *c*) L and C on the same side, as in thick-film technology [4]. *d*) L and C on different sides. *e*) Leadless components C_1 and C_2 on different sides.

[1] See for example R. J. Klein Wassink, *Soldering in electronics*, Electrochemical Publications Ltd., Ayr, Scotland, appearing shortly.
 [2] Block-shaped leadless components are generally referred to as chips, while cylindrical components with metal end caps are referred to as MELF components (metal electrode face bonding). In this article we are concerned with chip components.
 [3] See for example E. C. Munk and A. Rademakers, *Philips tech. Rev.* 27, 182, 1966.
 [4] See for example W. Funk, *Philips tech. Rev.* 35, 144, 1975.



Fig. 4. Philips placing machine for multiple assembly of taped leadless components. The board feed, the application of adhesive to the boards, the transfer of components from tape to boards and the transport and removal of the assembled boards are all automatic. The machine has a capacity of about 250 000 components per hour. *Right:* Detail of tapes with components.

it with hot air, a laser beam or a condensing vapour ('reflow soldering'). This method is used for a variety of professional applications. A third method is placing in droplets of non-conducting adhesive *between* the conductors. The adhesive is then cured and the electrical connections are made by soldering. This method is highly suitable for mass production with leadless and conventional components combined, and is therefore the subject of keen interest. This is the method we shall mainly be concerned with in the present article.

Conventional components are placed manually or with an automatic placement machine, depending on the number of components, the size of the printed boards and the batch size. A placement machine inserts the wires through the appropriate holes in the board, bends them over and cuts them to the right length. The bending and cutting process requires a

relatively large area around the holes to be free of leadless components.

For products using components of both types it is necessary to decide which are to be placed first. If conventional components are placed first, there are no problems in the automatic insertion of the lead wires and the area around the bent and cut wires can be used later for the leadless components. The wires already present do however limit the possibilities for applying the adhesive. If the leadless components are placed first, there are no restrictions on the application of the adhesive, but there are now limitations on the insertion, bending and cutting of the wires. The order to be preferred will depend on the ratio of the numbers of the two types. At present it is usual to start with the conventional components, but as leadless components become more widely used it will be desirable to place them first.

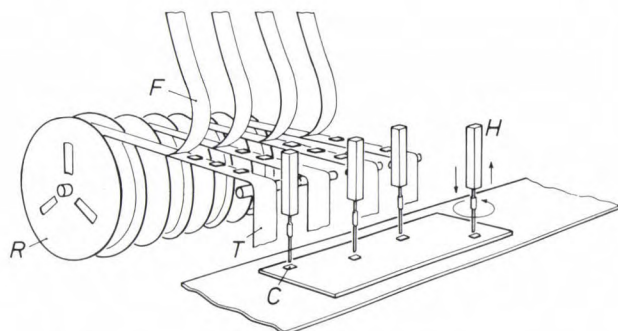


Fig. 5. Diagram of a method for automatic placement of leadless components on a printed board. *H* vacuum pipettes, which can be displaced horizontally or vertically. *R* reel. *T* tapes with leadless components *C*, which can be picked up after stripping off the plastic strips *F*.

Temporary bonding

For soldering, the leadless components are first bonded temporarily to the printed board. An adhesive is used that hardens by polymerization (curing). We shall first consider the application of the adhesive and then the subsequent curing.

Application of adhesive

The diagrams in *fig. 6* show three methods that can be used for applying the adhesive: screen printing, dispensing (with a syringe needle) and pin transfer. For screen printing a perfectly flat substrate is required, and therefore this method cannot be used when con-

ventional components are already on the board. Dispensing is only really feasible if fewer than about 5000 components per hour have to be placed. If many more components per hour (e.g. 100 000) have to be placed, then so many drops of adhesive are required that one dispenser is no longer sufficient. The use of several dispensers side by side can easily result in non-uniform droplets and thus cause all kinds of practical problems.

With the pin transfer method a stationary state is soon reached in which equal amounts of adhesive are taken up and deposited. At the locations where the pins pick up the adhesive the thickness of the layer of adhesive has to be strictly uniform. The disturbed layer is therefore smoothed flat with a squeegee after each pick-up. Surface tension and the force of gravity will cause a deposited drop of adhesive to take the form of a flattened cap of a sphere; but this effect is counteracted by the increase in viscosity resulting from the increase in 'structure' in the adhesive after the disappearance of the mechanical forces upon deposition. The rate at which the viscosity increases determines the extent to which the spherical cap shape is reached. This also implies that the ultimate shape of the droplet depends to some extent on the movement of the pin. The above process is illustrated in *fig. 7* with a few diagrams based on video images. The rheological behaviour of the adhesive can be influenced by the addition of various filler substances.

The amount of adhesive must be sufficient to bridge the gap between board and leadless component, but it must not be so large that the adhesive flows out over the conductors, impeding the soldering. In the bonding of components such as resistors and capacitors the height of the drop of adhesive should be greater than the sum of the thickness of the conductors on the board and that of the metallization on the underside of the component; see *fig. 8*. For the SOT-23 encapsulation the height of the adhesive droplet should be greater than the sum of the thickness of the conductors and the 'leg height'.

The thickness of the conductors may vary considerably. It may be between 30 and 37 μm for conductors produced by photo-etching, and between 50 and 200 μm for conductors on boards with metallized holes after a 'reflow' treatment. The thickness of the metallization for the various components is usually only about 20 μm . The leg height of the SOT-23 encapsulation, on the other hand, varies between 100 and 200 μm , since these components were developed for thick-film technology. The variation in the thickness of the conductors and, in the case of the SOT-23, in the leg height, entails a wide variation in the amount of adhesive required. For mass-production

applications it is of course desirable to use the same amount of adhesive. For this reason the conductors should have the same thickness everywhere, and should preferably be as thin as possible. The components should also be carefully matched to one another.

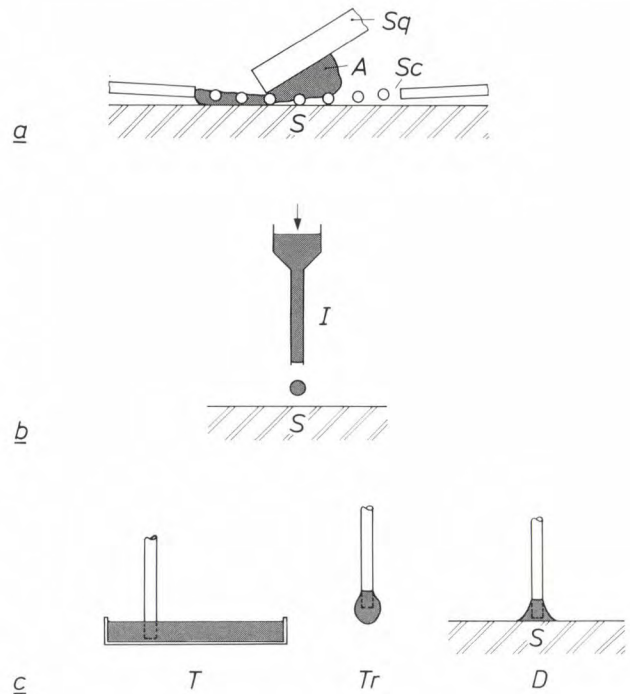


Fig. 6. Three methods for applying adhesive to a printed board. *a)* Screen printing, where the adhesive *A* is pressed through a screen *Sc* on to the substrate *S* by a squeegee *Sq*. *b)* Dispenser, where the adhesive is applied to the substrate by a pulse of compressed air through a syringe needle *I*. *c)* Pin transfer, consisting of the take-up *T*, the transfer *Tr* and the deposition *D* of the adhesive.

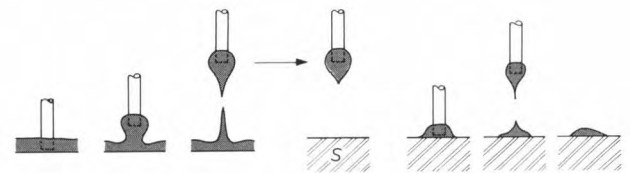


Fig. 7. Different stages in the pin transfer method. The diagram was made with the aid of high-speed video recordings.

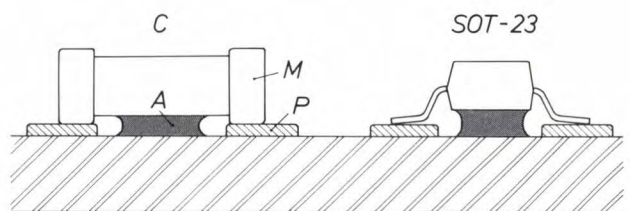


Fig. 8. Diagram of a printed board on which a leadless component *C* and an SOT-23 encapsulation have been bonded. *A* adhesive. *M* metallized bonding pads. *P* conductor patterns.

other in dimensions, so as to ensure, for example, that the SOT-23 leg height does not deviate too much from the thickness of the metallization on resistors and capacitors. For SOT encapsulations this entails a reduction of the leg height to below 100 μm .

Apart from the methods mentioned, there are other ways of applying adhesive. It can for instance be dispensed on the components (not on the bonding pads or terminations) instead of the printed boards. One method is to use pieces of double-sided adhesive crepe tape, or a two-component adhesive, applying one component to the electronic component and the other to the board. It is also possible to apply adhesive to the components before or during the packaging. Instead of adhesives that require to be cured, pressure-sensitive adhesives can be used. We shall not consider these methods further here.

Curing the adhesive

After the adhesive has been applied and the leadless components placed on it, the adhesive is left to harden by polymerization. Before this curing process the adhesion must be sufficient to withstand the effects of movement of the printed board. After curing, the adhesive bond must be strong enough to withstand the forces imposed at room temperature by any further addition of conventional components. Furthermore, the components must not be loosened by the thermal shock caused by the transition to the soldering temperature. The adhesives generally used are those that are cured by heating (thermosetting) or by ultraviolet irradiation. The difference between these two processes can be seen via the viscosity behaviour of the adhesive during the curing process; see *fig. 9*. In thermosetting the viscosity first decreases as a result of the rise in temperature. This could cause some of the adhesive to flow away, resulting in a weaker bond. An increase does not occur until the polymerization is well under way. For curing by ultraviolet irradiation the viscosity increases immediately.

The method of application imposes requirements on the rheology of the adhesive and hence on the composition. When thermosetting epoxy resins are used, slow polymerization starts at room temperature, so that the adhesive gradually becomes unusable. The time during which the adhesive can readily be worked is called the 'pot life'. An adhesive with a long pot life usually requires prolonged curing at high temperature. With some components, unable to withstand such prolonged heating, this can cause problems. The temperature has a considerable effect on the progress of the curing; see *fig. 10*. A high degree of curing not

only gives a firm bond but also ensures good resistance to moisture, thus helping to prevent corrosion. About half an hour is required for adequate curing at 80 to 100 °C of epoxy resins with a reasonable pot life (5 to 10 hours). A significant speeding up of the process is only possible by increasing the curing temperature.

For adhesives that cure well under ultraviolet irradiation it is general practice to use acrylates with an added photo-initiator [5]. After absorption of ultraviolet light (at about 350 nm) the initiator molecules dissociate to form radicals that start the polymerization. For good curing the adhesive must be sufficiently accessible to the radiation. During exposure, therefore, sufficient adhesive must protrude from beneath the components being bonded. It is also possible that the chain reactions occurring during polymerization may cause some curing of unexposed adhesive. The heat generated with a high irradiation intensity may also contribute to the polymerization. Curing with ultraviolet light can easily be completed in about 10 seconds on a conveyor belt.

Soldering

After the adhesive has been sufficiently cured, the leadless components are soldered to the board, usually with an alloy of 60% tin and 40% lead (melting point about 185 °C). The bonding pads of the components

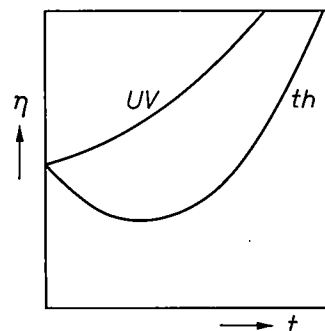


Fig. 9. Viscosity η of an adhesive as a function of time t (schematic) for thermal curing (*th*) and for curing with ultraviolet light (*UV*).

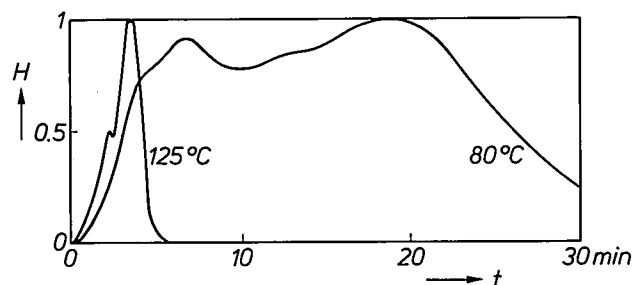


Fig. 10. Measured heat of reaction H (in arbitrary units) plotted against the time t for thermal curing of an adhesive, at 80 °C and 125 °C. The maximum of both curves is normalized to 1.

are thus connected electrically via the solder to the conductors on the board; see *fig. 11*. At the same time the conventional components are also soldered to the board. The joints required are obtained by wave soldering or drag soldering to bring the bonding pads and lead wires together with the conductors into contact with molten solder^[1]. In wave soldering this is done by pumping up a 'fountain' of solder, with the printed board travelling approximately horizontally; see *fig. 12*. In drag soldering the contact is produced by moving the board along a stationary bath filled with molten solder. The advantage of these methods is that the bonding pads and lead wires are heated quickly during soldering without the components becoming too hot^[6].

The different stages of the process take place automatically in a soldering machine. First, the flux required to ensure good wetting is applied. This is usually a resin dissolved in an organic solvent, such as isopropanol. Next the flux is dried until it has the desired viscosity, and at the same time the boards with the components to be soldered are preheated to permit faster wetting. When this has been done, the actual soldered joints are made.

Soldering is carried out at a temperature of 240 to 270 °C and the contact time is between 2 and 4 seconds. The components moved through the molten solder must be able to withstand this treatment, of course. This implies that they must be virtually unaffected by the high temperatures and the temperature gradients, i.e. that their characteristics are hardly affected at all. Moreover, the bonding pads and wire terminations must obviously not dissolve too greatly in the molten solder. It has been found that the components to be soldered generally meet these requirements satisfactorily.

Possible defects

After soldering leadless components various types of defects may appear. For example, a component may become displaced or even removed. The shifting of a component is not usually a direct result of the force exerted upon it by the flowing solder. A simple calculation shows that this force is only of the order of 10^{-4} N, which is at least a factor of 10^3 smaller than the force required to detach the component from the board. It is much more likely that the adhesive joint is broken because of thermal stresses, due to the contact with the hot solder.

Another possible fault is the absence of solder on a contact surface. This occurs when the flow of the solder is impeded to such an extent that it does not make contact with the solder land or the component termination. This obstruction may be caused by an un-

favourable shape of the components. It may occur, for example, when the short terminations of a transistor in the SOT-23 encapsulation lie in the shadow of the component. A further obstruction arises when gas generated by the flux solvent is trapped between the solder and the board. To some extent this can be avoided by making holes in the board to allow the gases to escape or by designing a conductor pattern in such a way that the supply of solder is improved by the use of particular patterns. It is evident that the supply of solder is impeded when part of the adhesive employed has flowed over on to a contact surface.

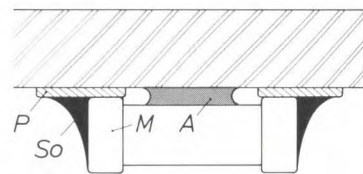


Fig. 11. Diagram of a printed board on which a leadless component has been bonded and soldered. *A* adhesive. *M* metallized bonding pads. *P* conductor patterns. *So* solder.

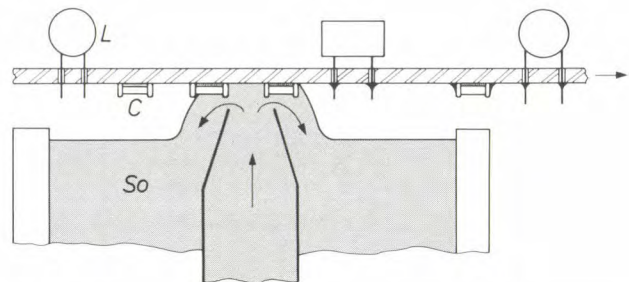


Fig. 12. Diagram of the wave soldering of leadless components *C* and conventional components *L* on a printed board. By pumping up the solder *So* and moving the board horizontally, the various contact surfaces and wire terminations of the components are bonded in turn to the conductors on the board.

A possible defect is also the occurrence of short-circuits between components situated close together. These short-circuits may be the result of incorrect process conditions, but usually they are due to an erroneous board design in which the distance between two components has been made too small relative to their height.

The last defect mentioned here is a bad soldered joint due to an insufficient amount of solder. This can happen if the process conditions and the board design are not properly matched, or if the metal of a contact surface cannot readily be wetted with solder.

^[5] The polymerization of acrylates has recently been described in this journal by J. G. Kloosterboer, G. J. M. Lippits and H. C. Meinders, *Philips tech. Rev.* **40**, 298, 1982.

^[6] R. J. Klein Wassink, *Philips tech. Rev.* **38**, 135, 1978/79.

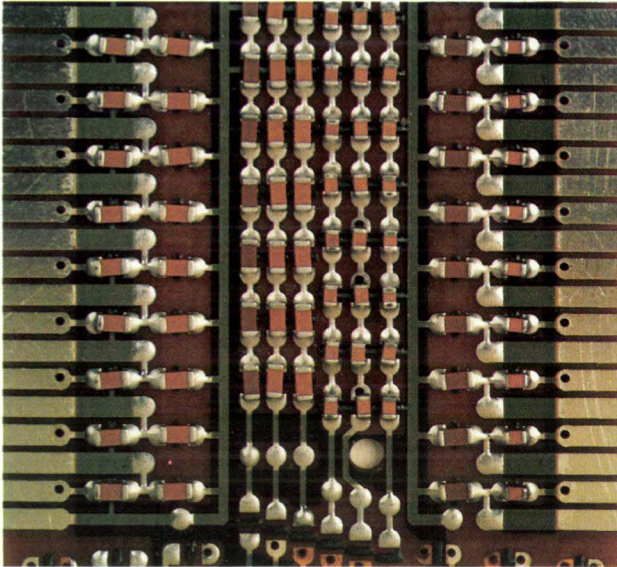


Fig. 13. Part of an experimental printed board of a phenolic-paper laminate (1.6 mm thick) with leadless components. The board was used for mechanically testing soldered joints and for investigating various configurations and conditions during bonding and soldering.

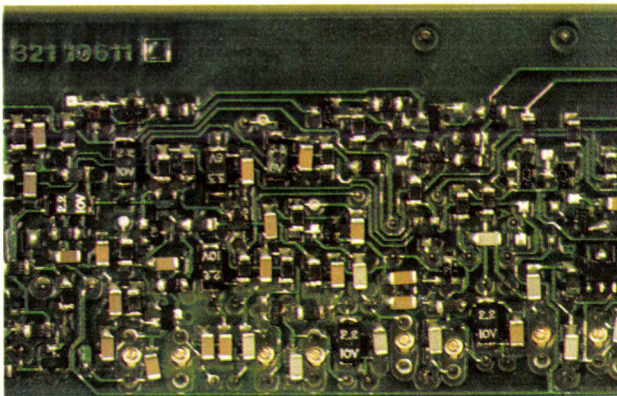


Fig. 14. Part of a double-sided printed board of glass-fibre-reinforced epoxy-resin laminate, with a large number of leadless components attached to the board by the methods described in this article.

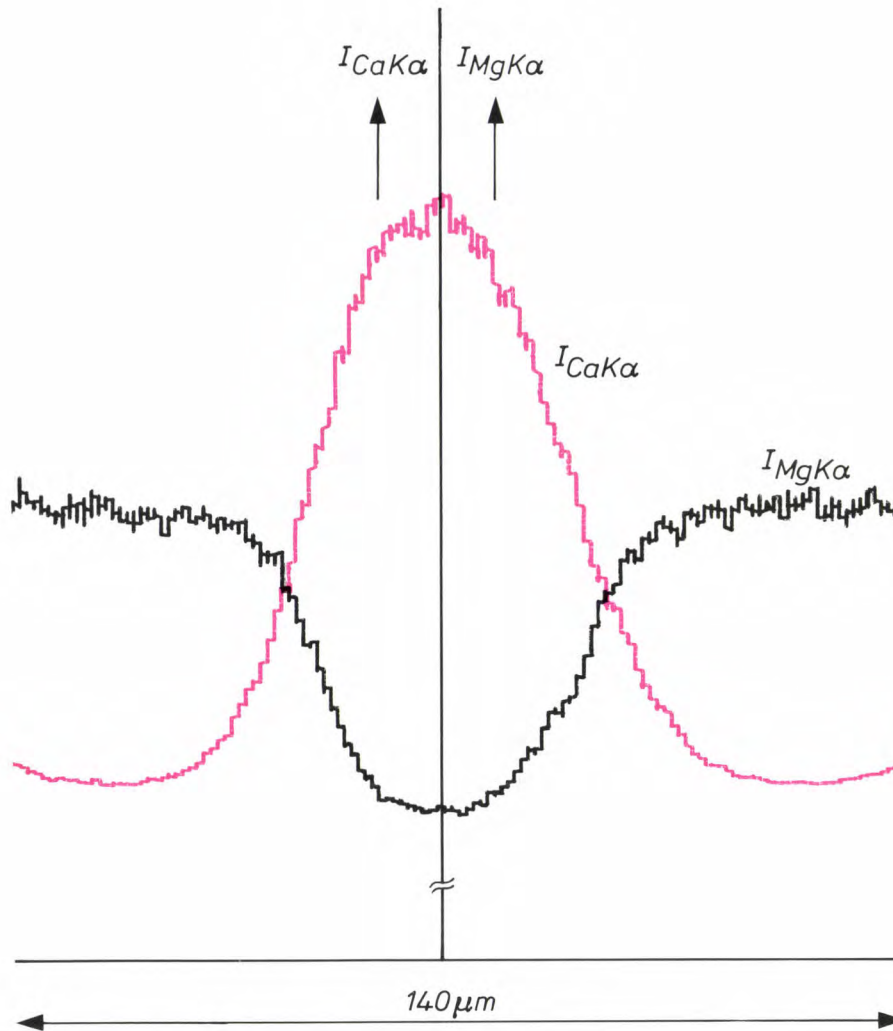
Quality of the soldered joints

With a suitable choice of materials, process conditions and pattern design, the defects mentioned above can mostly be avoided. The soldering of leadless components then gives reliable joints that will survive all kinds of thermal and mechanical tests. *Fig. 13* shows part of an experimental board used for mechanical tests on the soldered joints. Boards of this type have also been used for investigating various pattern configurations, as well as the effect of conditions during adhesive bonding and soldering.

The soldered joints are little affected by mechanical forces. They are relatively insensitive to vibration and will survive a severe bending test in which a board with components soldered to it is bent in a specified manner, and also various temperature-alternation tests. Intuitive objections to the attachment of rigid components to a relatively flexible material appear to be unfounded, provided the components are small. It may indeed be the difference in stiffness that makes the soldered joints so reliable mechanically. Flexible terminations are necessary for large components.

As an example of the many results of the attachment technology described here, *fig. 14* shows part of a printed board with a large number of leadless components.

Summary. The mounting of leadless electronic components on printed boards offers distinct advantages. They can easily be 'placed' automatically and they permit a high packing density. Since they are nearly always used in combination with components with wire terminations, a special technology is required for attaching them to the boards. The conventional components are usually placed on the boards first, and drops of adhesive are then applied at the appropriate places on the areas left free. This is preferably done by a pin transfer method. After the leadless components have been placed on the adhesive, the adhesive is cured by heating or by ultraviolet irradiation. All the components are then soldered to the board simultaneously by wave soldering or drag soldering. Defects such as shifting or disappearance of a component, no solder or too little solder and short-circuits between components can largely be avoided by a suitable choice of materials, process conditions and board design.



Electron microprobe analysis of a glass fibre

To determine whether glass fibres meet the requirements for optical communication it is necessary to use sensitive methods of analysis such as electron microprobe analysis^[1]. The figure shows the result of the analysis of a graded-index optical fibre made by the double-crucible system^{[2][3]}. The core has a diameter of about 50 μm and the cladding has an outer diameter of 140 μm . The original core material is a silicate glass with a refractive index of 1.530, containing 62.5 mol.% of SiO_2 and 12.5 mol.% of CaO ; it contains no MgO . The original cladding material is also a silicate glass, but with a refractive index of 1.515, containing 62.5 mol.% of SiO_2 and 12.5 mol.% of MgO ; it contains no CaO . The graded index is produced during glass transport in the double crucible by diffusion of Ca^{++} ions towards the cladding and diffusion of Mg^{++} ions towards the core. To determine how far the diffusion process has proceeded, the glass fibre is scanned radially in steps of about 1 μm by an

electron beam with a spot diameter of about 1 μm . The electron bombardment causes the material to emit X-rays. The black line in the figure gives the intensity $I_{\text{MgK}\alpha}$ of the $\text{MgK}\alpha$ radiation, and the red line gives the intensity $I_{\text{CaK}\alpha}$ of the $\text{CaK}\alpha$ radiation, both as a function of the radius. The measured values are approximately proportional to the gradients of the Mg and Ca concentrations. The zero levels of the curves are different because the peaks of the $\text{K}\alpha$ radiation are locally superimposed on a continuous 'Bremsstrahlung' background X-ray spectrum.

^[1] M. Klerk, The electron microprobe, Philips tech. Rev. 34, 370-374, 1974.

^[2] H. M. J. M. van Ass, P. Geitner, R. G. Gossink, D. Küppers and P. J. W. Severin, The manufacture of glass fibres for optical communication, Philips tech. Rev. 36, 182-189, 1976.

^[3] H. H. Brongersma, C. M. G. Jochem, T. P. M. Meeuwse, P. J. W. Severin and G. A. C. M. Spierings, The preparation of alkali-germanosilicate optical fibres using the double crucible system, Acta Electronica 22, 245-254, 1979.

A refrigerator-freezer with heat pipe

G. A. A. Asselman and A. J. van Mensvoort

Until recently, combined refrigerator-freezers with the 'four-star' rating have required two separate refrigeration circuits each with its own compressor, for optimum operation. The authors have succeeded in making one of the two compressors unnecessary by including a heat pipe that transfers heat from the refrigerator compartment to the refrigeration circuit of the freezer compartment. The problem of temperature control in the refrigerator compartment has been solved by using an adsorbant in the heat pipe, combined if necessary with an inert buffer gas.

Introduction

ISO standards specify that a combined domestic refrigerator-freezer with a 'three-star' rating should be able to maintain a temperature lower than $-18\text{ }^{\circ}\text{C}$ in the freezer compartment. At the same time the temperature in the refrigerator compartment should remain between 0 and $5\text{ }^{\circ}\text{C}$. If the refrigerator-freezer has four stars, this means that food can be deep-frozen in the freezer compartment. The refrigeration circuit that extracts heat from the freezer compartment, and is normally controlled by a thermostatic 'on/off' switch, must then be continuously switched on. During deep-freezing, the temperature in the refrigerator compartment must also remain between 0 and $5\text{ }^{\circ}\text{C}$.

Until recently it has only been possible to meet these conditions by providing separate refrigeration circuits for the refrigerator compartment and the freezer compartment, each circuit with its own compressor and its own temperature control; see *fig. 1a* and *b*. Refrigerator-freezers do exist with one refrigeration circuit containing two evaporators in cascade, but here only the temperature in the refrigerator compartment is controlled; see *fig. 1c*. Deep-freezing is not possible with this arrangement, since the temperature in the refrigerator compartment would become too low with the refrigeration circuit in continuous operation. Since the temperature in the freezer compartment is not controlled, the temperature in this compartment is usually much lower than necessary. In principle, it is

also possible to have two evaporators in parallel in a single refrigeration circuit: the flow of the refrigerant must then be separately controlled for each evaporator; see *fig. 1d*. However, this requires the use of electromechanically operated valves, and these make the system less reliable and relatively expensive.

We have found a better solution to the problem by using the principle of the heat pipe^{[1][2]}. There is just one single-compressor refrigeration circuit (called the primary refrigeration circuit); see *fig. 1e*. A temperature sensor in the freezer compartment switches the compressor on and off. The heat pipe has the task of transferring heat from the refrigerator compartment to the evaporator E_F of the primary circuit. A heat pipe is able to transfer heat because a liquid evaporates in one part of the pipe and condenses in another. In most heat pipes the condensed vapour is returned by means of a 'wick', which is a porous layer (usually of gauze) located close to the outer wall of the pipe. In the application described here the condensed liquid is returned by gravity. The part E_R where the evaporation takes place is situated at the bottom of the heat pipe, in the refrigerator compartment. The part C_R where the condensation takes place is situated at the upper end of the heat pipe, at the place where it is in contact with the primary evaporator E_F . The working fluid in the heat pipe is the same as that in the primary refrigeration circuit, i.e. CCl_2F_2 ('Freon 12').

In theory our method gives a lower efficiency (the ratio of the heat to be removed to the energy used in removing it) than if a second compressor were to be

Ir G. A. A. Asselman is with Philips Research Laboratories, Eindhoven; A. J. van Mensvoort, formerly with these laboratories, is now with the Philips Plastics and Metalware Factory, Eindhoven.

used, since the theoretical maximum efficiency, the Carnot efficiency, is inversely proportional to the temperature difference through which the heat must be transferred. Since the heat to be extracted from the refrigerator compartment is brought to the temperature level of the primary evaporator E_F , this temperature difference is greater in our case. This undesirable

sorbant in the buffer space, the quantity of inert gas in the heat pipe is varied, since the adsorbant — which is selective — takes up inert gas but not the refrigerant. In this way it is possible to control the effective area of the condenser and hence the heat transfer in the heat pipe. In the second system no inert gas is used, and the heat pipe only contains the refrigerant. The amount

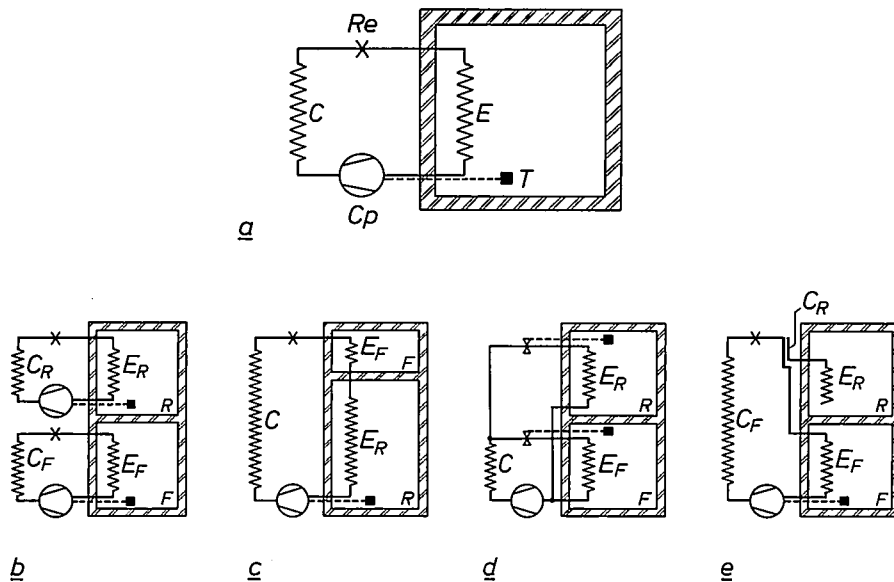


Fig. 1. a) The conventional refrigeration circuit for domestic refrigerators. C_p compressor, C condenser, E evaporator, Re capillary restriction, T temperature sensor. The refrigerant, usually a compound of fluorine and carbon ('Freon'), evaporates in E ; the vapour is compressed by the compressor and condenses in C ; the pressure of the condensate is reduced in Re ; the fluid finally enters the evaporator. The temperature sensor switches the compressor motor on and off as required. b) Combination of refrigerator compartment R and freezer compartment F . Each has its own refrigeration circuit. E_R and E_F evaporators, C_R and C_F condensers. c) Combined system using a single refrigeration circuit with one condenser C and two evaporators E_R and E_F in cascade. The temperature is measured in R , but not in F . In this system the temperature in F is therefore usually lower than necessary and deep-freezing is also not possible. d) Refrigerator-freezer with one refrigeration circuit with two evaporators E_R and E_F in parallel. The temperatures in both F and R are measured. The flow of refrigerant in the two evaporators is controlled by means of two electromechanical valves. e) Refrigerator-freezer with one refrigeration circuit (the primary refrigeration circuit) with compressor. The heat from R is conducted to the evaporator E_F of the primary circuit by a heat pipe with evaporator E_R and condenser C_R .

effect is more than compensated, however, because a single compressor (of greater capacity) has a higher efficiency than two separate compressors. In practice, therefore, a refrigerator-freezer with our system uses the same amount of energy — or a little less — than a system with two compressors. The cost of materials is also lower with only one compressor.

One problem is the control of the temperature in the refrigerator compartment R . Electromechanical valves cannot be used, since they are rather expensive and insufficiently reliable. The two systems that we shall now describe ultimately proved to be the most useful. The first system, discussed earlier in this journal [2], uses a quantity of inert gas, N_2 or CO_2 , contained in a buffer space connected with the heat pipe. By controlling the temperature of a quantity of ad-

of active refrigerant is controlled by adsorbing some of it in a different type of adsorbant. Here again, the heat transfer is controlled by varying the temperature of the adsorbant. If a large amount of refrigerant is adsorbed, the remainder cannot reach the bottom of the evaporator because it has completely evaporated half-way down. The evaporator then only operates over part of its surface.

In the following we shall first discuss the operation of the heat pipe in this application. We shall then deal with the two methods of controlling the temperature in the refrigerator compartment.

[1] G. A. A. Asselman and D. B. Green, Heat Pipes, I. Operation and characteristics, Philips tech. Rev. 33, 104-113, 1973.
 [2] G. A. A. Asselman and D. B. Green, Heat Pipes, II. Applications, Philips tech. Rev. 33, 138-148, 1973.

The heat pipe

The heat pipe without the control system is shown in the diagram of *fig. 2*. The evaporator E_R is at the bottom and the condenser C_R is at the top. The condenser is in contact with the evaporator E_F of the primary refrigeration circuit that cools the freezer compartment. The refrigerant condenses in C_R and is returned by gravity to the evaporator, so that the heat pipe does not require a wick, as mentioned earlier.

The actual appearance of the heat pipe is shown in *fig. 3*. Since the heat-transfer coefficient for the transition from the evaporator to the air in the refrigerator compartment is relatively small, the evaporator is part of the large evaporator plate EP , which is attached to the rear wall of this compartment.

The heat pipe is a means of transferring large quantities of heat across a small temperature difference. The heat pipe in this application, with an inside diameter of 6 mm and a length of nearly 3 m, can transfer up to 70 W with a temperature difference of no more than 12 °C between condenser and evaporator. If the same amount of power was transferred by a solid copper rod of these dimensions the temperature difference between the two ends would have to be about 9000 °C.

A heat-transport capability of 70 W is amply sufficient. For a refrigerator compartment of volume 160 litres — the volume for which our system was

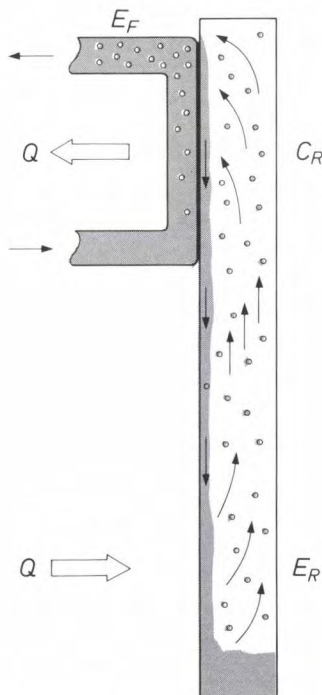


Fig. 2. The heat pipe without control system, schematic. E_F evaporator of the primary refrigerating system for the freezer compartment. C_R condenser of the heat pipe. E_R evaporator of the heat pipe. Q heat flow through the heat pipe, from the refrigerator compartment to the primary evaporator.

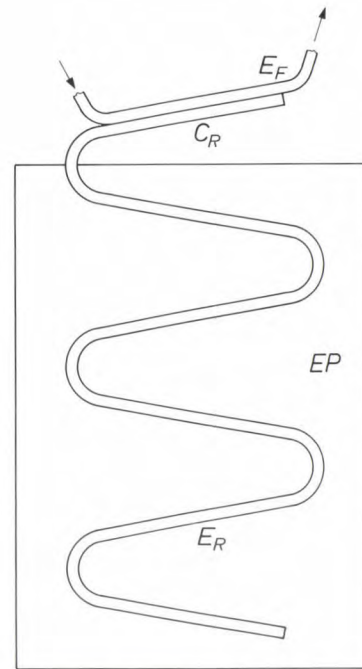


Fig. 3. Construction of the heat pipe. EP evaporator plate, incorporating the evaporator E_R . The pipe is nearly 3 m long and has an inside diameter of 6 mm. See caption to *fig. 2*.

designed — about 1 W of refrigeration capacity is required per degree of temperature difference between the ambient temperature and the temperature inside the refrigerator compartment. For a maximum ambient temperature of 32 °C and a temperature of 5 °C in the refrigerator compartment, the heat pipe therefore has to provide 27 W of heat-transport capability. There is thus a sufficient surplus of refrigeration capacity to compensate for losses when the door is opened, and to permit the required temperature to be reached within a reasonable time after switching on.

The limits to the heat-transport capability of a heat pipe are connected with flow effects and the quantity of refrigerant in the pipe. This is illustrated in *fig. 4*, which gives the results of a number of measurements on the heat pipe in *fig. 3*. The maximum heat-transport capability Q_{max} was measured for different quantities M_F of the refrigerant. It is found that the heat-transport capability is limited by three distinct effects. With normal steady-state flow in the pipe, a closed-loop integral for the pressure differences must have a value of zero. Neglecting the pressure differences in the radial direction, we find that the closed-loop integral over the length l of the pipe, inclined at an angle α , is given by the expression

$$g \sin \alpha \int_0^l \rho_x dx + \int_{x=0}^{x=l} dp_L + \int_{x=l}^{x=0} dp_v = 0, \quad (1)$$

where g is the acceleration due to gravity, ρ_x the density of the liquid, p_L the pressure in the liquid, p_v

the vapour pressure and x the position in the pipe, with $x = 0$ for the condenser end. If equation (1) cannot be satisfied, because there is too little fluid present, an equilibrium arises for part of the pipe: the liquid does not reach the bottom end of the evaporator. The capacity of the evaporator decreases, and so therefore does Q_{max} , as shown by curve 1 in fig. 4.

If, on the other hand, there is too much refrigerant in the pipe, then again no steady-state equilibrium as given by (1) is reached: the level of the liquid is raised by the vapour bubbles. The condenser becomes obstructed by the liquid and the heat-transport capability again decreases: curve 3. The curve is shown as a dashed line because relatively few measurements were carried out in this region, which is of lesser interest.

with the production of waves in the sea by the wind blowing over the water [3].

The region in which the heat pipe can be used, and in which the steady-state equilibrium of (1) is present, lies between the M_F -axis and the curves 1, 2 and 3 in fig. 4. At a given value of M_F , the values for the instantaneous heat transfer lie on a line — for example the chain-dotted line in the figure — between the M_F -axis and one of the limiting curves. The operating point depends on the temperature difference between the evaporator and the condenser. By using a small variable filling (area under curve 1), it is possible to vary Q_{max} and hence control the temperature of the refrigerator compartment. We shall return to this in more detail towards the end of this article.

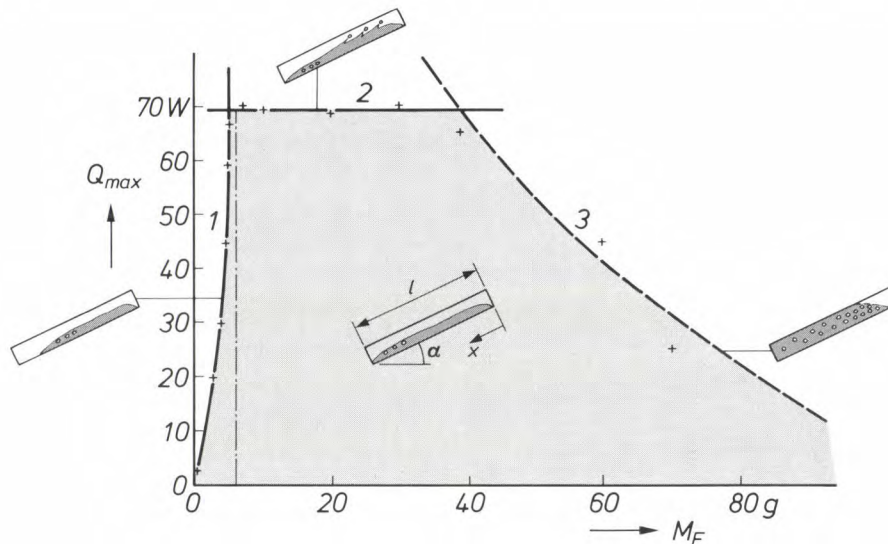


Fig. 4. The region in which the heat pipe can be used. The maximum heat-transport capability Q_{max} is plotted as a function of the filling M_F of the heat pipe in fig. 3, as the result of a number of measurements. Curve 1 is the limit when too little Freon is present: the liquid does not reach the bottom of the pipe. Curve 3 (shown dashed because not many measurements were made in this less-interesting region) is the limit when there is too much Freon in the pipe: the liquid level is raised by the vapour bubbles and the condenser becomes obstructed by the liquid. Line 2 is the limit when the vapour velocity is so high that liquid droplets are carried along with the vapour. This effect occurs if Weber's number exceeds a critical value. The figure in the middle represents a normally operating heat pipe; l length of heat pipe, α angle of inclination of the heat pipe, x co-ordinate of position. The diagrams by curves 1, 2 and 3 illustrate the flow pattern at the appropriate limits. The chain-dotted line indicates one of the possible operating characteristics of the heat pipe.

Line 2, parallel to the M_F -axis, corresponds to a third deviation from the steady-state equilibrium of (1). If a large amount of heat must be transferred, liquid is pulled along by the fast-moving vapour. This 'entrainment' effect, which does not depend on the amount of liquid in the pipe, is characterized by Weber's number:

$$Wb = \frac{\rho_x v^2 l}{\gamma}, \quad (2)$$

where v is the axial velocity of the vapour and γ is the surface tension of the liquid. The effect is comparable

Controlling the temperature in the refrigerator compartment

Control with an inert buffer gas

One of the ways in which the heat flow through the heat pipe can be varied is shown in fig. 5. Above the actual heat pipe there is a space filled with inert buffer gas. The quantity of buffer gas determines the effective area of the condenser C_R . The vapour flowing in the heat pipe sweeps along any buffer gas that may be

[3] See § 268 in H. Lamb, Hydrodynamics, Dover Publications, New York 1945.

present. This gives rise to an interface S , with the buffer gas and superheated Freon vapour above it. The pressure above the interface is equal to the sum of the pressure of the buffer gas p_g and the saturated vapour pressure p_{t_p} of Freon at the temperature t_p of the primary circuit. The pressure below the interface is equal to the saturated vapour pressure p_{t_e} of Freon at the temperature t_e of the evaporator. Since we are considering a steady-state situation, we have equilibrium, so that:

$$p_{t_e} = p_{t_p} + p_g \quad (3)$$

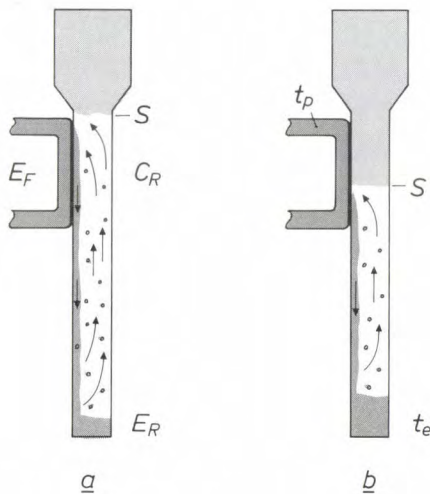


Fig. 5. Principle of heat-transfer control by means of an inert gas. Variation of the effective volume of inert gas changes the effective area of the condenser C_R . S interface between Freon at the saturated vapour pressure at the temperature t_e of the evaporator and a mixture of inert gas with Freon at the saturated vapour pressure at the temperature t_p of the primary circuit. Other symbols are as explained in the caption of fig. 2. *a*) Condenser fully operational. *b*) Condenser only partly operational.

The way in which the quantity of buffer gas can be varied can be seen from *fig. 6*. The temperature t_R in the refrigerator compartment R is measured by means of the temperature sensor T . The current through a heater element in the space with the adsorbent material is switched on and off by the control unit CU , depending on the temperature measured. When the adsorbent is heated, buffer gas is released. The interface S in the condenser (see *fig. 5*) is then displaced downwards and the quantity of heat removed from the refrigerator compartment R decreases. As a result the temperature t_R increases. The current is then switched off by CU and buffer gas is again adsorbed. *Fig. 7* shows what the various parts of the system of *fig. 6* look like in practice.

Suitable buffer gases include nitrogen and carbon dioxide. The adsorbent chosen is a zeolite of the type MS4A ('4Å molecular sieve'). This type of zeolite has pores of 0.40 nm diameter and is thus capable of ad-

sorbing CO_2 molecules (0.330 nm diameter) or N_2 molecules (0.315 nm diameter). The larger CCl_2F_2 molecules cannot be absorbed, however. After trials with both types of buffer gas, we decided to use CO_2 , largely because the same weight of adsorbent can adsorb ten times as much CO_2 as N_2 . Furthermore, the curve of the gas pressure as a function of the temperature of the adsorbent is steeper for CO_2 than for N_2 , which means that the control responds more rapidly. *Fig. 8* gives the adsorption curves for the combination MS4A/ CO_2 , with the temperature t_A of the adsorbent as a parameter. It can be seen that the adsorbent must be heated to between 125 and 150 °C to produce sufficient desorption of buffer gas.

It is very important to fill the buffer space with the optimum quantity of CO_2 . This can be determined as follows. The maximum volume of CO_2 required is equal to that of the buffer space plus that of the condenser, and is reached when the refrigerating action is stopped and no further heat can be transferred by the heat pipe — for example if the refrigerator compartment has to be cleaned. In such a situation the highest

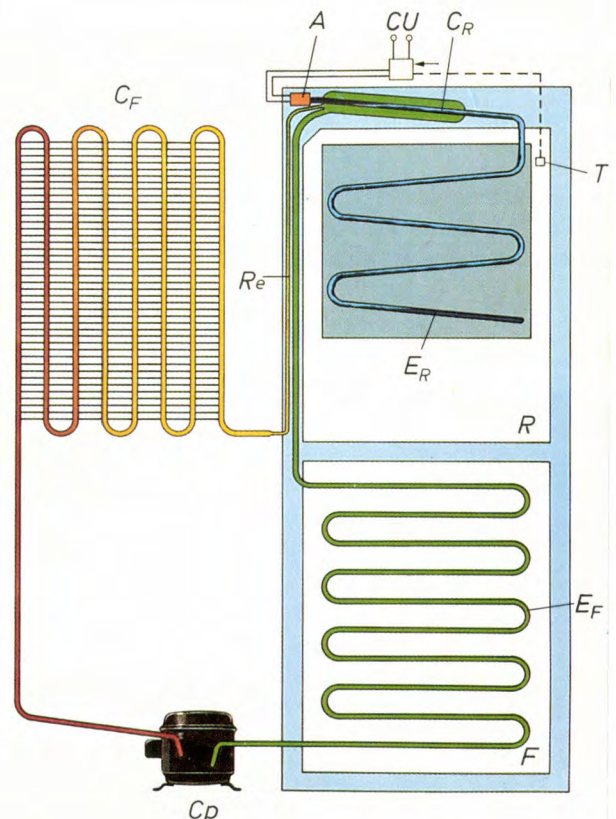


Fig. 6. Design of the temperature-control system with a buffer gas. T temperature sensor. A adsorbent, which can be heated by a heater element. The current through the element is switched on and off by the control unit CU , depending on the temperature measurement by T . Other symbols are as in *fig. 1*.

CO₂ pressure required is reached when the primary refrigeration circuit is permanently switched on ('deep-freeze' setting) at a maximum ambient temperature of 32 °C (as quoted in the ISO standard). In this extreme case the primary evaporator, and therefore the condenser of the heat pipe as well, reach a temperature $t_p = -32$ °C, and the evaporator of the heat pipe is then at a temperature equal to the ambient temperature: $t_e = 32$ °C. By employing equation (3) and tables for the saturation vapour pressure of CCl₂F₂, we find that the pressure of the buffer gas is $p_g = p_{32} - p_{-32} = 7.84 - 0.92 = 6.92$ bars. At this pressure the optimum quantity of CO₂ is 4.2 g for the volume in our system. Curve 1 in fig. 9 gives the same optimum value as the result of a number of measurements. The temperature t_R of the refrigerator compartment was measured under the conditions stated above while the quantity M_G of the buffer gas was varied. At $M_G \geq 4.2$ g the temperature is indeed equal to the ambient temperature: the heat pipe is out of action. If the quantity of buffer gas is smaller, the temperature is lower and obviously heat is still removed from the refrigerator compartment by the heat pipe. Curve 2 in fig. 9 relates to the measurement made on a fully operating heat pipe, i.e. with the adsorbant at a low temperature, and also at an ambient temperature of 32 °C. The measurements show that the heat-transport capability of the heat pipe with a quantity of 4.2 g of buffer gas is indeed sufficient to meet the requirement that the temperature in the refrigerator compartment should be between 0 and 5 °C at a maximum ambient temperature of 32 °C.

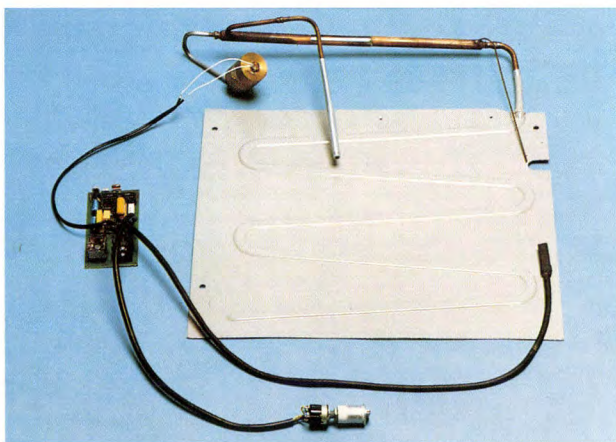


Fig. 7. Photograph of the evaporator plate incorporating the evaporator of the heat pipe. At the upper end the heat exchanger formed by E_F and C_R can be seen (see fig. 3) and the buffer gas with adsorbent connected to it. The electronic temperature-control circuit is on the left. The temperature sensor, attached to the evaporator plate, is on the right. The potentiometer for setting to the required temperature in the refrigerator compartment is at the lower end.

We can also calculate the maximum permissible quantity of Freon. When the refrigerator-freezer cabinet is assembled in the factory and the two refri-

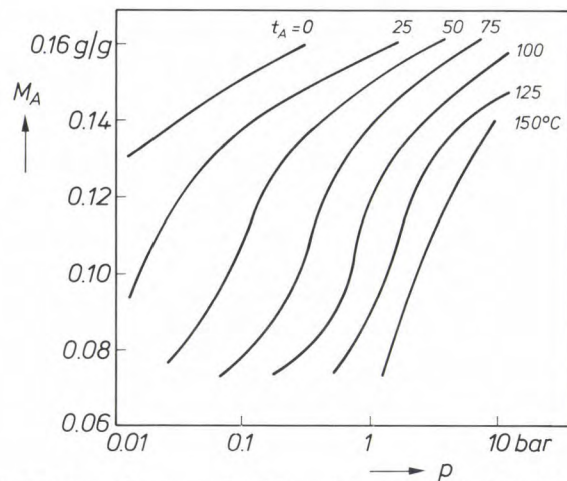


Fig. 8. Adsorption curves for CO₂ and the adsorbant zeolite MS4A. The quantity M_A , in grams of adsorbed CO₂ per gram of zeolite, is plotted as a function of the pressure p of the non-adsorbed CO₂, for different adsorbant temperatures t_A .

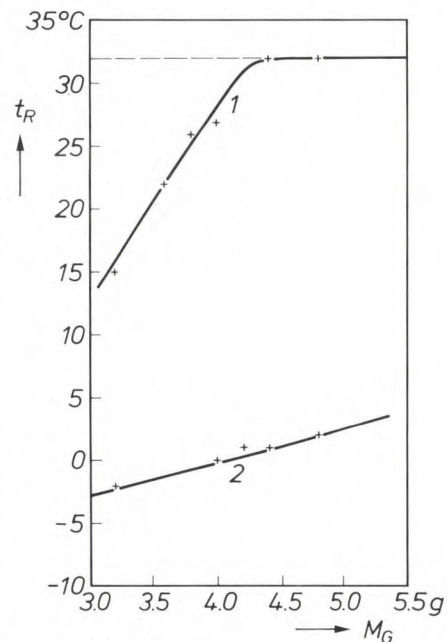


Fig. 9. Determining the optimum quantity of CO₂ buffer gas in the heat pipe. For two situations the temperature t_R in the refrigerator compartment is shown as a function of the quantity M_G of buffer gas. Curve 1 relates to the extreme case where the primary refrigeration circuit is permanently switched on ('deep-freeze' setting), the ambient temperature is 32 °C (the maximum value specified by the ISO standard) and the adsorbant is heated. It can be seen that with a filling of 4.2 g of CO₂ the refrigeration of the refrigerator compartment can just be switched off (by heating the adsorbant), as would be necessary for cleaning. Curve 2 relates to maximum operation of the heat pipe, again at an ambient temperature of 32 °C. The adsorbant is not heated in this case. It can be seen that with a filling of 4.2 g of CO₂ the heat-transport capability of the heat pipe is still high enough to meet the requirements that the temperature in the refrigerator should be between 0 and 5 °C. Curve 2 is not parallel to the M_G -axis because, with a larger CO₂ filling, a rather higher pressure is necessary for the CO₂ to be taken up by the adsorbant. This situation corresponds to a higher evaporator temperature t_e and hence to a higher temperature t_R .

geration systems have been filled, the double walls of the cabinet are filled with polyurethane foam. This material is produced by a chemical reaction *in situ* while the entire cabinet is heated to a temperature of 75 °C in an oven. If liquid Freon was present in the heat pipe at this temperature, the pressure in the heat pipe would rise to the saturated vapour pressure of 20.9 bars. For structural reasons, however, the pressure in the heat pipe should not exceed 9 bars. In *fig. 10* curve 1 is a plot of the saturated vapour pressure of Freon 12 as a function of temperature. The point *P* relates to the safety limit: 9 bars at 75 °C.

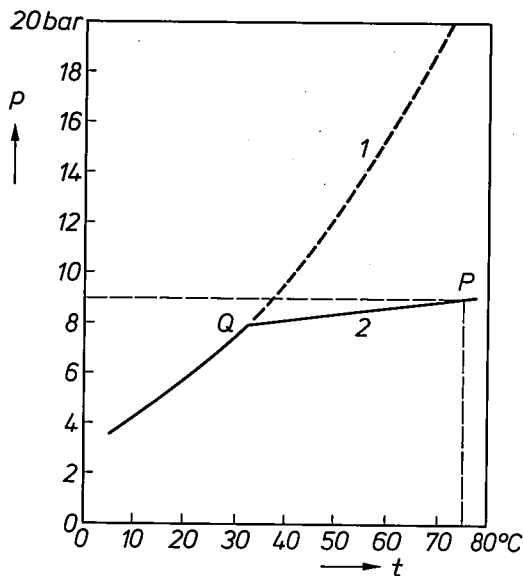


Fig. 10. Determining the optimum quantity of Freon in the heat pipe in the system with buffer gas. Curve 1 gives the saturated vapour pressure p of CCl_2F_2 (Freon 12) as a function of temperature t . At a temperature of 75 °C (the temperature reached when the walls of the cabinet are being filled with polyurethane foam) the pressure would rise to an excessive level if liquid Freon were present. The pressure should then be no higher than 9 bars, however: the point *P*. Curve 2, relating to the function $p/(273 + t) = \text{constant}$, relates to a situation that is only just safe. At point *Q* the liquid has then just vaporized during the heating-up period. The density of the saturated vapour at the associated temperature of 32 °C gives the maximum permissible quantity of Freon (6 grams). The corresponding operating line (the chain-dotted line) in *fig. 4* indicates that with this filling the heat pipe can operate at maximum heat-transport capability.

Curve 2 gives a plot of the function $p/(273 + t) = \text{constant}$. At the point *Q*, where curves 1 and 2 intersect, the liquid must have just disappeared during the heating-up period. The density of the saturated Freon vapour at the corresponding temperature of 32 °C is 0.045 g/cm³. At temperatures up to 75 °C the buffer gas is almost completely taken up by the adsorbant, and therefore in calculating the quantity of Freon vapour we must allow for the entire volume of the heat pipe, including the space for the buffer gas. This

total volume is 133 cm³, and with the density quoted above, it corresponds to a maximum permissible quantity of Freon of 6 g. This is also nearly equal to the limit filling at which the heat pipe can just operate at maximum heat-transport capability, as indicated by the chain-dotted line in the graph of *fig. 4*. The quantities of Freon and CO₂ used for filling the heat pipe must therefore be determined very accurately.

Control without buffer gas

The system shown in *fig. 6* can be greatly simplified by using an adsorbant capable of taking up Freon. The inert buffer gas is then no longer necessary. The heat-transport capability is varied by operating the heat pipe in the region under curve 1 in *fig. 4*. The effective part of the evaporator is then varied, and not the effective part of the condenser, as in the previous system. The control action then has to be reversed. If the temperature of the refrigerator compartment rises, the heating of the adsorbant must be switched on; in the previous system the heating of the adsorbant then had to be switched off.

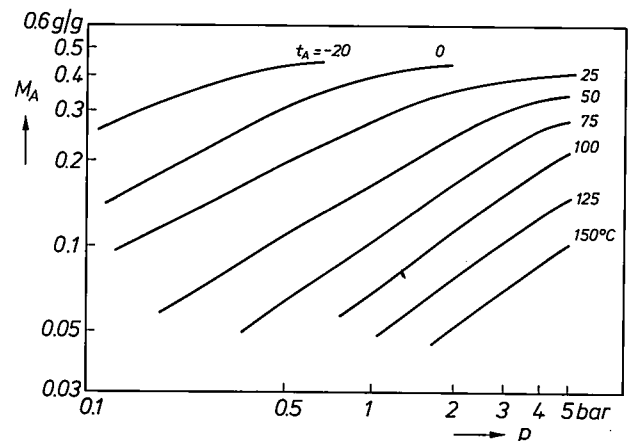


Fig. 11. Adsorption curves for CCl_2F_2 and the adsorbant silica gel. The quantity M_A , in grams of adsorbed CCl_2F_2 per gram of silica gel, is given as a function of the pressure p of the non-adsorbed CCl_2F_2 , for different temperatures t_A of the adsorbant. The scales are logarithmic for both axes, and not as in *fig. 8*.

A good adsorbant for Freon is silica gel. It has the advantage of being much cheaper than zeolite, and it is also much easier to activate (this is done mainly to remove the water). Zeolite has to be activated in a vacuum at a pressure of 10 to 100 Pa for about four hours at a temperature of 450 °C. Silica gel, at the same pressure, requires only six minutes at about 200 °C. Silica gel is also less sensitive to residual traces of water. The adsorption characteristics of the combination CCl_2F_2 / silica gel are given in *fig. 11*.

Since the refrigerant is mostly adsorbed in the silica gel during the filling of the double walls of the refrigerator with polyurethane foam, the pressure in the heat pipe is lower. The quantity of Freon is therefore subject to less critical limits than in the previous system, and this considerably simplifies manufacture.

A disadvantage of the buffer-gas system is that, when a new refrigerator is put into operation, the mixing of CO₂ and Freon in the heat pipe causes some delay. It takes some time before the refrigerant and the buffer gas are separated to give the ideal steady-state situation indicated in fig. 5. This difficulty does not occur in the system without buffer gas, and the full heat-transport capability of the heat pipe is directly available.

The factory at Cassinetta di Biandromo, Italy — part of the Philips Major Domestic Appliances Division — has produced 30 prototype refrigerator-freezers, operating with the heat pipe with CO₂ buffer gas as described here. These prototypes were tested and

operated satisfactorily for a year and met the requirements of the ISO standards comfortably. In the meantime some of these refrigerators have been fitted with a heat pipe without buffer gas. These also operated satisfactorily, with the advantageous features mentioned above.

Summary. Until recently refrigerator-freezers with separate temperature controls required two distinct refrigeration circuits, each with its own compressor, for optimum operation. The use of a heat pipe, which transfers heat from the refrigerator to the refrigeration circuit of the freezer, makes a second compressor unnecessary. A heat pipe can transmit a very high heat flow with a small temperature gradient between condenser and evaporator. The heat-transport capability has an upper limit, however, connected with flow effects in the pipe and with the quantity of refrigerant. To control the temperature in the refrigerator compartment it is necessary to be able to reduce the heat-transport capability of the heat pipe as required. This can be done by limiting the action of the condenser or evaporator. In the first case it is necessary to use an inert buffer gas, which can be taken up in an adsorbant (a zeolite). In the second case the refrigerant itself is adsorbed in a different type of adsorbant (silica gel).

Electromechanical transducers with no hysteresis

A. Schnell

The electrostrictive effect has never been regarded as a serious competitor of the piezoelectric effect as the working principle in electromechanical transducers, since it did not seem to offer sufficient sensitivity. However, the author now shows that electrostrictive transducers with ceramics of the barium-titanate type, with the Curie point just below the temperature of operation, are sufficiently sensitive. Their great advantage is that they are almost completely free of the hysteresis found in piezoelectric transducers, an effect that can be a nuisance.

Introduction

The piezoelectric effect, discovered in 1880 by Pierre and Jacques Curie, is nowadays the subject of renewed interest. The effect can be used for converting a voltage directly into a displacement. Although this can also be done using a coil with an iron core, these are relatively large and heavy, the magnetic field can create interference, and heat can be continuously dissipated in the energized coil without any external work being performed. Piezoelectric materials are used in servomechanisms of various kinds, e.g. for controlling the magnetic heads in video cassette recorders in the Video 2000 system, and also in experimental relays and matrix printers.

When piezoelectric materials are used in control circuits of analog rather than 'on/off' type, hysteresis is undesirable. The materials generally used, however, are ferroelectric ceramics of the BaTiO₃ type, which, as the term 'ferroelectric' indicates, do give hysteresis in their piezoelectric behaviour. While there are piezoelectric materials that are not ferroelectric and therefore give hardly any hysteresis, their sensitivity (the ratio of the relative elongation to the applied field-strength) is relatively low. The sensitivity of quartz, for example, is about 50 times less than that of the ferroelectric ceramics in general use. Electronic circuits can sometimes be used to compensate for the consequences of hysteresis, but they make the servomechanism unnecessary complicated.

One answer to this problem is to use ferroelectric ceramics whose composition is such that the Curie point is just below the temperature of operation

Dr A. Schnell is with Philips GmbH Forschungslaboratorium Aachen, Aachen, West Germany.

(usually room temperature). Above the Curie point, as in ferromagnetic materials, the hysteresis effects have virtually disappeared. (By analogy with the term 'paramagnetic', materials in this temperature range are referred to as 'paraelectric'.) To obtain a displacement dependent on the electric potential it is necessary to make use of the electrostrictive effect, which is more general. Normally this effect is barely perceptible and therefore of little technological significance. However, since ferroelectric materials have a very high permittivity at temperatures close to the Curie point, relatively large electrostrictive displacements can be achieved. A difficulty is that in this temperature range the permittivity — and hence the sensitivity — is highly temperature-dependent. This difficulty can however be overcome reasonably well by slightly adapting the composition of the ceramic, depending on the application.

Another complication encountered with electrostriction is that the effect is proportional to the square of the field-strength, and the displacement is therefore always in the same direction. If a constant electric bias field is applied, however, displacements from the biased position in both positive and negative directions can be obtained. To increase the useful displacements, a ceramic rod composed of two layers — rather like a bimetallic strip — can be used, with opposing bias fields in the two layers. A configuration of this type is also employed in piezoelectric applications, but with opposite polarizations in the two layers.

In the following we shall first look at the general relationship of the electrostrictive and the piezoelectric effect (and other effects) to the crystal structure. We

shall then consider in more detail the electrostrictive effect in ceramics and in barium-titanate type materials above the Curie point. Finally, we shall consider the results of measurements on such materials and show how an electromechanical transducer can be constructed from these materials.

Relationship of the piezoelectric and other effects to the crystal structure

Every insulator, whether crystalline or amorphous, undergoes slight deformation when it is subjected to an external electric field. This change of shape is caused by the forces that the charges induced by the electric field exert upon each other. The effect is known as electrostriction.

The piezoelectric effect is much less common. It occurs in all crystalline insulators that possess a particular asymmetry in their crystal structure. The best known example of these is quartz. Deformation produces an electric polarization in these materials, and conversely the application of an electric potential produces deformation. In the piezoelectric effect the deformation and the applied voltage always have the same sign. In electrostriction the deformation is always in the same direction, as we noted earlier.

There are 32 classes of crystal symmetry, 11 of them with a centre of symmetry. Materials with this centrosymmetric crystalline lattice are not piezoelectric. The materials of the other crystal classes do give the piezoelectric effect (with one exception: the cubic class 432, which has other symmetries that compensate for the absence of a centre of symmetry).

Of the twenty non-centrosymmetric piezoelectric crystal classes there are ten in which the crystals have a spontaneous electric polarization. These are known as polar crystals. In eight of the ten classes the direction of polarization (the polar axis) is also an axis of symmetry of the crystal; one example is tourmaline. In the two other classes (classes l and m) the direction of polarization cannot be predicted from the lattice geometry. The spontaneous polarization is in general not measurable, since most crystals are not perfect insulators. The effect is temperature-dependent and can therefore be observed when the temperature of the crystal is changed. For this reason the spontaneous electric polarization effect is frequently referred to as pyroelectricity [1].

Another effect is found in some of the crystals in the ten classes that give spontaneous polarization: the sign of the spontaneous polarization suddenly changes when an applied external electric field exceeds a critical value. The structure of the material then changes into an associated twinned crystal structure, which is geo-

metrically almost identical with the first. The material then has the opposite electric polarization. Because of this reversal, the curve of polarization as a function of the electric field-strength has the same characteristic shape as the hysteresis curve of ferromagnetic materials. These insulators are therefore said to be 'ferroelectric' [2]. The longest-known ferroelectric material is Rochelle salt. Barium titanate and materials with crystal structures that resemble it (e.g. with Pb substituted for Ba and Zr substituted for Ti) are other ferroelectric materials discovered more recently.

When ground titanates and zirconates of lead and barium and compounds of the form $Ba_{1-x}Pb_xTi_{1-y}Zr_yO_3$ are sintered at high temperature, ceramics are obtained whose properties can be changed by adapting the composition. Since the spontaneous polarizations of the individual crystals are randomly oriented, the resulting polarization is equal to zero. When an external electric field is applied, however, the polarizations can be aligned in the same direction. A material is thus given the desired piezoelectric properties [3] and the associated hysteresis, which is usually undesirable.

A typical feature of all ferroelectric materials, including the ceramics, is that, like ferromagnetics, they lose their ferroelectric properties at temperatures above the Curie point, because they undergo a phase transition. For $BaTiO_3$ this temperature is about $120^\circ C$ and the tetragonal crystal structure changes into a cubic structure. In the region of the Curie point the relative permittivity ϵ_r is very high. The solid curve in *fig. 1* is a plot of ϵ_r as a function of temperature for pure $BaTiO_3$. In practice such a steep curve is

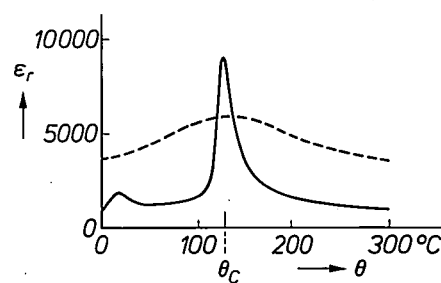


Fig. 1. Relative permittivity ϵ_r as a function of temperature θ ; the solid curve relates to pure single-crystal $BaTiO_3$, the dashed curve to slightly impure ceramic $BaTiO_3$. Above the Curie point θ_C there is a phase transition from tetragonal to cubic crystal structure, the ferroelectric properties of the material disappear and there is a peak in the ϵ_r curve. The peak for the impure $BaTiO_3$ is spread out over a larger temperature range, and the maximum for ϵ_r has a lower value than for pure $BaTiO_3$.

[1] J. F. Nye, *Physical properties of crystals*, Clarendon Press, Oxford 1957.

[2] F. Jona and G. Shirane, *Ferroelectric crystals*, Pergamon Press, Oxford 1962.

[3] J. C. B. Missel, *Piezo-electric materials*, Philips tech. Rev. 11, 145-150, 1949/50.

usually undesirable. The width and height of the peak can be varied by modifying the chemical composition and the physical properties such as the grain size and porosity, see the dashed curve in fig. 1. Materials of this type, with a high relative permittivity (up to 14 000) have been used for some time in capacitors [4]. As we shall see, these materials, with their high permittivity and strong electrostrictive effect, which is related to the piezoelectric properties below the Curie point, can be very useful in electromechanical transducers.

The electrostrictive effect in barium-titanate-type ceramics above the Curie point

The relative deformation caused by electric polarization in an isotropic insulator can in general be described by the equation

$$s_i = Q_{ij}(P_{sp,j} + P_{ind,j})^2, \quad (1)$$

where the subscripts i and j can assume the values 1, 2 and 3, corresponding to the directions x , y and z in an orthogonal coordinate system. In the equation $P_{sp,j}$ represents the spontaneous polarization and $P_{ind,j}$ represents the polarization induced by an external electric field, both in the j -direction. The quantity s_i represents the relative deformation in the i -direction caused by the polarizations. Q_{ij} is a constant of the material. If in (1) we substitute the relation between the induced polarization and the applied field-strength E_j ,

$$P_{ind,j} = \epsilon_0(\epsilon_r - 1)E_j, \quad (2)$$

where ϵ_0 is the permittivity of free space, and if $P_{sp,j}$ is assumed to be zero, we obtain the square-law relation of the electrostrictive effect, between the electric field-strength and the mechanical deformation.

In the piezoelectric materials of the ten crystal classes in which there is spontaneous polarization, the spontaneous polarization is usually much larger than the induced polarization. Equation (1) therefore approximates to

$$s_i = 2Q_{ij}P_{sp,j}P_{ind,j} + s_{sp,i}, \quad (3)$$

where $s_{sp,i} = Q_{ij}P_{sp,j}^2$ represents the constant deformation caused by the spontaneous polarization. In piezoelectric materials in general use the deformation is therefore linearly dependent on the induced polarization and hence on the electric field-strength. The term 'induced piezoelectricity' is therefore used [8].

From equations (1) and (2) it follows that the relative permittivity ϵ_r of the non-spontaneously polarized insulators must have a high value if large deformations are to be produced by means of the electrostrictive

effect. This condition is satisfied by the ferroelectric materials if they are used above the Curie point θ_c in the temperature range

$$\theta_c \leq \theta < \approx 1.3 \theta_c. \quad (4)$$

If the constant Q_{ij} also has a reasonably high value, it is then possible to achieve electromechanical deformation values that are comparable with those of conventional piezoelectric materials.

Measurements on piezoelectric materials and ferroelectric materials above the Curie point

Fig. 2. gives the results of measurements on a sample of polarized piezoelectric ceramic, Philips PXE 52 [6], with dimensions of 5 by 5 mm and a thickness of 0.3 mm, with electrodes on opposite faces. Fig. 2a shows the curves obtained by applying 50 Hz alternating voltages of different amplitude and measuring the relative deformation s_3 perpendicular to the ap-

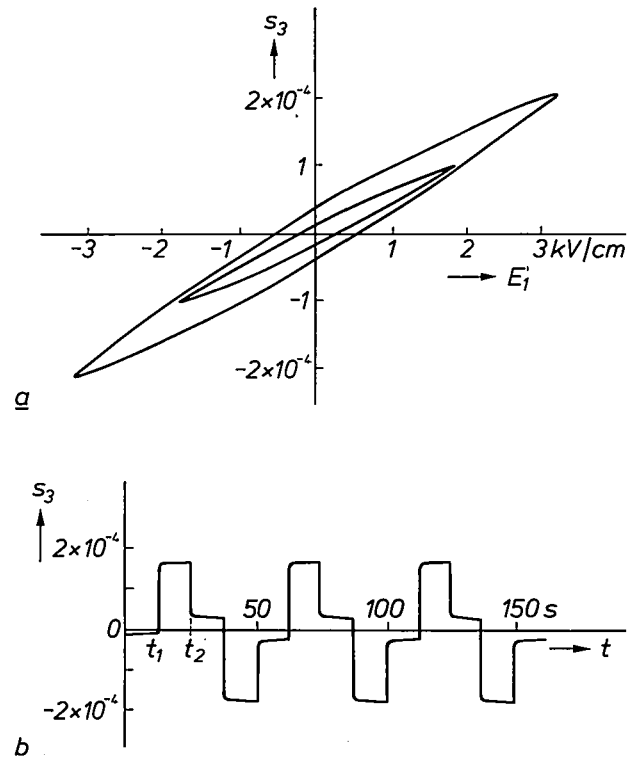


Fig. 2. Results of measurements on a sample of polarized piezoelectric ceramic, Philips PXE 52, 5 mm by 5 mm and 0.3 mm thick, with electrodes on opposite faces. a) The hysteresis loops obtained when alternating voltages at 50 Hz and of different amplitude are applied across the thickness of the sample, and the relative deformation s_3 is measured in a direction perpendicular to the field-strength E_1 . The mean slope of the hysteresis loop depends on the amplitude. The hysteresis effect is due to the ferroelectric properties of the ceramic. b) Plot of the relative deformation s_3 when a direct voltage producing a field-strength E_1 is applied and kept positive for 12.5 s (from t_1 to t_2), kept at zero for 12.5 s, made negative for 12.5 s, and so on. During the periods when the voltage is zero there is still considerable deformation. This means that there is not only dynamic hysteresis (fig. 2a) but also appreciable static hysteresis.

plied fields. It can be seen that the average slope of the curves depends on the amplitude, and that a large hysteresis effect is present. Fig. 2b shows that the hysteresis is not only a dynamic effect. If a direct voltage is applied at the instant t_1 , a relative extension

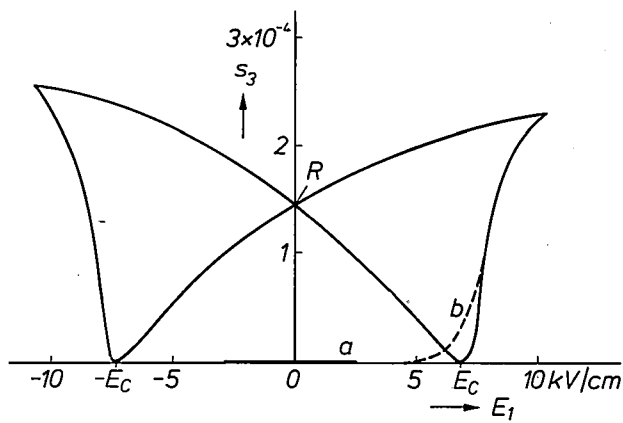


Fig. 3. The 'butterfly curve' with the relative deformation s_3 plotted as a function of the field-strength E_1 in a direction perpendicular to it. The figure shows the results of measurements on a piezoelectric sample similar to the one in fig. 2, but with a thickness of 0.5 mm. The zero of the s_3 -axis is not the same as in fig. 2. If the material is still 'virgin', i.e. not yet polarized, and the coercive field-strength E_C is not exceeded, the straight line a is described when an alternating voltage is applied. The deformations of the individual crystals are arbitrarily oriented and give a resultant deformation of zero. If the voltage is increased and the field-strength exceeds the value E_C , the dashed curve b is described first: the material now becomes polarized and the polarizations of the individual crystals become aligned. When the field-strength reaches a value of about $2E_C$, the voltage is reduced again and another branch of the butterfly curve is described, leaving the remanent deformation R when the field-strength has fallen to zero. When the field-strength reaches the value $-E_C$ in the other direction, the polarization of the material reverses and causes the direction of the deformation to change. In practical applications of the material, therefore, the field-strength must remain between the values E_C and $-E_C$.

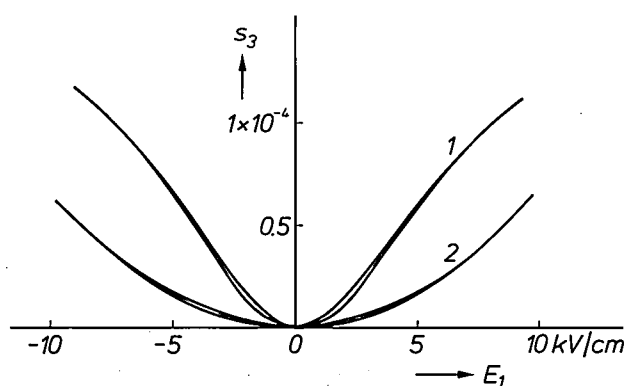


Fig. 4. Results of similar measurements on a sample with the same dimensions as in fig. 3, but made of ceramic capacitor material. Measurements were made at temperatures above the Curie point of the material; the Curie point is just below room temperature. Since the ferroelectric effect has almost disappeared because of the transition to a different crystal structure, hardly any hysteresis is perceptible and the remanent deformation has fallen to virtually zero. The electrostrictive effect now makes the deformation consistently positive, irrespective of the sign of the field-strength. Curve 2 relates to a higher temperature than curve 1. At the higher temperature the hysteresis effects are even smaller, but the deformations are also smaller, since the relative permittivity ϵ_r is lower.

is produced that does not entirely disappear when the direct voltage is switched off at the instant t_2 . The same effect occurs in the opposite direction. This means that the length of the sample at zero applied voltage is not accurately known, and this is highly undesirable in applications such as an electromechanical transducer in control circuits. Both the dynamic and the static hysteresis are due to the ferroelectric properties of the ceramic material. The effects can be reduced slightly by optimizing the composition of the material, but cannot be completely eliminated.

It follows from fig. 3 that the extension cannot be increased merely by increasing the amplitude of the alternating voltage applied to a polarized piezoelectric. The figure shows the results of measurements on the same block of material as in fig. 2, but this time 0.5 mm thick. The line a is described in the 'virgin material', before it is polarized. The polarizations of the different crystals are still differently oriented and give zero resultant displacement. When the field-strength exceeds the coercive field-strength E_C , the material becomes polarized — via the dashed curve b — and the individual polarizations of the crystals become identically aligned. When the field-strength is reduced again, after reaching a value of about $2E_C$, another branch of the familiar 'butterfly curve' [6] is described, leaving a remanent deformation R when the voltage has fallen to zero. The deformation R corresponds to the term $s_{sp,i}$ in equation (3). Whenever the field-strength goes beyond E_C or $-E_C$, the direction of the polarization reverses and the direction of the deformation changes. These situations have to be avoided in practice and the field-strength must not be allowed to reach the values E_C or $-E_C$.

Fig. 4 shows the results of measurements on ferroelectric material in the electrostrictive temperature range. The material is a sample of a standard Philips barium-titanate-based ceramic, used as a dielectric in capacitors. The Curie point of this material is below the temperature at which the material is used — room temperature or a little above. It has a very high relative permittivity, about 14 000, so that capacitors of small dimensions can be made. In our application, an electromechanical transducer, the relatively large electrostrictive deformations possible with this high value of ϵ_r are of particular interest; see equations (1) and (2). The curves in fig. 4 show that the hysteresis is very small and that the extensions resulting from the electrostrictive effect are always in the same direction, whatever the polarity of the field. Curve 2 relates to a

[4] G. H. Jonker and J. H. van Santen, The ferro-electricity of titanates, Philips tech. Rev. 11, 183-192, 1949/50.

[5] Philips Data handbook, Components and materials, Part 16, 1982.

[6] See for example page 148 in [2].

higher temperature than curve 1 and indicates a smaller sensitivity, because of the decrease in the permittivity. Curve 2 shows even less hysteresis. The square-law relation between E_1 and s_3 clearly applies at small field-strengths, in this case up to about 4 kV/cm. When the field-strength increases, the curves start to deviate slightly from the parabolic shape, since ϵ_r varies slightly with the field-strength at higher values.

The square-law relation between the relative deformation and the electric polarization is also seen very clearly in fig. 5, which gives the results of measurements on the same sample as in fig. 4 and for the temperature of curve 1 in that figure. With the quantity P_1^2 plotted along the horizontal axis the resultant curve is very nearly linear and indicates that there is virtually no hysteresis.

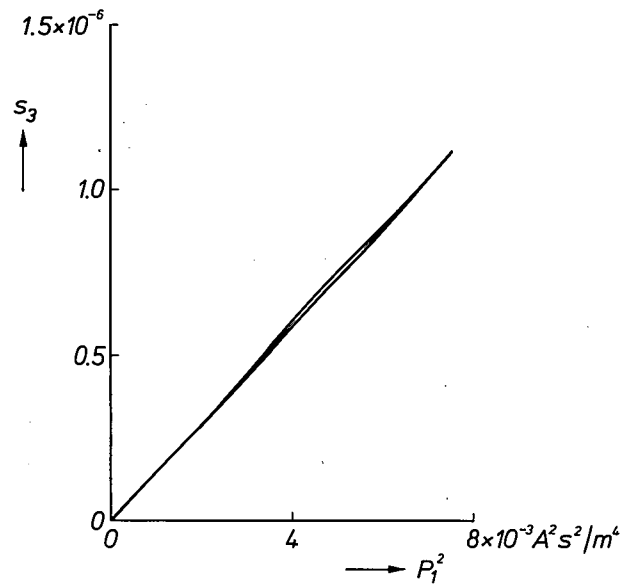


Fig. 5. The relative deformation s_3 as a function of the square of the polarization P_1 for the same sample and the same conditions as for curve 1 in fig. 4. The linear curve confirms the square-law relation between deformation and polarization that is inherent to the electrostrictive effect. Both branches of curve 1 in fig. 4 coincide in this figure.

The use of ferroelectric materials above the Curie point in electromechanical transducers

The main properties of a ferroelectric used in the electrostrictive range follow from fig. 4. The hysteresis is small and the extension relatively large — at least equal to that of piezoelectric materials — but the relation between extension and field-strength is not linear, and moreover the extension is always in the same direction. An additional difficulty is the temperature dependence of the permittivity and hence of the sensitivity, as demonstrated by the difference between curves 1 and 2 in fig. 4. Investigations are still in progress to see how the temperature dependence might be reduced by changing the composition of the commonly used ceramics. For many applications, however, a broad temperature range is not necessary, and the standard ceramics for capacitors as mentioned earlier are perfectly satisfactory.

If ceramic materials are used as either electrostrictive or piezoelectric transducers in the form of small blocks, the deformations that can be obtained with the usual dimensions are no more than about 1 μm . With flexure elements, consisting of two strips, with electrodes on opposite faces, and bonded together, deformations ranging from 10 to 100 μm are obtainable. The action of such an element is comparable with that of a thermal bimetallic strip. When piezoelectric ceramics are used, the two strips have to be polarized in opposite directions, resulting in a 'piezoelectric bimorph' [7]. When a voltage is applied to the common electrode, one strip becomes longer and the other shorter, in a direction perpendicular to the applied fields. The amplitude of the deflection f of one end of a bending element clamped at the other end is given by the equation [8]

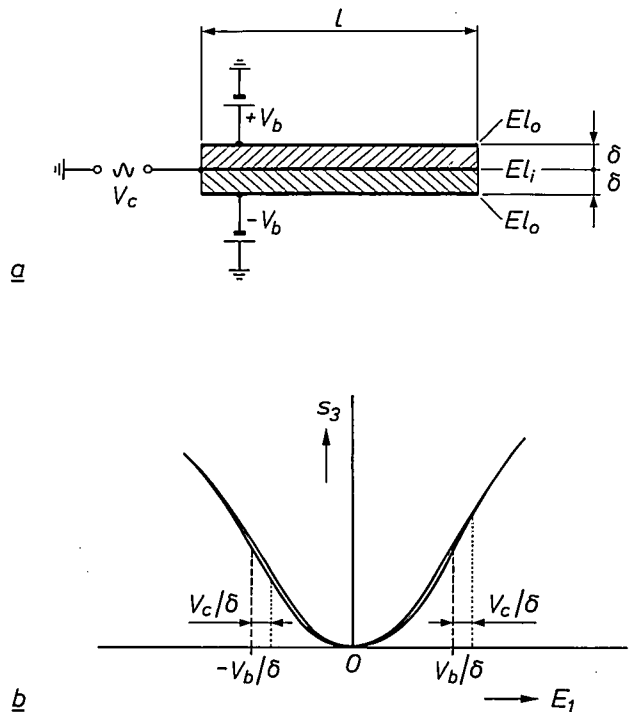


Fig. 6. a) An electromechanical transducer, designed as an electrostrictive bending element. Bias voltages of opposite polarity V_b and $-V_b$ are applied between the common electrode El_i and the two outer electrodes El_o . Since the transducer is constructed as a bimetallic element, the deformation is increased by a factor $3/4\delta$; l is the length and δ the thickness of the strips. The control voltage V_c is applied to the electrode El_i . b) The magnitude of the bias voltages V_b and $-V_b$ must be such as to bring the material into the linear part of the electrostrictive characteristic. The bias voltages V_b and $-V_b$ produce equal extensions in the two strips and cause no deflection at the end. The control voltage V_c , however, makes one strip extend and the other contract, and does produce a deflection of the element.

$$f \approx \frac{3}{4} \frac{l^2}{\delta} d_{31} E_1, \quad (5)$$

where l is the length of the element, δ the thickness of each of the two strips and d_{31} is the piezoelectric constant of the material. The amplitude obtainable is thus increased by a factor $3l/4\delta$.

It is also possible to make a flexure element of this type with a ceramic that deforms electrostrictively. A constant direct-voltage bias must then be applied to the electrodes of the two strips; see fig. 6a. This bias gives the two strips a constant elongation, but no bending takes place. The bias V_b must be large enough to ensure that the operating point lies approximately

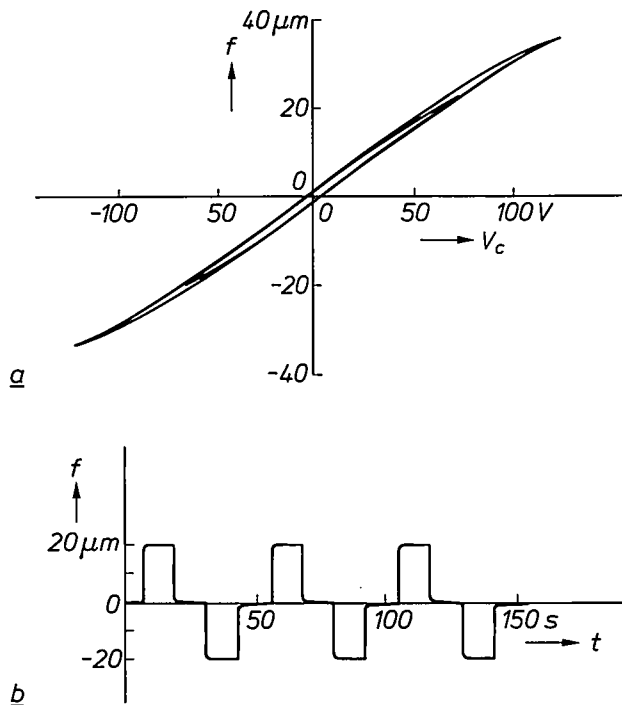


Fig. 7. Measurements on an electrostrictive bending element as in fig. 6, with a length of 12 mm and a thickness of 0.3 mm per strip, giving a total thickness of 0.6 mm. Ceramic capacitor material was used with a Curie point a little below room temperature and a relative permittivity of about 14000. *a*) The deflection f at the end of the clamped element as a function of an applied alternating voltage V_c at 50 Hz, for two different amplitudes. It can be seen that the dynamic hysteresis effects, like those of the piezoelectric material in fig. 2, have virtually disappeared. The magnitude of the deflection is comparable with that of a piezoelectric bending element of the same dimensions. An incidental advantage is that depolarization effects cause no long-term deterioration in performance of the electrostrictive transducer. *b*) Measurement of the static deflection f as a function of time t by applying a direct voltage of alternating polarity with zero-voltage periods in between. There is hardly any remanent deformation during the zero-voltage periods. The static hysteresis has disappeared; this is a great advantage in the application of such elements in servomechanisms.

on the linear part of curve 1 in fig. 4; see fig. 6b. A control voltage V_c applied to the common electrode E_1 increases the field-strength in one strip and decreases it in the other. The electrostrictive element is then deflected at one end; the deflection is analogous to the piezoelectric deflection given by (5).

Fig. 7a gives two characteristic curves for a flexure element of this type, made from the ceramic capacitor material mentioned earlier, with a relative permittivity of 14000. The strips are each 0.3 mm thick and 12 mm long; the frequency of the applied alternating voltage is 50 Hz. An electromechanical transducer has thus been obtained that has a greatly reduced hysteresis, has about the same sensitivity as a piezoelectric material and has a deflection nearly linear with the applied voltage. Fig. 7b shows a waveform quite different from the one in fig. 2b: the static hysteresis has now been reduced to practically zero. Nor can the characteristics of the transducer change as a result of a slow progressive depolarization, as sometimes happens in piezoelectrics. In spite of the essentially square-law relation between voltage and deformation with the electrostrictive effect, a linear electromechanical transducer is obtained. These hysteresis-free elements are suitable for applications where only a very small force is required, with a displacement between 10 and 100 μm . Possible applications include micromanipulators, the deflection of laser beams, optical 'switches' in glass fibre optics and the control of magnetic heads in audio and video recorders.

[7] J. van Randerat and R. E. Setterington (eds), Piezoelectric ceramics, Philips Application book, 1974.

[8] M. R. Steel, F. Harrison and P. G. Harper, The piezoelectric bimorph: An experimental and theoretical study of its quasi-static response, J. Physics D 11, 979-989, 1978.

Summary. A disadvantage of piezoelectric transducers, as used in servomechanisms, is that they give hysteresis effects. This disadvantage can be overcome by making the transducers of bariumtitanate-type ceramics, operated above the Curie point, and using the electrostrictive effect. This effect is generally very small, but can nevertheless be turned to practical advantage because of the high permittivity of ferroelectric materials just above the Curie point. Measurements of the characteristic curves of ferroelectric ceramic materials show that when the temperature is increased the familiar 'butterfly curve' changes to the small loop typical of a material without hysteresis. A flexure element made from two strips of electrostrictive ceramic capacitor material, each strip with an opposite voltage bias, forms an electromechanical transducer that is a good match for a piezoelectric 'bimorph'. A great advantage of this electrostrictive transducer is that it is virtually free of hysteresis.

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 following research laboratories:

Philips Research Laboratories, Eindhoven, The Netherlands	<i>E</i>
Philips Research Laboratories, Redhill, Surrey RH1 5HA, England	<i>R</i>
Laboratoires d'Electronique et de Physique Appliquée, 3 avenue Descartes, 94450 Limeil-Brévannes, France	<i>L</i>
Philips GmbH Forschungslaboratorium Aachen, Weißhausstraße, 51 Aachen, Germany	<i>A</i>
Philips GmbH Forschungslaboratorium Hamburg, Vogt-Kölln-Straße 30, 2000 Hamburg 54, Germany	<i>H</i>
Philips Research Laboratory Brussels, 2 avenue Van Becelaere, 1170 Brussels (Boitsfort), Belgium	<i>B</i>
Philips Laboratories, N.A.P.C., 345 Scarborough Road, Briarcliff Manor, N.Y. 10510, U.S.A.	<i>N</i>

D. J. H. Admiraal (Institute for Perception Research, Eindhoven): A dropout annoyance measuring apparatus 'DAMA' to check magnetic tapes. *J. Audio Engng Soc.* **27**, 788-792, 1979 (No. 10).

D. E. Aspnes*, **J. B. Theeten** & **R. P. H. Chang*** (*Bell Laboratories, Murray Hill, N.J.): Nondestructive characterization of interface layers between Si or GaAs and their oxides by spectroscopic ellipsometry. *J. Vac. Sci. Technol.* **16**, 1374-1378, 1979 (No. 5). *L*

D. E. Aspnes (Bell Laboratories, Murray Hill, N.J.), **J. B. Theeten** & **F. Hottier**: Investigation of effective-medium models of microscopic surface roughness by spectroscopic ellipsometry. *Phys. Rev. B* **20**, 3292-3302, 1979 (No. 8). *L*

S. E. Bradshaw (Philips Semiconductor Development Lab., Nijmegen) & **J. Goorissen**: Silicon for electronic devices. *J. Crystal Growth* **48**, 514-529, 1980 (No. 4). *E*

P. Branquart, **G. Louis** & **P. Wodon**: Aspects du langage CHILL. *Bull. Gr. Progr. et Lang. AFCET* **8**, 80-92, 1979. *B*

P. Branquart, **G. Louis** & **P. Wodon**: Une méthode de description systématique pour langages algorithmiques. *Bull. Gr. Prog. et Lang. AFCET* **9**, 170-180, 1979. *B*

H. J. L. Bressers & **J. G. Kloosterboer**: Thermally and light-induced polymerization of ethyl acrylate and methyl methacrylate, studied by DSC. *Polymer Bull.* **2**, 201-204, 1980 (No. 3). *E*

R. J. Brewer: The low light level potential of a CCD imaging array. *IEEE Trans.* **ED-27**, 401-405, 1980 (No. 2). *R*

K. H. J. Buschow & **N. M. Beekmans**: On the thermal stability and the electrical resistivity of amorphous Th-Fe alloys. *Phys. Stat. sol. (a)* **56**, 505-511, 1979 (No. 2). *E*

R. P. H. Chang, **C. C. Chang**, **W. T. Chiang**, **S. Darack**, **T. T. Sheng** (all with Bell Laboratories, Murray Hill, N.J.) & **J. B. Theeten**: Some physical properties of plasma grown SiO₂. *Extended Abstracts Electrochem. Soc.* **79-2**, 886-887, 1979. *L*

R. Clasen: Low-dose electronradiographic multilayer system with PbO binder layers. *J. photogr. Sci.* **28**, 226-230, 1980 (No. 6). *A*

S. Colak & **E. H. Stupp**: Reverse avalanche breakdown in gated diodes. *Solid-State Electronics* **23**, 467-472, 1980 (No. 5). *N*

W. J. Dallas: Tomosynthesis and computer tomography: a continuous description with examples. *Appl. Optics* **19**, 2472-2476, 1980 (No. 14). *H*

C. B. Dekker: The application of tamed frequency modulation to digital transmission via radio. *Conf. Rec. 1979 National Telecomm. Conf.*, Washington, DC, 7 pp. *E*

M. Delfino & **P. S. Gentile** (Fordham Univ., New York): Approximate nonlinear optical susceptibility of cubic boracites. *J. appl. Phys.* **51**, 2264-2266, 1980 (No. 4). *N*

M. Delfino, **G. M. Loiacono**, **W. A. Smith** & **P. S. Gentile** (Fordham Univ., New York): Calorimetric investigation of the ferroelectric 43m-mm2 phase transition in boracite crystals. *J. solid State Chem.* **33**, 107-114, 1980 (No. 1). *N*

P. A. Devijver: New error bounds with the nearest neighbor rule. *IEEE Trans.* **IT-25**, 749-753, 1979 (No. 6). *B*

R. J. Dewey: High-gain flat-plate phase-monopulse antenna. *Electronics Letters* **16**, 30-32, 1980 (No. 1). *R*

W. F. Druyvesteyn, **L. Postma** & **G. Somers**: Wafer testing of thin film record and reproduce heads. *IEEE Trans.* **MAG-15**, 1613-1615, 1979 (No. 6). *E*

- C. Ducot & G. J. Lubben:** A typology for scenarios. *Futures* **12**, 51-57, 1980 (No. 1). *L, E*
- L. D. J. Eggermont:** Oversampling in waveform coding. *T. Ned. Electronica- en Radiogen.* **44**, 243-256, 1979 (No. 5/6). *E*
- K. Enke & D. Mateika:** Refractive index and optical absorption of bariumhexa-aluminate $\text{BaAl}_{12}\text{O}_{19}$. *J. Mat. Sci.* **15**, 1066-1067, 1980 (No. 4). *H*
- W. van Erk, H. J. G. J. van Hoek-Martens & G. Bartels:** The effect of substrate orientation on the growth kinetics of garnet liquid phase epitaxy. *J. Crystal Growth* **48**, 621-634, 1980 (No. 4). *E, H*
- J.-M. Goethals:** Association schemes. Algebraic coding theory and applications, ed. G. Longo (CISM Courses and Lectures No. 258), pp. 243-283; Springer, Wien 1979. *B*
- R. G. Gossink, H. van Doveren & J. A. T. Verhoeven:** Decrease of the alkali signal during Auger analysis of glasses. *J. non-cryst. Sol.* **37**, 111-124, 1980 (No. 1). *E*
- S. Gourrier, A. Mircea & M. Bacal** (Ecole Polytechnique, Palaiseau Cédex, France): Oxidation of GaAs in an oxygen multipole plasma. *Thin Solid Films* **65**, 315-330, 1980 (No. 3). *L*
- H. C. de Graaff & J. G. de Groot:** The SIS tunnel emitter: a theory for emitters with thin interface layers. *IEEE Trans. ED-26*, 1771-1776, 1979 (No. 11). *E*
- H. B. Haanstra & H. Ihrig:** Transmission electron microscopy at grain boundaries of PTC-type BaTiO_3 ceramics. *J. Amer. cer. Soc.* **63**, 288-291, 1980 (No. 5-6). *E, A*
- H.-J. Hagemann, A. Hero* & U. Gonser*** (* Univ. Saarbrücken): The valence change of Fe in BaTiO_3 studied by Mössbauer effect and gravimetry. *Phys. Stat. sol. (a)* **61**, 63-72, 1980 (No. 1). *A*
- J. Hallais, F. Hottier, J. B. Theeten & G. Laurence:** Vapour epitaxial growth of (Ga,Al)As - GaAs monitored by in-situ ellipsometry. *Extended Abstracts Electrochem. Soc.* **79-2**, 1446-1447, 1979. *L*
- H. Heitmann, P. Hansen, R. Spohr** (Ges. f. Schwerionenforschung, Darmstadt) & **K. Witter:** Properties of magneto-optic $(\text{Gd,Bi})_3(\text{Fe,Ga})_5\text{O}_{12}$ films irradiated with high-energy heavy ions. *J. Magn. magn. Mat.* **15-18**, 1543-1544, 1980 (Part III). *H*
- K. Holford:** Microwave intruder detector, 1 & 2. *Wireless World* **86**, Feb. 1980, 34-38, & March 1980, 79-84 (Nos 1530 & 1531). *R*
- F. Hottier & J. B. Theeten:** Surface analysis during vapour phase growth. *J. Crystal Growth* **48**, 644-654, 1980 (No. 4). *L*
- H. Ihrig, J. H. T. Hengst & M. Klerk:** Conductivity-dependent cathodoluminescence in BaTiO_3 , SrTiO_3 and TiO_2 . *Z. Physik B* **40**, 301-306, 1981 (No. 4). *A, E*
- U. Killat, C. Clausen & G. Rabe:** Binary phase gratings for couplers used in fiber-optic communications. *Fiber and integr. Opt.* **3**, 221-235, 1980 (No. 2/3). *H*
- J. G. Kloosterboer & H. J. L. Bressers:** Evidence for two gel effects during the bulk polymerization of ethyl acrylate from DSC, Rayleigh and Brillouin Scattering. *Polymer Bull.* **2**, 205-210, 1980 (No. 3). *E*
- J. E. Knowles, R. F. Pearson & A. D. Annis:** The angular distribution of the magnetization in magnetic recording tapes. *IEEE Trans. MAG-16*, 42-44, 1980 (No. 1). *R*
- J. E. Knowles:** Magnetic measurements on single acicular particles of $\gamma\text{Fe}_2\text{O}_3$. *IEEE Trans. MAG-16*, 62-67, 1980 (No. 1). *R*
- C. J. Koomen:** Information processing in system design. *Proc. Int. Conf. on Cybernetics and society, Denver 1979*, pp. 8-12. *E*
- C. J. Koomen:** Reducing model complexity in system design. *Proc. Int. Conf. on Cybernetics and society, Denver 1979*, pp. 830-833. *E*
- M. H. Kuhn & R. Geppert:** A low cost speaker verification device. *Proc. 1980 Carnahan Conf. on Crime countermeasures, Lexington, Kentucky*, pp. 57-61. *H*
- P. I. Kuindersma & R. M. van der Heij:** Dynamic methods for determination of the charge density on unipolar electrets. *1979 Annual Rep. Conf. on Electrical insulation and dielectric phenomena, Whitehaven, PA, 1979*, pp. 325-333. *E*
- D. Meignant, D. Boccon-Gibod & J. M. Bourgeois:** Trap localisation in the active layer of GaAs microwave f.e.t.s. *Electronics Letters* **15**, 779-780, 1979 (No. 24). *L*
- A. Milch & P. Tasaico:** The stability of tellurium films in moist air; a model for atmospheric corrosion. *J. Electrochem. Soc.* **127**, 884-891, 1980 (No. 4). *N*
- L. Minnema, H. A. Barneveld & P. D. Rinkel:** An investigation into the mechanism of watertreeing in polyethylene high-voltage cables. *1979 Annual Rep. conf. on Electrical insulation and dielectric phenomena, Whitehaven, PA, 1979*, pp. 480-489. *E*
- A. Mitonneau, J. P. Chané & J. P. André:** Defect characterization at the growth interface in GaAs epitaxy by metallorganic and chloride depositions. *J. electronic Mat.* **9**, 213-229, 1980 (No. 2). *L*
- J. H. Neave, P. Blood & B. A. Joyce:** A correlation between electron traps and growth processes in *n*-GaAs prepared by molecular beam epitaxy. *Appl. Phys. Letters* **36**, 311-312, 1980 (No. 4). *R*

- G. F. Neumark, B. J. Fitzpatrick, P. M. Harnack, S. P. Herko, K. Kosai & R. N. Bhargava:** Potential profiling as a means to determine conductivity type: application to ZnSe.
J. electrochem. Soc. **127**, 983-985, 1980 (No. 4). *N*
- G. F. Neumark:** Are impurities the cause of "self"-compensation in large-band-gap semiconductors?
J. appl. Phys. **51**, 3383-3387, 1980 (No. 6). *N*
- E. de Niet & R. Vreeken:** A magnetoresistive head with magnetic feedback.
IEEE Trans. **MAG-15**, 1625-1627, 1979 (No. 6). *E*
- S. G. Nootboom & G. J. N. Doodeman** (both with Institute for Perception Research, Eindhoven): Production and perception of vowel length in spoken sentences.
J. acoust. Soc. Amer. **67**, 276-287, 1980 (No. 1).
- D. H. Paxman & K. R. Whight:** Observation of lifetime controlling recombination centres in silicon power devices.
Solid-State Electronics **23**, 129-132, 1980 (No. 2). *R*
- P. Piret:** Sliding block implementation of block codes.
IEEE Trans. **IT-25**, 725-728, 1979 (No. 6). *B*
- P. Piret:** Good linear codes of length 27 and 28.
IEEE Trans. **IT-26**, 227, 1980 (No. 2). *B*
- R. D. Plättner, W. W. Krühler** (both with Siemens AG, München), **W. K. Zwicker, T. Kovats & S. R. Chinn** (M.I.T., Lexington, Mass.): The growth of large, laser quality $\text{Nd}_x\text{RE}_{1-x}\text{P}_5\text{O}_{14}$ crystals.
J. Crystal Growth **49**, 274-290, 1980 (No. 2). *N*
- D. Pons & S. Makram-Ebeid:** Phonon assisted tunnel emission of electrons from deep levels in GaAs.
J. Physique **40**, 1161-1172, 1979 (No. 12). *L*
- M. J. Powell:** Site percolation in randomly packed spheres.
Phys. Rev. B **20**, 4194-4198, 1979 (No. 10). *R*
- M. J. Powell:** Computer-simulated random packing of spheres.
Powder Technol. **25**, 45-52, 1980 (No. 1). *R*
- J. M. Robertson, J. P. M. Dameñ & H. A. Algra:** Liquid phase epitaxial growth and properties of spinel thin films.
IEEE Trans. **MAG-15**, 1870-1872, 1979 (No. 6). *E*
- P. Röschmann:** Redistribution kinetics of Ga and Al substitutions in yttrium iron garnet.
J. Magn. magn. Mat. **15-18**, 1305-1306, 1980 (Part III). *H*
- J. C. Rosier, R. Polaert, T. N'Guyen-Trong* & B. Sidoruk*** (* Laboratoire d'Astronomie Spatiale, Marseille): An image tube with a curved microchannel plate and its use in a photon counting imaging system.
Adv. in Electronics and Electron Phys. **52**, 369-378, 1979. *L*
- S. K. Salmon:** Practical aspects of surface-acoustic-wave oscillators.
IEEE Trans. **MTT-27**, 1012-1018, 1979 (No. 12). *R*
- P. J. Severin & H. van Esveld:** The measurement of the total loss and scattering loss coefficients in optical fibres, particularly fibres of compound glasses.
Phys. Chem. Glasses **21**, 58-66, 1980 (No. 1). *E*
- G. A. C. M. Spierings, C. M. G. Jochem, T. P. M. Meeuwssen, F. Meyer & P. J. W. Severin:** Some aspects of the preparation of alkali lime germanosilicate optical fibres.
Phys. Chem. Glasses **21**, 30-33, 1980 (No. 1). *E*
- B. Steinmüller:** The two-solarimeter method for insolation on inclined surfaces.
Solar Energy **25**, 449-460, 1980 (No. 5). *A*
- B. Strocka, G. Bartels & R. Spohr** (Ges. f. Schwerionenforschung, Darmstadt): Lattice strain in garnet single crystals caused by high-energy heavy ion irradiation.
Appl. Phys. **21**, 141-149, 1980 (No. 2). *H*
- T. J. B. Swanenburg:** Machine acquisition of language from examples.
Proc. Int. Conf. on Cybernetics and society, Denver 1979, pp. 17-20. *E*
- T. J. B. Swanenburg:** Self-organization in artificial systems.
Proc. Int. Conf. on Cybernetics and society, Denver 1979, pp. 826-829. *E*
- J. B. Theeten, R. P. H. Chang*, D. E. Aspnes* & T. E. Adams*** (* Bell Laboratories, Murray Hill, N.J.): *In situ* measurement and analysis of plasma-grown GaAs oxides with spectroscopic ellipsometry.
J. electrochem. Soc. **127**, 378-385, 1980 (No. 2). *L*
- M. Urner-Wille:** Magneto-optical properties of amorphous ferrimagnetic GdFeBi thin films.
J. Magn. magn. Mat. **15-18**, 1339-1340, 1980 (Part III). *H*
- N. A. M. Verhoeckx, H. C. van den Elzen, F. A. M. Snijders & P. J. van Gerwen:** Digital echo cancellation for baseband data transmission.
IEEE Trans. **ASSP-27**, 768-781, 1979 (No. 6, Part II). *E*
- H. Verweij:** Raman study of the reactions in a glass-forming mixture with molar composition: $30\text{K}_2\text{CO}_3\text{-}70\text{SiO}_2\text{-}1\text{As}_2\text{O}_3$.
J. Amer. cer. Soc. **62**, 450-455, 1979 (No. 9-10). *E*
- A. B. Voermans, D. J. Breed, W. van Erk & F. M. A. Carpay:** Bubble device materials with orthorhombic anisotropy.
J. appl. Phys. **50**, 7827-7829, 1979 (No. 11, Part II). *E*
- R. Ward:** Developments in electron image projection.
J. Vac. Sci. Technol. **16**, 1830-1833, 1979 (No. 6). *R*
- R. de Werdt & W. A. M. Meeuwissen:** Structuring of bubble overlays by reactive sputter etching in an Ar-H₂O atmosphere.
J. Vac. Sci. Technol. **16**, 2093-2095, 1979 (No. 6). *E*
- H. W. Werner:** Modern methods for thin film and surface analyses.
Mat. Sci. Engng **42**, 1-12, 1980 (No. 1). *E*

Volume 40, 1982, Supplement

PHILIPS TECHNICAL REVIEW

INDEX
VOLUMES 31-40



PHILIPS

Subject index, Volumes 31-40

Figures in bold type indicate the volume number, figures in ordinary type indicate the page number. Subjects dealt with in volumes 1-30 are listed in the indexes included with volumes 10, 15, 20, 25, 30 and 35. The asterisk * indicates that the entry refers to a photograph and caption.

- Abstraction in programming . . . 40,225
- Acoustics:**
- acoustic surface-wave filters . { 32,179
36, 29
 - magnetoacoustic effects in bismuth 32,233
 - acoustic waves in piezoelectric semiconductors 33,336
 - vibration patterns and radiation behaviour of loudspeaker cones 36, 1
 - non-rectangular reverberation chamber 37,176
 - see also **Audio**
- Actuator for anti-lock braking system 36, 74
- Adsorption on Si and Ge surfaces . 32,131
- Aerial, see **Antenna systems**
- Aircraft, measurement of ozone in . 38,131
- Airfields:**
- short-range radar for ground services 32, 13
 - MADGE guidance system for aircraft landing { 34,225
35,271
 - displays showing processed radar echoes * 40,218
- Air pollution, see **Environmental science**
- Air traffic, see **Airfields**
- Alkali-antimonide films for photocathodes 40, 19
- Alkali-germanosilicate glass, manufacture of fibres from 36,182
- Alkali-vapour dispenser 36, 16
- Alloys:**
- alloying behaviour of transition metals 33,149
 - influence of ionicity 33,196
 - ferromagnetic, electrical conduction in 35, 29
 - intermetallic compounds 36,136
 - heat of formation 36,217
 - NiFe, for read-out of magnetic tape 37, 48
 - amalgams for fluorescent lamps 38, 83
- Alumina, sintered, light transmission of 36, 47
- Aluminate host lattices for phosphors 37,221
- Aluminium:**
- electrodeposition of 39, 87
 - optically smooth, machined on high-precision lathe * 39,183
 - single-crystal, spark-machined to form test spheres 40,202
- Aluminium-gallium arsenide:**
- laser 39, 37
 - laser, CQL10, microscope photograph of * 39,324
- Amalgams for fluorescent lamps . . 38, 83
- Amplifier:**
- small-signal, MOS transistor as 31,216
 - integrated audio, with high input impedance and low noise . . 31,245
 - power, for h.f. band, MOS transistor as 31,251
 - integrated amplifier circuits . 32, 8
 - parametric, for interferometer . 32, 20
- Amplitude and phase contrast, simultaneous, in STEM 37, 1
- Amplitude modulation (AM) 36,309
- Analog computer for simulating one-dimensional antenna arrays 31, 2
- Analog signals:**
- shift register for ('bucket-brigade memory') { 31, 97
31,266
 - linear basic ICs for 32, 1
 - quantization and coding of . . . 36,337
- Anisotropy of electrical resistivity . 35, 29
- ANS, the Netherlands astronomical satellite:**
- general 33,117
 - attitude control 33,162
 - onboard computer 34, 1
 - UV experiment (Groningen University) 34, 33
 - X-ray experiment (Utrecht University) 34, 43
 - reaction wheels 34,106
 - sun sensors 34,208
 - horizon sensor 34,213
 - star sensor 34,218
- Antenna systems:**
- one-dimensional, simulation of 31, 2
 - electronically controlled, with *P-I-N* diodes 32,405
- Anti-lock braking system, fast actuator for 36, 74
- Argon-plasma torch, r.f., for emission spectroscopy 33, 50
- Arsenic oxides, action in glass fining 40,310
- Artificial languages in PHLIQA I . 38,269
- ASKA (Automatic System for Kinematic Analysis) used in stress calculations for TV tubes 37, 56
- Assembly robot, experimental . . . 40, 33
- Astronomy:**
- window for ultrasoft X-rays . 40, 12
 - see also **ANS** and **Radio-astronomy**
- Asynchronous motors, see **Induction motors**
- Audio:**
- cassettes 31, 77
 - integrated audio amplifier . . 31,245
 - annoyance due to modulation noise and drop-outs 37, 29
 - station and programme identification in FM broadcasting . 39,216
 - manipulation of speech sounds 40,134
 - Compact Disc Digital Audio . 40,149
 - see also **Acoustics**
- Auger electron spectroscopy 34,359
- AUROS, system for speaker recognition by computer 37,207
- Automation:**
- flexible, an experiment in 38,329
 - system that learns to recognize two-dimensional shapes 38,356
 - self-organizing systems 38,364
 - experimental assembly robot . 40, 33
 - flexible, software for 40,237
- Avalanche photodiode:**
- development, use of fast scanning microscope 35, 23
 - detector for optical communication 36,205
- AVOID, short-range high-definition radar 32, 13
- _____
- Background illumination for displays 40, 79
- Barium titanate:**
- ceramic, microstructure of . . . 32, 92
 - PTC effect in 38, 73
- Bearings:**
- helical-groove, grease-lubricated 34,103
 - spiral-groove, grease-lubricated { 35,137
39,184
 - 'push-pull' spiral-groove 35, 11
 - air, tapered, in metrology . . . 40,338
- Bi-aspheric lenses, made on COLATH 39,243
- Bioceramic of sintered hydroxylapatite 37,234
- Bismuth, magnetoacoustic effects in 32,233
- Bismuth-silicon-oxide crystals, large, pulling from the melt 37,250
- 'Bitter water', making tracks on video tape visible with 40,129
- Bloch walls in garnet films for fast bubbles 38,211
- 'BOL', nuclear measurement system at IKO 39,302
- Bonding, thermocompression, with IC bonding bit 40,200
- Bonding leadless electronic components to printed boards 40,342
- Boron filament as constructional material 35,125
- Bourdon gauge for hot and corrosive gases 39,344
- Braking system, anti-lock, actuator for 36, 74
- Bridgman anvils for very high pressures 36,251
- Brittle materials, grinding 38,105
- Broadband circulators for VHF and UHF 36,255

- Brushless d.c. motors:**
 with integrated Hall device . . . 31,366
 for ANS reaction wheels . . . 34,106
 for ANS horizon sensor . . . 34,213
 steady-state performance . . . 35,106
- Bubble memories:**
 preparation of garnet films for . . . 35, 1
 single-mask, with rotating-field control . . . 36,149
- Bubbles, see Magnetic bubbles**
- Bucket-brigade memory, see CTD**
- CAD, see Computer-Aided Design**
- Camera tubes, TV, new concept for** 39,201
- Capacitor motor, supply-voltage speed control** . . . 34,180
- Carbon foam** . . . 36, 93
- Cassettes, audio** . . . 31, 77
- Cathodoluminescence for investigating crystal defects** . . . 35,239
- Cavity resonator:**
 for MIP for emission spectroscopy . . . 39, 66
 for microwave measurement of moisture content . . . 40,116
- CCD, see CTD**
- Ceramic technology:**
 renaissance in (Prof. A. L. Stuijts's inaugural lecture) . . . 31, 44
 influence of microstructure on properties of electroceramics . . . 32, 79
 isostatic hot pressing without metallic encapsulating layer . . . 35, 65
 thick-film technology for hybrid-IC interconnection patterns . . . 35,144
 Sintering process . . . 35,188
 ceramic-to-metal bonding . . . 35,209
 sintered Al₂O₃ for sodium lamps . . . 36, 47
 manufacture of Ferroxdure . . . 37,165
 sintered hydroxylapatite as bioceramic . . . 37,234
 grinding brittle materials . . . 38,105
- Channel multipliers:**
 investigation of microchannel plates by scanning electron microscopy . . . 34,270
 image intensifier for hard X-rays . . . 37,124
 microchannel-plate photomultipliers with subnanosecond characteristics . . . 38,240
- Charge-coupled device, see CTD**
- Chemical analysis, see Chromatography and Inorganic chemical analysis**
- Chemical Vapour Deposition:**
 in the manufacture of silica-glass fibres . . . 36,185
 in the manufacture of products of pyrolytic graphite . . . 37,189
 in applying wear-resistant coatings to tool steel . . . 40,204
- Chemiluminescence, detectors based on** . . . 34, 73
- Chlorine, organo-, residues in milk, determination of** . . . 36,284
- Chromatography:**
 for determination of organochlorine residues in milk . . . 36,284
 combination of MIP with gas chromatograph . . . 39, 70
 chromatofocusing for separation of proteins . . . 39,125
 liquid, electrochemiluminescence cells as detector for . . . 40, 77
- Circulators, broadband, for VHF and UHF** . . . 36,255
 { 37,266
 { 37,272
- Clean rooms for LSI *** . . . { 37,272
- Clock radio with LSI *** . . . 37,302
- COD value for water** . . . 34,123
- Code modulation with digitally controlled companding** . . . 31,335
- COLATH, numerically controlled precision lathe** . . . 39,229
- Colour rendering:**
 of discharge lamps . . . 35,347
 optimum spectra for light sources . . . 35,361
- Colour television, see Television**
- Communication satellites, modulation in** . . . 36,359
- Commutator motors, see D.C. motors**
- Compact Disc Digital Audio:**
 special issue . . . { 40,149
 { -180
 general . . . 40,151
 system aspects and modulation . . . 40,157
 error correction and concealment . . . 40,166
 conversion from digital to analog . . . 40,174
- Companding, digitally controlled, for code modulation** . . . 31,335
- Compensation wall in ferrimagnetic materials** . . . 34, 96
- Complementary MOS circuits, LOCOS technology** . . . 34, 19
- Composites, growth of *** . . . 32,102
- Computer-Aided Design:**
 with interactive display . . . 36,162
 design of LSI circuits . . . 37,278
 data-base management for . . . 40,245
- Computer applications:**
 calculation of Gunn effect . . . 32,385
 { 33, 89
 reading handwritten numerals { 33,130
 analysis of emission spectra . . . 34,322
 { 34,330
 neutron activation analysis . . . { 34,351
 measurements from X-ray phototographs . . . 35,170
 testing digital circuits . . . 35,261
 designing a loudspeaker cone . . . 36, 14
 designing workpieces . . . 36,162
 simulation of control system for water supply . . . 36,273
 optical inspection of connecting-lead patterns for ICs . . . 37, 77
 speaker recognition . . . 37,207
 designing LSI circuits . . . 37,278
 designing TV deflection coils . . . 39,154
 document handling with Megadoc system . . . 39,329
 control of experimental assembly robot . . . 40, 41
 see also **Picture processing and Software**
- Computers:**
 analog, for simulating one-dimensional antenna arrays . . . 31, 2
 ANS onboard computer . . . 34, 1
 microprogramming for P1000 family of computers . . . 34,132
 modelling and simulation in design of . . . 39,134
 Constant-current source, integrated . . . 32, 4
- Contrast enhancement:**
 by histogram transformation . . . 38,300
 with DOT Scan CCTV . . . 38,319
- Control, see Measurement and Control**
- 'Controlled cascading' in frequency dividers** . . . 32,103
- Correction plates for projection TV** . . . 39, 15
- CQL10:**
 information read-out with . . . 39, 37
 microscope photograph . . . 39,324
- 'Cryogem', miniature Stirling refrigerator** . . . 32, 51
- Cryopumps in industrial vacuum technology** . . . 39,246
- Crystal defects:**
 vacancy clusters in dislocation-free Si and Ge . . . 34,244
 transport of, in solids . . . 35,181
 investigation of cathodoluminescence . . . 35,239
 influence on luminescence of GaP . . . 38, 41
- Crystals:**
 liquid, for numerical displays . . . 37,131
 bismuth-silicon-oxide, pulling from the melt . . . 37,250
 DKDP, for TITUS tube . . . 39, 50
 AlGaAs, for laser . . . 39, 37
 TTC-cut, for quartz-crystal resonators . . . 40, 1
- CTD, Charge-Transfer Device:**
 shift register for analog signals . . . 31, 97
 integrated, with MOS tetrodes . . . 31,266
 image sensor . . . 37,303
 P²CCD in storage oscilloscope . . . 40, 55
- Curvature of reflecting surfaces, measurement of** . . . 40,338
- CVD, see Chemical Vapour Deposition**
- DAMA (Drop-out Annoyance Measuring Apparatus) in magnetic sound recording** . . . 37, 36
- Data base:**
 computer program for data-base consultation in English . . . 38,228
 { 38,269
 management, for CAD and CAM . . . 40,245
- Data links, fast automatic equalizer for** . . . 37, 10
- Data modem, Sematrans 102 *** . . . 36,356
- Data processing, see Computer applications**
- Data transmission:**
 in MADGE aircraft-guidance system . . . 35,271
 modulation systems for . . . 36,349
 digital modulation stage for . . . 37,291
 baseband, echo canceller for . . . 39,102
 error control in mobile-radio data communication . . . 39,172
- D.C. motors:**
 general . . . 33,230
 with speed control . . . 34,163
 small, designing . . . 35, 96
 linear, with permanent magnets . . . 40,329
 see also **Brushless d.c. motors**
- Deflection, electro-optic, of laser beam** . . . 36,117
- Deflection coils:**
 for TV picture tubes . . . 32, 61
 for 30AX colour system . . . { 39,154
 { 39,277
- Delta modulation** . . . { 31,338
 { 36,339
 { 37,313
- Demonstration and test vehicle for exterior lighting *** . . . 35,307
- Diameter variations in wire, measurement of** . . . 31,111
- Diamond, spark machining of** . . . 40,202
- Diamond die, model of *** . . . 40,133
- Differential pulse-code modulation (DPCM)** . . . { 31,338
 { 36,338
- Diffractionmeter, single-crystal X-ray, PW 1100** . . . { 33, 61
 { 38,246
- Diffusion of donors and acceptors in silicon** . . . 35,181
- Digital circuits:**
 integrated, with MOS transistors . . . 31,277
 integrated, with low dissipation . . . 35,212
 in video telephone . . . 36,233
 see also **LSI**

- Digital signals:**
 transmission of 36,343
 error control in mobile radio . 39,172
 station and programme identification in FM sound broadcasting 39,216
 representation of documents by, with Megadoc 39,329
- Digital systems:**
 control of companding for code modulation 31,335
 tone generation for keyboard instrument 31,354
 automatic equalizer for data links 37, 10
 echo canceller for data transmission 39,102
 equalizer for echo reduction in Teletext 40,319
 see also **Picture processing**
- Digital-to-analog conversion in Compact Disc** 40,174
- Dilution refrigerators for continuous cooling in the millikelvin range** . . 36,104
- Dip soldering electronic components** 38,135
- Discharge lamps:**
 general 35,308
 with metal halides 35,347
 electrodes 35,356
 see also **Mercury-vapour lamps and Sodium lamps**
- Dislocations, formation and movement** 32,250
- Dispenser, sodium-vapour** 36, 16
- Dispersion measurements on optical fibres for communication** 36,211
- Distance measurement with microwaves (MADGE)** 35,271
- Distributed computer systems:**
 software aspects 40,262
 computations on arrays of processors 40,270
- DKDP crystals for TITUS tube** . . . 39, 50
- Document handling with Megadoc system** 39,329
- DOD (Droplet On Demand) principle in ink-jet printing** 40,192
- Domains:**
 bubbles, see **Magnetic bubbles**
 observation, in ferroelectrics and ferromagnetics, with SEM . . . 36, 18
- Domain structure:**
 compensation wall 34, 96
 bubble memories 35, 1
- Doping by ion implantation** 31,267
 39, 1
- DOR (Digital Optical Recording) in Megadoc system** 39,329
- Dosage of oxygen in gas systems** . . 31,112
- Double-crucible method for making alkali-germanosilicate glass fibres** 36,182
- Double-sideband modulation (DSB)** 36,311
- Droplet interferometry for investigating smooth surfaces** 33, 74
- 'Drop-outs' in magnetic sound recording** 37, 29
- Ear-lobe clip for physiological transducers** 33,102
- Echo compensator:**
 for data transmission 39,102
 for Teletext 40,319
- Echography, medical, transducer for** 38,195
- EDDY, computer program for design of electromechanical devices** . . . 39, 78
- Edge enhancement:**
 OCR camera with 37,180
- digital image enhancement by** . 38,304
- two-dimensional, with Dot Scan CCTV** 38,317
- Elastic and electro-elastic waves** . . 33,309
- Electric motors:**
 small, special issue { 33,213
 -271
 general 33,215
 moving-coil motors 33,244
 electromagnetic vibrators 33,249
 frequency-analog speed control 33,260
 small, special issue II { 34,153
 -189
 speed control for capacitor motors 34,180
 small, special issue III { 35, 77
 -123
 early history of 35, 77
 applications of stepping motors 35,104
 scaling laws 35,116
 see also:
D.C. motors
Induction motors
Synchronous motors
- Electrochemistry:**
 galvanic effects in wet-chemical etching 38,149
 energy production by photoelectrochemical processes 38,160
 electrodeposition of aluminium 39, 87
 electrochemiluminescence in electrolyte-free solutions . . . 40, 69
- Electrodeposition of aluminium** . . . 39, 87
- Electromagnetism:**
 theory 32,150
 scattering characteristics of crossed waveguides 32,165
 electromagnetic waves 33,309
- Electromechanical devices, design of, by computer** 39, 78
- Electromechanical transducers** . . . 40,358
- Electron accelerator, linear, at IKO** 39,325
- Electron-beam pattern generator for LSI** 37,334
- Electronic components:**
 innovation in electronic devices 32,117
 thermal behaviour of, during soldering 38,135
 leadless, attachment to printed boards 40,342
- Electronic keyboard instrument with digital tone generation** 31,354
- Electronics at IKO for nuclear physics** 39,312
- Electron-image projector for LSI** . . . 37,347
- Electron lithography for LSI** 37,334
- Electron microprobe:**
 general 34,370
 analysis of glass fibres * 40,349
- Electron microscope, see Scanning electron microscope**
- Electron multipliers, see Channel multipliers**
- Electron resists:**
 for manufacture of ICs 35, 41
 negative, for VLSI 39,346
- Electrophoretic coating of cathode with LaNi₄Cu powder** 36,267
- Electroradiography, medical** 39, 19
- Electrostatic printing** 36, 57
- Electrostriction, transducers based on** 40,358
- Emission spectrometry:**
 r.f. argon-plasma torch for 33, 50
 multielement analysis by 34,305
 automatic analysis of emission spectra 34,322
 microwave-induced plasma for AES 39, 65
- Endoscope, experimental, with miniature TV camera** 35,166
- Energy production by photoelectrochemical processes** 38,160
- Engineering technology at IKO** . . . 39,315
- Environmental science:**
 SO₂ monitoring network for Rhine estuary region 32, 33
 monitoring centre for air pollution * 33,194
 ozone, NO and NO₂ detectors . . . 34, 73
 monitoring quality of surface water 34,113
 measuring oxygen demand of water 34,123
 automatic extractor 35,196
 control of water-purification plant 36,273
 measurement of ozone in an aircraft 38,131
- Epitaxial GaAs:**
 growth of, for microwave electronics 32,380
 growth of, for lasers 36,194
- Equalizer:**
 fast automatic, for data links . . . 37, 10
 see also **Echo compensator**
- Ergonomic lathe*** 36,160
- Error correction:**
 in mobile-radio data communication 39,172
 and concealment in Compact Disc system 40,166
- ESCA (Electron Spectroscopy for Chemical Analysis)** 34,361
- Etching:**
 with ions 35,199
 experimental etching equipment 38, 51
 wet-chemical, galvanic effects in 38,149
 plasma, in IC technology 38,200
- Evacuated solar collector with heat pipe** 40,181
- Evaporation of thin films, cryopump for** 39,254
- EVOLUON** 31,187
- Extractor:**
 automatic, bubble-train 35,196
 bubble-train, in pesticide monitor for milk 36,286
- Extrusion of glass** 32, 96
- FEM (Finite-Element Method), in stress calculations on TV picture tubes** 37, 56
- Ferrimagnetic materials:**
 compensation wall in 34, 96
 epitaxial growth of iron garnets for bubble memories 35, 1
 garnets for magneto-optic memories 37,197
- Ferrites:**
 thermogravimetric analysis 31, 24
 in kicker magnets for proton beams 31, 66
 control of properties via microstructure 32, 83
 microwave integrated circuits on ferrite substrate 32,315
 sintering 35,191
 see also **Ferroxdure**
- Ferroelectrics:**
 ceramic, microstructure of 32, 90
 for detecting infrared 35,247
 domains in, observation with SEM 36, 18
 ceramic, for analog memory 37, 51
 LiNbO₃ for holographic information storage 37,109
 PTC effect of BaTiO₃ 38, 73
 DKDP crystals for TITUS tube 39, 50
 for electromechanical transducers with no hysteresis 40,358

- Ferromagnetic materials:**
 electrical conduction in ferromagnetic metals 35, 29
 domains in, observation with SEM 36, 18
- Ferroxdure:**
 in d.c. motor for washing machine 34,163
 research and technology 37,157
- Field-effect transistors:**
 development of FET electronics (Prof. Klaassen's inaugural lecture) 33,203
 GaAs, microwave 39,269
 see also **MOS transistors**
- Field emission of electrons and ions 33,277
 Filament, boron 35,125
- Filters:**
 based on acoustic surface waves { 32,179
 YIG, for microwaves 32,322
 inductorless 33,294
 Fining of glass 40,310
 Finite-element method, in calculating stress in TV tubes 37, 56
 Fluidics in dishwasher control 33, 29
 Fluorescence analysis, X-ray 34,339
- Fluorescent lamps:**
 long, electronic starter for 31, 54
 low-pressure discharges, general optimum colour rendering 35,361
 suspension technology for applying fluorescent coating 36,269
 behaviour of aluminate phosphors for 37,221
 amalgams for 38, 83
- Fluorine:**
 cyclic processes in halogen lamps { 35,302
 profiling filaments 35,332
- FM sound broadcasting, station and programme identification in 39,216
 Fracture faces of optical glass fibres { 37, 89
 39,245
- Frequency-analog systems:**
 for speed control 33,260
 general 34,288
- Frequency dividers:**
 with 'controlled cascading' 32,103
 for ultra-high frequencies 38, 54
- Frequency-division multiplex (FDM) 36,316
 Frequency modulation (FM) 36,318
 Frequency stabilization with TTC-cut crystals 40, 1
- Gallium antimonide/metal contacts,**
 tunnelling in 32,214
- Gallium arsenide:**
 tunnelling in contacts with metal epitaxial 32,380
 diode, combined with IR phosphors 34, 31
 investigation of crystal defects 35,239
 laser diodes of, for optical communication 36,190
 microwave field-effect transistor 39,269
- Gallium arsenophosphide, zinc-diffusion profiles in 37,121
 Gallium nitride, light-emitting diodes based on 37,237
- Gallium phosphide:**
 radiating centres in 32,261
 investigation of crystal defects 35,239
 zinc-diffusion profiles in 37,121
 influence of crystal defects on luminescence 38, 41
- Gallium selenide, stoichiometric analysis of 34,350
- Galvanic effects in wet-chemical etching 38,149
- Gamma spectrometry, neutron activation analysis { 34,330
 34,351
- Garnets:**
 YIG microwave filters 32,322
 iron-garnet films for bubble memories 35, 1
 magnetic anisotropy and magnetostriction of Ru and Ir ions in YIG 35,225
 ferrimagnetic, for magneto-optic memories 37,197
 garnet films for fast magnetic bubbles 38,211
- Gas-analysis system with cryopump 39,250
- Gas-discharge lamps, see **Discharge lamps**
- Gas discharges:**
 general 35,310
 low-pressure 35,321
 investigation by light scattering 35,344
- Gas kinetics 33, 43
 Gas laser, optical polarization in 32,190
- Gas-phase epitaxy:**
 growth of gallium arsenide 32,380
 fabrication of GaN diodes 37,237
- Geiger-Müller counters from IKO 39,296
 'Gems' of lithium niobate * 31, 23
- Germanium:**
 chemical behaviour of clean surfaces 32,131
 dislocation-free, vacancy clusters in 34,244
- Getters:**
 for metal-iodide high-pressure mercury-vapour lamps 35,354
 Th-Ce-Al system 36,138
- Glass:**
 extrusion of 32, 96
 transparent single-point turning of 39, 92
 fining 40,310
- Glass fibres:**
 optical communication with 36,178
 manufacture of 36,182
 coupling of lasers to 36,202
 index profiling 36,211
 obtaining smooth fracture faces 37, 89
 machine for hot splicing 38,158
 precision fracture of * 39,245
 components for glass-fibre circuits 40, 46
 analysis of, with electron microprobe * 40,349
- Glassy carbon * 31,369
- Graphite, pyrolytic, products of 37,189
- Greases for spiral-groove bearings 39,189
- Grinding brittle materials 38,105
- Group theory 32,152
- Gunn effect:**
 Gunn diode 32,370
 computer calculations 32,385
 oscillators and amplifiers 32,397
- Gyrator 33,305
- Hall device:**
 integrated, for brushless d.c. motors 31,366
 for measuring induction-motor torque 34,153
- Halogen lamps, see **Incandescent lamps**
- Handwriting:**
 analysis and synthesis 32, 73
 image processing of, with OCR camera 37,180
- Heartbeat, recording with ear-lobe clip 33,102
- Heat-flux transformers 33,141
 Heat of formation of alloys 36,217
- Heat pipes:**
 in Stirling engine 31,176
 operation and characteristics 33,104
 applications 33,138
 in evacuated solar collector 40,181
 in refrigerator-freezer 40,350
- HEIS (High-Energy Ion Scattering) for determining implantation profiles 39, 5
- Helical-groove bearings, see **Bearings**
- Helicon waves 33,323
 Helium II, vortices in 32,256
- High-frequency amplifiers, MOS transistors as 31,251
- High pressures, apparatus for research at 36,245
- High-pressure sodium lamps { 35,334
 39,211
- High-speed spark-machining 40,199
- Historical:**
 40 years of workshop technology 31,127
 electric motors 35, 77
 20 years of research on intermetallic compounds 36,136
 IKO, the Institute for Nuclear Physics Research 39,286
 Holst, Gilles, pioneer of industrial research in the Netherlands 40,121
- Holography:**
 holographic strain analysis 35, 53
 holographic display of vibration patterns 36, 8
 LiNbO₃ for holographic information storage 37,109
- Holst, Gilles, pioneer of industrial research in the Netherlands 40,121
- Hot-gas engine 31,169
- Hot pressing of ceramics, without metallic encapsulating layer 35, 65
- Hot-spot model for burn-out of incandescent lamps 35,296
- HPI lamps 35,347
- Hybrid ICs, connection patterns for 35,144
- Hydrogen:**
 in discharge lamps 35,354
 storage in LaNi₅ 36,145
 absorption in intermetallic compounds 36,226
- Hydroxylapatite, sintered, as bioceramic 37,234
- IC, see **Integrated circuits**
- I²L, see **Injection Logic**
- Image intensifier for hard X-rays 37,124
- Image processing, see **Picture processing**
- Image projector, electron, for LSI 37,347
- Image sensor with resistive electrodes 37,303
- IMPATT diodes:**
 general 32,328
 in microwave oscillators 32,345
 anomalous oscillations in 32,361
- Incandescent lamps:**
 non-destructive measurement of pressures in { 31, 93
 32,155
 burn out { 35,296
 temperature distribution in 32,206
 research 35,295
 halogen, cyclic processes in 35,302
 carbon-filament, with halogen filling 35,316
 temperature profiling of filament 35,332

- INDA, software for flexible automation 40,237
- Index profiling of glass fibres 36,211
- Induction motors:**
- general 33,226
 - torque measurements on 34,153
 - high-speed, with solid rotor 34,170
- Inductors, simulation of, in filters 33,294
- Industrial research:**
- innovation in electronic devices 32,117
 - and theoretical physics 32,149
 - Holst, Gilles, pioneer of, in the Netherlands 40,121
- Infrared:**
- phosphors 34, 24
 - horizon sensor in ANS satellite 34,213
 - pyroelectric detectors 35,247
 - thermography 37,241
- Injection logic:**
- integrated (I²L) 33, 76
 - low-dissipation 35,217
 - CAD system for I²L circuits 37,290
 - I²L circuit for digital modulation stage 37,291
- Ink-jet printing 40,192
- Inorganic chemical analysis:**
- detectors for ozone, NO and NO₂ 34, 73
 - special issue { 34,297
-374
 - general 34,298
 - multielement analysis by optical emission spectrometry 34,305
 - automatic analysis of emission spectra 34,322
 - neutron activation analysis { 34,330
34,351
 - line-intensity calculations for X-ray fluorescence analysis 34,339
 - mass-spectrographic determination of C, O and N 34,344
 - stoichiometric analysis of GaSe 34,350
 - determining the content of active oxygen in oxides * 34,356
 - surface analyses 34,357
 - electron microprobe 34,370
 - simple extractor 35,196
 - measurement of ozone in aircraft 38,131
 - emission spectrometry with MIP 39, 70
- Institute for Nuclear Physics Research (IKO):**
- special issue { 39,285
-328
 - historical 39,286
 - cyclotron 39,290
 - radioisotopes 39,294
 - Geiger-Müller counters 39,296
 - semiconductor detectors 39,298
 - BOL 39,302
 - electronics for nuclear physics 39,312
 - engineering technology 39,315
 - ion-beam technology 39,320
 - linear electron accelerators 39,325
- Integrated circuits:**
- bucket-brigade delay line, shift register for analog signals 31, 97
 - Hall device for brushless d.c. motors 31,366
 - linear basic 32, 1
 - integrated injection logic (I²L) 33, 76
 - mounted with leads on tape 34, 85
 - making photomasks with Optycograph 34,257
 - hybrid, connection patterns for 35,144
 - P-N-P-N elements for telephone switching 36,291
 - optical inspection of connecting-lead patterns 37, 77
 - abbreviations and acronyms 37,277
 - in PM 2517 digital multimeter 38,181
 - plasma etching in IC technology 38,200
 - echo canceller for data transmission 39,102
 - equalizer for echo reduction in Teletext 40,319
 - see also:
 - LOCOS technology**
 - LSI**
 - MOS technology**
 - VLSI**
- Integrated microwave circuits:**
- general 32,292
 - with lumped elements 32,305
 - on ferrite substrate 32,315
- Interactive display, in computer-aided design 36,162
- Interconnection patterns:**
- on tape 34, 85
 - in thick-film technology 35,144
 - for ICs, optical inspection of 37, 77
- Interferometry:**
- droplet, for investigating smooth surfaces 33, 74
 - in MADGE aircraft-landing guidance system 34,225
- Intermetallic compounds, 20 years of research 36,136
- Ion-beam system (200 kV) at FOM * 39,319
- Ion etching 35,199
- Ion implantation:**
- in MOS transistors 31,267
 - in semiconductors 39, 1
- Ionography with compressed gas, system for medical electroradiography 39, 21
- Isothermal spaces (with heat pipe) 33,139
- Isotopes, radioactive, for determination of zinc-diffusion profiles in GaP and GaAsP 37,121
-
- Kerr cells in deflector for laser beam 36,117
- Keyboard instrument, transposing, with digital tone generation 31,354
- Kicker magnets for proton beams 31, 66
-
- Laboratories, Philips Research ('Nat.Lab'), at Waalre 31,153**
- Lacquers, photopolymerizable, for LaserVision discs 40,298
- Lamps, see **Discharge lamps and Incandescent lamps**
- Langmuir trough for building monomolecular layers 36, 44
- Large-screen projector with laser beam 36,117
- Lasers:**
- gas, optical polarization in 32,190
 - electro-optic deflection of laser beam 36,117
 - semiconductor, for optical communication 36,190
 - semiconductor, coupling to glass fibre 36,202
 - semiconductor, for information read-out 39, 37
 - microscope photograph of CQL * 39,324
 - see also **LaserVision**
- LaserVision:**
- Philips system { 33,177
-193
 - general 33,178
 - signal processing 33,181
- scanning system of player 33,186
- control mechanism of player 33,190
- modulation method 36,332
- scanner for player * 37, 90
- disc manufacture 40,287
- photopolymerizable lacquers for discs 40,298
- Lateral skin effect in a flat conductor 32,221
- Lathe:**
- ergonomic * 36,160
 - high-precision, COLATH 39,229
- Lead-oxide/gas-layer system for medical electroradiography 39, 25
- Lectures:**
- H. B. G. Casimir 40,121
 - Bernard Dixon 38, 17
 - F. M. Klaassen 33,203
 - W. Martienssen 38, 25
 - A. E. Pannenberg 38, 33
 - G. W. Rathenau 32,117
 - A. L. Stuijts 31, 44
 - H. J. Vink 35,181
- LED, see **Light-emitting diodes**
- LEED (Low-Energy Electron diffraction) 34,358
- Levitation theory 34, 67
- Life of incandescent lamps 32,155
- Light:**
- special issue { 35,293
-371
 - research on incandescent lamps 35,295
 - discharge lamps 35,308
 - carbon-filament lamps with chemical transport cycle 35,316
 - low-pressure gas discharges 35,321
 - profiling tungsten filaments 35,332
 - high-pressure sodium lamp 35,334
 - scattering, in gas discharges 35,344
 - metal-halide discharge lamps 35,347
 - getter for metal-iodide high-pressure mercury-vapour lamps 35,354
 - electrodes in discharge lamps 35,356
 - optimum spectra for light sources 35,361
- Light beam, deflection of, with Kerr cells 36,117
- Light-emitting diodes:**
- physics of centres in GaP 32,261
 - combination of IR phosphors with GaAs diode 34, 31
 - for optical communication 36,190
 - of GaN 37,237
- Lighting, see:
- Discharge lamps**
 - Incandescent lamps**
 - Light**
 - Road lighting**
- Light modulation in optical-fibre transmission systems 36,201
- Light transmission of sintered alumina 36, 47
- Liquid chromatography, see **Chromatography**
- Liquid crystals for numerical displays 37,131
- Liquid-phase epitaxy:**
- for GaAs 32,382
 - for growing iron-garnet films 35, 1
 - in making semiconductor lasers 36,194
- Lithium niobate:**
- 'gems' of * 31, 23
 - for holographic information storage 37,109
- Lithography:**
- methods for IC production 37,270
 - electron-beam pattern generator 37,334
 - electron-image projector 37,347
- LOCOS technology:**
- for complementary MOS circuits 34, 19
 - CAD for 37,289

- LOCOS technology:**
 for MOS transistors 31,234
 for bipolar transistors * 31,276
 see also **LOCOS technology**
- Logic analyser, PM 3543 *** 40,286
- Logic circuits:**
 fluidics, for dishwasher control 33, 29
 large, structural test methods for 35,261
 LSI 37,274
 see also **Injection logic**
- Loudspeaker cones, vibration patterns and radiation behaviour of** 36, 1
- LPE, see Liquid-phase epitaxy**
- LSI:**
 testing large digital circuits 35,261
 special issue { 37,265
 -356
 revolution in electronics 37,267
 abbreviations and acronyms 37,277
 design by CAD 37,278
 digital modulation stage for data transmission 37,291
 image sensor with resistive electrodes 37,303
 delta modulation/PCM converter 37,313
 Silicon Repeater 37,330
 electron-beam pattern generator 37,334
 electron-image projector 37,347
- Luminescence:**
 characteristic, special issue { 31,303
 -332
 absorption and emission spectra of important activators 31,304
 efficiency of phosphors excited in activator 31,314
 energy transfer and efficiency 31,324
 electro- { 32,261
 37,237
 screen with 'pillar structure' * 33,161
 chemi-, detectors based on 34, 73
 cathodo-, in investigating crystal defects 35,239
 of GaP, influence of crystal defects on 38, 41
 electrochemi-, in electrolyte-free solutions 40, 69
 see also **Phosphors**
- Lumped-element microwave components** 32,305
- Machine-print characters, image processing of, with OCR camera** 37,180
- Machining, transparent, of glass** 39, 92
- MADGE guidance system for aircraft landing** { 34,225
 35,271
- Magnetic bubbles:**
 detection of 36,157
 generation and manipulation by microwaves 37, 38
 fast, garnet films for 38,211
- Magnetic deflection in TV picture tubes** 32, 61
- Magnetic domains, see Magnetic bubbles**
- Magnetic fields of TV deflection coils, measurement of** 39,277
- Magnetic materials:**
 for permanent magnets 34,193
 anisotropy and magnetostriction of Ru and Ir ions in YIG 35,225
 see also:
Ferrimagnetic materials
Ferrites
Ferromagnetic materials
- Magnetic recording:**
 audio tape cassettes 31, 77
 of video signals 36,326
 read-out of magnetic tape by magnetoresistance effect 37, 42
 making video tracks visible 40,129
- Magnetic sound recording, annoyance due to modulation noise and drop-outs in** 37, 29
- Magnetic tape, production of** 36,268
- Magnetism:**
 photomagnetic effects 31, 33
 magnetoacoustic effects in bismuth 32,233
 magnetic liquids * 33,293
 levitated magnets 34, 67
 the compensation wall in ferromagnetic materials 34, 96
 measuring magnetic properties of microscopic particle 39, 48
 see also **Bubble memories and Magnetic bubbles**
- Magneto-optical memories** 34, 96
- Magneto-optic storage wafer MOPS** 37,197
- Magnetometer, vibrating-reed** 31, 40
- Magnetoresistance effect, read-out of magnetic tape by** 37, 42
- Magnetostriction of Ru and Ir ions in YIG** 35,225
- Magnetron sputtering, cryopump for Mammograms, computerized processing of** 39,252
 38,347
- Manometers, see Pressure measurements**
- Masks for integrated circuits:**
 making photomasks with Optycograph 34,257
 for fabrication of digital modulation stage 37,300
 electron-beam pattern generator for LSI 37,334
 removal of, by plasma etching 38,208
 checking computer design of * 38,290
- Mass spectrometry:**
 determination of C, O and N in semiconductors 34,344
 analysis of solid surfaces 34,362
- Materials, non-destructive examination of** 38,298
- Measurement and control:**
 attitude control of ANS satellite { 33,162
 34,208
 control systems in LaserVision player { 33,186
 33,190
 with frequency-analog signals { 33,260
 34,288
 reaction wheels for ANS satellite 34,106
 servocontrol of Optycograph 34,263
 control of water-purification plant 36,273
 control of experimental robot 40, 41
- Mechanization:**
 mounting ICs and conductor patterns on tape 34, 85
 see also **Automation**
- Medical technology:**
 ear-lobe clip for physiological transducers 33,102
 experimental endoscope with miniature TV camera 35,166
 measurements on X-ray pictures transducer for medical echography 38,195
 digital enhancement of X-ray pictures 38,298
 flashing tomosynthesis of X-ray pictures 38,338
 mammogram processing 38,347
 electroradiography 39, 19
 TOMOSCAN 310 * 40,253
- Megadoc, modular system for electronic document handling** 39,329
- Memories:**
 bucket-brigade delay line, shift register for analog signals { 31, 97
 31,266
 MOS store with discretionary wiring 31,286
 analog, ferroelectric 37, 51
 holographic, LiNbO₃ for 37,109
 magneto-optic, of discretionary type 37,197
 DOR disc in Megadoc modular system 39,329
 storage oscilloscope with P²CCD 40, 55
 see also **Bubble memories**
- Mercury-vapour lamps:**
 low-pressure, with fluorescence, see **Fluorescent lamps**
 sorting quartz tubing for 32, 57
 high-pressure, metal-iodide, hydrogen getter for 35,354
- Metal-halide lamps** 35,347
- Metallurgical quantities, values of** 38,257
- Metals:**
 model for alloys { 33,149
 33,196
 surface imaging by field emission ferromagnetic, electrical conduction in 35, 29
 metallization technique for ceramic-to-metal bonding 35,209
 20 years of research on intermetallic compounds 36,136
 heat of formation of alloys 36,217
 interconnection in ICs 37,273
 atom as building block 38,257
- Microchannel plate in subnanosecond photomultiplier tubes** 38,240
- Microcomputer Development System PM 4421 *** 40,269
- Microcomputer for intercom system *** 37,275
- 'Micro Cryogem' miniature Stirling refrigerator** 32, 51
- Microfilms by PD process** 33, 1
- Microprobe, electron** 34,370
- Microprogram control of P1000 family of computers** 34,132
- Microscopy:**
 imaging metal surfaces by field emission 33,277
 fast SEM in development of avalanche photodiodes 35, 23
 see also **Scanning electron microscope**
- Microstructure of electroceramic materials** 32, 79
- Microwaves:**
 cross-coupled waveguides 32,165
 special issue { 32,281
 -412
 solid-state electronics for 32,282
 integrated microwave circuits 32,292
 lumped-element circuits 32,305
 integrated circuits on ferrite substrate 32,315
 YIG filters 32,322
 IMPATT diodes 32,328
 IMPATT oscillators 32,345
 anomalous oscillations in IMPATT diode 32,361
 epitaxial GaAs for microwave devices 32,380
 P-I-N diodes in phase-shifters for antenna arrays 32,405
 MADGE aircraft-guidance system { 34,225
 35,271
 circulators, broadband 36,255
 modulation in microwave links 36,357
 generation and manipulation of magnetic domains 37, 38

- microwave-induced plasma for emission spectroscopy . . . 39, 65
 GaAs microwave FET . . . 39,269
 TRAPATT oscillator . . . 40, 99
 measurement of moisture content in process materials . . . 40,112
 see also **Gunn effect and Radar**
 MIG welding with plasma . . . 33, 21
 Milk, determination of organochlorine residues in . . . 36,284
 MIP (Microwave-Induced Plasma) for emission spectroscopy . . . 39, 65
 Mobile-radio data communication, error control in . . . 39,172
Modulation:
 code, with digitally controlled companding . . . 31,335
 light, in optical-fibre transmission systems . . . 36,201
 special issue . . . { 36,305
 -362
 introduction . . . 36,307
 of sinusoidal carrier . . . 36,309
 of pulse trains . . . 36,329
 quantization and coding of analog signals . . . 36,337
 transmission of digital signals . . . 36,343
 in telecommunication . . . 36,353
 in data transmission . . . 37, 13
 digital modulation stage for data transmission . . . 37,291
 delta modulation/PCM converter with LSI . . . 37,313
 in Compact Disc . . . 40,157
 Modulation noise in magnetic sound recording . . . 37, 29
 Moisture content of process materials, measured with microwaves . . . 40,112
 Monitoring centre for air pollution * 33,194
 Monitoring network for SO₂ in Rhine estuary region . . . 32, 33
 Monomolecular layers, Langmuir trough for . . . 36, 44
 MOPS, magneto-optic storage wafer 37,197
MOS technology:
 general . . . 31,225
 with ion implantation . . . 31,267
 in thin monocrystalline Si layers 31,271
 store with discretionary wiring . . . 31,286
 LOCOS, for complementary circuits . . . 34, 19
 ICs with low dissipation . . . 35,212
 see also **MOS transistors**
MOS transistors:
 special issue . . . { 31,205
 -295
 general . . . 31,206
 operation and d.c. behaviour . . . 31,209
 as small-signal amplifier . . . 31,216
 LOCOS technology . . . 31,234
 carrier mobility . . . 31,237
 in integrated audio amplifiers . . . 31,245
 in integrated chopper circuit . . . 31,248
 as r.f. power amplifiers . . . 31,251
 MOS tetrodes for UHF . . . 31,259
 bucket-brigade delay line with MOS tetrodes . . . 31,266
 in digital ICs . . . 31,277
 see also **Field-effect transistors**
 Moving-coil motors . . . 33,244
 MRH (Magneto-Resistive Head) for read-out of magnetic tape . . . 37, 42
 Multimeter, PM2517 automatic digital . . . 38,181
-
- Network theory for inductorless filters 33,294
 Neutron activation analysis . . . { 34,330
 34,351
- NIRMS (inert-gas-ion reflection mass spectrometry) . . . 34,363
 NO and NO₂ detectors . . . 34, 73
Noise:
 in X-ray films, determining . . . 31,117
 annoyance due to modulation noise in magnetic sound recording . . . 37, 29
 Numerals, handwritten, optical character recognition of . . . 33,130
 Numerical display with liquid crystals 37,131
-
- O-BUS, system for flexible public transport by on-call buses . . . 40,231
 OCR, see **Optical character recognition**
 Open-air laboratory for road lighting * . . . 35,320
 Ophthycograph, for optical drawing of photomasks . . . 34,257
Optical character recognition:
 automatic reading of hand-written numerals . . . 33, 89
 edge-enhancing double-focus OCR camera . . . 37,180
 see also **Picture processing and Recognition**
Optical communication via glass fibres:
 special issue . . . { 36,177
 -216
 general . . . 36,178
 manufacture of glass fibres . . . 36,182
 semiconductor lasers for . . . 36,190
 light modulation and injection . . . 36,201
 avalanche photodiode as detector . . . 36,205
 testing optical fibres by dispersion measurements . . . 36,211
 Optical inspection of connecting-lead patterns for ICs . . . 37, 77
Optical technology:
 turning glass for aspheric lenses 39, 92
 making bi-aspheric lenses with COLATH . . . 39,243
 Organochlorine residues in milk, determination of . . . 36,284
 Orthogonal transformation of TV pictures in real time . . . 38,119
 Oscilloscope, digital storage, with P²CCD . . . 40, 55
Oxygen:
 monitoring low partial pressures 31,112
 oxygen demand of water . . . 34,123
 determining active-oxygen content of oxides * . . . 34,356
Ozone:
 detector based on chemiluminescence . . . 34, 73
 measurement of, in aircraft . . . 38,131
-
- Parallel computer programs . . . { 40,254
 40,278
 Parametric amplifiers for radio astronomy . . . 32, 20
 Passivation layers, measurement of stress in . . . 39,130
 Pattern generator, electron-beam, for LSI . . . 37,334
 P²CCD in storage-oscilloscope . . . 40, 55
PD process:
 general . . . 33, 1
 connection patterns for ICs . . . 34, 85
- Perception:**
 analysis and synthesis of handwriting . . . 32, 73
 measurement of visual conspicuity 36, 71
 speaker recognition by computer 37,207
 OCR camera with edge enhancement for image processing . . . 37,180
 manipulation of speech sounds 40,134
Permanent magnets:
 materials for . . . 34,193
 intermetallic compounds for . . . 36,140
 Ferroxdur[®] for . . . 37,157
 linear d.c. motor with . . . 40,329
 Personal computer, P 2000 * . . . 40,261
 Pesticides, milk monitor for . . . 36,284
 Phase and amplitude contrast, simultaneous, in STEM . . . 37, 1
 Phase diagrams for ferrites . . . 31, 24
 Phase modulation . . . 36,318
 Phase-shifters with *P-I-N* diodes, for antenna arrays . . . 32,405
 Phenol synthesis and photomorphogenesis . . . 38, 89
 PHIDAS, data-base management system for CAD/CAM . . . 40,245
 PHIDIAS, software aspects of . . . 40,262
PHLIQA 1:
 organization and performance . . . 38,229
 artificial languages and translation operations . . . 38,269
Phosphors:
 for colour TV . . . 32,125
 infrared . . . 34, 24
 aluminate, for fluorescent lamps, behaviour of . . . 37,221
 pigmentation of, for colour TV 40, 48
 see also **Luminescence**
Photocathodes:
 caesium-iodide, for electron-image projector . . . 37,348
 alkali-antimonide films for . . . 40, 19
 Photodiodes, avalanche . . . { 35, 23
 36,205
 Photoelectrochemical energy production . . . 38,160
 Photoemission of alkali-antimonide films . . . 40, 19
 Photography, PD process . . . 33, 1
 Photomagnetic effects . . . 31, 33
 Photomasks, see **Masks for integrated circuits**
 Photomorphogenesis . . . 38, 89
 Photomultipliers, ultra-fast, with microchannel plate . . . 38,240
Photopolymerization:
 in manufacture of LaserVision discs . . . 40,287
 photopolymerizable lacquers for LaserVision discs . . . 40,298
 Phototitit optical converter . . . 34,274
 Physiological transducers, ear-lobe clip for . . . 33,102
Picture processing:
 Phototitit, an optical converter for . . . 34,274
 camera for, with edge enhancement . . . 37,180
 special issue . . . { 38,289
 -371
 a challenge . . . 38,291
 digital image enhancement . . . 38,298
 Dot Scan CCTV for, in real time 38,310
 experiment in flexible automation . . . 38,329
 digital, universal instrument for 38,326
 flashing tomosynthesis . . . 38,338
 computerized mammogram processing . . . 38,347
 system that learns to recognize two-dimensional shapes . . . 38,356
 self-organizing systems . . . 38,364

- Picture tubes, see **Television**
- Piezoelectrics:**
ceramic, microstructure of . . . 32, 91
waves in 33,310
- Pigmentation of TV phosphors . . . 40, 48
- P-I-N* diodes in phase shifters for antenna arrays 32,405
- PINPOINT** radio system for locating and monitoring vehicles 35, 15
- Plants, phenol synthesis and photomorphogenesis in 38, 89
- Plasma:**
MIG welding 33, 21
r.f. argon-plasma torch for emission spectroscopy 33, 50
torch for heating quartz glass * 34, 60
plasma etching in IC technology 38,200
microwave-induced, for emission spectrometry 39, 65
- Plasma method for making glass fibre 36,185
- Plumbicon TV camera tube, new concepts 39,201
- Pockels effect 34,274
- Polarization, optical, in gas laser . . 32,190
- Power amplifiers, high-frequency, MOS transistors as 31,251
- Pressing, hot, of ceramics 35, 65
- Pressure measurements:**
non-destructive, in lamps . . . 31, 93
monitoring oxygen partial pressures down to 10^{-30} bar . . . 31,112
miniature pressure transducers with silicon diaphragm . . . 33, 14
fast gauge for hot and corrosive gases 39,344
- Pressure vessels for solid-state research 36,245
- Printing:**
electrostatic 36, 57
ink-jet 40,192
- Products of pyrolytic graphite . . . 37,189
- Projection TV, correction plates for 39, 15
- Proteins, separation of, by chromatofocusing 39,125
- Proton synchrotrons:**
kicker magnets for 31, 66
stabilization of proton orbits . . 34, 64
- PSD (Picture Store and Display) instrument for digital picture processing 38,326
- PTC effect of BaTiO_3 38, 73
- Pulse-amplitude modulation (PAM) 36,330
31,335
36,337
37,313
- Pulse-code modulation (PCM) . . . 36,337
37,313
- Pulse-frequency and pulse-duration modulation (PFM and PDM) . . . 36,332
- Purification plant for water, control of 36,273
- PW 1100 single-crystal X-ray diffractometer 38,246
- Pyroelectric infrared detectors . . . 35,247
- Quartz-crystal resonators with TTC-cut crystals 40, 1
- Quartz-glass fibre, manufacture . . . 36,185
- Quartz tubing, sorting, for mercury-vapour lamps 32, 57
- Question-answering system for data-base consultation in English { 38,229
38,269
- Radar:**
simulation of antenna arrays by analog computer 31, 2
traffic-flow analysis by 31, 17
AVOID, short-range high-definition radar 32, 13
- display screens of 'Signal' Automatic Radar Processing System * 40,218
- Radio astronomy:**
simulation of linear antenna arrays by analog computer . . . 31, 2
parametric amplifiers for radio interferometer 32, 20
- Radioisotopes from IKO 39,294
- Radio system for vehicle location (PINPOINT) 35, 15
- Reaction wheels for ANS satellite . . 34,106
- Reading, automatic, of handwritten numerals 33,130
- Recognition:**
of speakers by computer . . . 37,207
of two-dimensional shapes, by learning system 38,356
picture, with self-organizing systems 38,364
- Refractometer for thin films 35,142
- Refrigeration technology:**
cooling of parametric amplifier 32, 27
miniature refrigerators for electronic devices 32, 49
continuous cooling in the millikelvin range 36,104
critical review 37, 91
refrigerator-freezer with heat pipe 40,350
- Regenerative solar cells 38,166
- Research Laboratories, Philips, at Waalre 31,153
- Resistive electrodes for image sensor 37,303
- Resists:**
electron, for the manufacture of ICs 35, 41
electron, negative, for VLSI . . . 39,346
- Resonators with TTC-cut quartz crystals 40, 1
- Reverberation chamber, non-rectangular 37,176
- R.F. generator for argon-plasma torch 33, 54
- Rhine estuary region, SO_2 monitoring network for 32, 33
- Road lighting:**
open-air laboratory * 35,320
demonstration and test vehicle * 35,307
- Robot, experimental, for assembly . 40, 33
- Rotating-field control, single-mask bubble memory with 36,149
- SAMEN/SAMO 39,134
- Satellite, astronomical, see ANS
- Satellite television, 12 GHz receivers for 39,257
- Scaling laws for electric motors . . 35,116
- Scanning electron microscope:**
investigating microchannel plates with 34,270
electron microprobe { 34,370
40,349
PSEM 500 35,153
investigation of crystal defects by cathodoluminescence . . . 35,239
observations of domains in ferroelectrics and ferromagnetics . . 36, 18
scanning transmission electron microscope 37, 1
- Scanning microscope for development of avalanche photodiodes . . 35, 23
- Scanning transmission electron microscope, simultaneous phase and amplitude images in 37, 1
- SEM, see **Scanning electron microscope**
- Sematrans 102 * 36,356
- Semiconductors:**
ceramic, effect of microstructure 32, 92
tunnelling in metal-semiconductor contacts under pressure . . 32,211
in microwave electronics . . . 32,282
mass-spectrographic determination of C, N and O in 34,344
light-emitting diodes of GaN . . 37,237
for photoelectrochemical energy production 38,160
ion implantation in 39, 1
detectors from IKO 39,298
alkali-antimonide films for photocathodes 40, 19
see also:
Aluminium-gallium arsenide
Gallium arsenide
Gallium phosphide
Germanium
Silicon
- Semiconductor technology, see:
Gas-phase epitaxy
Ion implantation
Liquid-phase epitaxy
LOCOS technology
MOS technology
Shape memories of NiTe 36,143
Shaver, electric, vibrator for . . . 33,251
Shift register for analog signals (bucket-brigade delay line) . . { 31, 97
31,266
Sigma-delta modulation 38,186
Silica-glass fibre, manufacture . . . 36,185
- Silicon:**
surfaces, clean, behaviour of . . 32,131
diaphragm, in miniature pressure transducers 33, 14
dislocation-free, vacancy clusters in 34,244
diffusion of donors and acceptors in 35,181
avalanche photodiode of, for optical communication . . . 36,205
see also **Integrated circuits**
- Silicon Repeater for LSI 37,330
- SIMS:**
for surface analysis 34,366
for determining doping profiles 39, 8
- Single-mask bubble memory with rotating-field control 36,149
- Single-sideband modulation 36,313
- Sintering, see **Ceramic technology**
- Skin effect, lateral, in flat conductors 32,221
- SO_2 monitoring network in Rijnmond 32, 33
- Sodium lamps:**
low-pressure 35,321
high-pressure 35,334
high-pressure, new generation . . 39,211
- Sodium-vapour dispenser 36, 16
- Software:**
calculations on cross-coupled waveguides 32,165
control of X-ray diffractometer 33, 61
microprogram control for P1000 computer family 34,132
ASKA, for stress calculations on TV tubes 37, 56
for data-base consultation in English { 38,229
38,269
EDDY, for design of electromechanical devices 39, 78
SAMEN/SAMO, for simulating computer systems 39,134
for simulating telephone cables 40, 85
special issue { 40,217
-286
general 40,219
abstraction 40,225
O-BUS, system for public transport with on-call buses . . . 40,231
INDA, software tool for production engineer 40,237

- data-base management system for CAD/CAM 40,245
- parallel programs 40,254
- software aspects of PHIDIAS system 40,262
- PM 4421, aid in microcomputer development * 40,269
- distributed computations on arrays of processors 40,270
- transformation methods for improving parallel programs 40,278
- Solar cells, regenerative 38,166
- Solar collector with heat pipe 40,181
- Soldering:**
 - thermal behaviour of electronic components during soldering 38,135
 - of leadless components to printed boards 40,342
- Solidification, directional, in growing composites * 32,102
- Solid-state electronics for microwaves 32,282
- Solid-state research at low temperatures 37, 96
- SON lamps (high-pressure sodium lamps) { 35,334
39,211
- Sound, see **Acoustics and Audio**
- Space-division multiplex (SDM) system for telephone switching 36,294
- Space science:**
 - boron filament as constructional material 35,125
 - grease-lubricated spiral-groove bearings 35,137
 - see also **ANS**
 - Spark machining 40,199
 - SPARX system, speech studies with 40,134
 - Speaker recognition by computer 37,207
 - Spectra, optimum, for light sources 35,361
- Spectrometry:**
 - gamma, neutron activation analysis { 34,330
34,351
 - HEIS for determining doping profiles 39, 5
 - Raman, for investigating glass fining 40,310
 - see also:
 - Emission spectrometry**
 - Mass spectrometry**
 - SIMS**
 - Speech sounds, manipulation of 40,134
 - SPI (Station and Programme Identification) in FM sound broadcasting 39,216
 - Spiral-groove bearings, see **Bearings**
 - Splicing glass fibres for optical communication 38,158
 - Stabilization by oscillation 34, 61
 - Starter for long fluorescent lamps 31, 54
 - Station and programme identification 39,216
 - STEM (Scanning Transmission Electron Microscope) 37, 1
 - Stepping-motor applications 35,104
 - Stirling engine, prospects for vehicular propulsion 31,169
 - Stoichiometric analysis of GaSe 34,350
- Strain gauges:**
 - miniature, in Si diaphragm 33, 14
 - thin-film 39, 94
- Structural testing of digital circuits 35,261
- Superconductivity:**
 - vorticity in superconductors 32,252
 - in intermetallic compounds 36,139
 - and superfluidity 37, 98
- Surfaces:**
 - of clean Si and Ge, chemical behaviour of 32,131
 - metal, imaging by field emission 33,277
 - analysis of, general 34,357
 - surface segregation in alloys 36,228
 - aluminium, optically smooth * 39,183
 - reflecting, curvature of 40,338
 - Surface tension of metals 38,257
- Surface waves:**
 - acoustic, in filters { 32,179
36, 27
 - acoustic, in piezoelectric materials 33,310
 - Suspension technology 36,264
 - Synchrocyclotron at IKO { 39,291
39,308
- Synchronous motors:**
 - general 33,217
 - small single-phase, stability of 33,235
- Tape recorders:**
 - for audio cassettes 31, 87
 - moving-coil motors for 33,247
- Telecommunication:**
 - broadband circulators for VHF and UHF 36,255
 - modulation in 36,353
 - see also **Data transmission and Telephony**
- Telephony:**
 - video telephony, see separate entry
 - switching telephone signals 36,291
 - delta modulation/PCM converter 37,313
 - computer-aided research on multiwire cables 40, 85
 - digital exchange * 40,224
 - see also **Optical communication via glass fibres**
- Teletext:**
 - LSI converter for reception of * 37,312
 - automatic equalizer for echo reduction 40,319
- Television:**
 - magnetic deflection for picture tubes 32, 61
 - miniature TV camera in experimental endoscope 35,166
 - stress calculations for picture tubes 37, 56
 - solid-state image sensor 37,303
 - real-time orthogonal transformation of colour-TV pictures 38,119
 - Dot Scan CCTV for real-time image processing 38,310
 - correction plate for projection TV 39, 15
 - deflection coils of 30AX colour-picture system 39,154
 - new concept for camera tubes 39,201
 - 12 GHz receivers for satellite TV 39,257
 - measuring magnetic fields of TV deflection coils 39,277
 - games with computer * 40,230
- Temperature distribution in gas-filled incandescent lamps 32,206
- Tetrodes, MOS { 31,259
31,266
- Thermal insulation, high, by double glazing 34,242
- Thermography, macro- and micro- 37,241
- Thermogravimetric analysis of ferrites 31, 24
- Thick-film technology 35,144
- Thin films and coatings:**
 - Si, MOS transistors in 31,271
 - evaporated, fast refractometer for 35,142
 - strain gauges 39, 94
 - cryopump for magnetron sputtering 39,252
 - cryopump for evaporation 39,254
 - CVD for applying wear-resistant coatings to tool steel 40,204
- Thin-layer cells for electrochemiluminescence 40, 69
- Time-division multiplex (TDM) with pulse-amplitude modulation 36,330
- Tin-halide lamps 35,347
- Titanium carbide and nitride, wear-resistant coatings on tool steel 40,204
- Titration, complexometric, of GaSe 34,353
- TITUS tube:**
 - single, in colour data display 34,129
 - Phototitus optical converter 34,274
 - DKDP crystals for 39, 50
- Tomosynthesis, flashing 38,338
- Tone generation, digital, for musical instruments 31,354
- Tool steel, application of wear-resistant coatings to 40,204
- Traffic control:**
 - traffic-flow analysis by radar 31, 17
 - radio system for vehicle location (PINPOINT) 35, 15
 - actuator for anti-lock braking system 36, 74
 - see also **Airfields and Road lighting**
- Transducers:**
 - with miniature strain gauges in Si diaphragm 33, 14
 - for medical echography 38,195
 - with thin-film strain gauges 39, 94
 - electromechanical, with no hysteresis 40,358
- Transformation methods for improving parallel programs 40,278
- Transistors:**
 - bipolar, LOCOS technology for * 31,276
 - see also **Field-effect transistors and MOS transistors**
- Transition metals, alloying behaviour of { 33,149
36,217
- Translation operations in PHLIQA 1 38,269
- Transmission:**
 - of speech by code modulation with digitally controlled companding 31,335
 - of simple pictures 32, 42
 - see also **Data transmission**
- Transmitting-valve grids of pyrolytic graphite 37,194
- Transport reactions, chemical, in incandescent lamps { 35,296
35,302
35,316
35,322
36,133
- TRAPATT oscillator 40, 99
- Travelling-wave divider for UHF 38, 54
- Triglycine sulphate (TGS) for infrared detectors 35,247
- TTC (Thermal-Transient Compensated) cut for quartz crystals 40, 1
- Tungsten transport by water cycle 36,133
- Tuning of musical instruments 31,354
- Tunnelling in metal-semiconductor contacts under pressure 32,211
- Turn-to-turn diffusion in incandescent lamps 35,296
- Undercutting in wet-chemical etching 38,149
- Vacancy clusters in dislocation-free Si and Ge 34,244
- Vacuum technology:**
 - flow of highly rarefied gases 33, 43
 - industrial, cryopumps in 39,246

- Vehicle, demonstration and test, for exterior lighting* 35,307
- Vehicular propulsion, prospects of Stirling engine for 31,169
- Vibrating-reed magnetometer for microscopic particles 31, 40
- Vibrators, electromagnetic 33,249
- Video cassette recorder:**
magnetic recording in 36,326
grease-lubricated spiral-groove bearings for 39,184
- Video discs, see **LaserVision**
- Video tape, making tracks visible 40,129
- Video telephony:**
experimental system 36, 85
digital circuits in video telephone 36,233
signal switching 36,291
- Visual conspicuity, measurement of 36, 71
- VLP, see **LaserVision**
- VLSI:**
negative electron resists 39,346
for Compact Disc* 40,173
- Vortices 32,247
-
- Washing machine, speed-controlled d.c. motor for 34,163
- Water cycle, transport of tungsten by 36,133
- Water pollution, see **Environmental science**
- Water supply, control of treatment plant for 36,273
- Waveguides, cross-coupled 32,165
- Waves, special issue, with list of contents on p. 310 { 33,309 -349
- Wave soldering electronic components { 38,135 40,342
- Wear-resistant coatings, deposition on tool steel 40,204
- Welding, plasma-MIG 33, 21
- Window for ultrasoft X-rays from space 40, 12
- Windows, double-glazed, with good thermal insulation 34,242
- Wire:**
measurement of variation in diameter 31,111
Filament, boron 35,125
- Workshop technology:**
40 years of 31,127
engineering technology at IKO 39,315
- Writing, muscular activity in 32, 73
-
- X-rays:**
noise in X-ray films 31,117
PW 1100 single-crystal diffractometer 33, 61
fluorescence analysis { 34,339 34,353
automatic measurement of medical X-ray photographs 35,170
hard, channel-plate image intensifier for 37,124
crystal-structure research at high pressures with PW 1100 38,246
- digital image enhancement 38,298
flashing tomosynthesis 38,338
computerized mammogram processing 38,347
electroradiography for medical X-rays diagnosis 39, 19
window for ultrasoft X-rays from space 40, 12
TOMOSCAN 310* 40,253
-
- Young Scientists and Inventors, 10th European Philips Contest for:**
special issue { 38, 1 - 39
inquiry 38, 2
pollution of the Shannon in Limerick 38, 6
lichens 38, 8
hovercraft sprayer 38, 11
local functional equations 38, 13
instrument for measuring cave cross-section 38, 14
lecture 'Science and Education' 38, 17
lecture 'Science and Society' 38, 25
lecture 'Science and World Problems' 38, 33
-
- Zinc-diffusion profiles in GaP and GaAsP 37,121

Author index, Volumes 31-40

Figures in bold type indicate the volume number, and those in ordinary type the page number. Articles published in volumes 1-30 are given in the author indexes at the end of volumes 10, 15, 20, 25, 30 and 35.

- Aagaard, E. A.**, P. M. van den Avoort and F. W. de Vrijer
An experimental video-telephone system **36**, 85
- Aalders, J. W. G.**, R. J. van Duinen and P. R. Wesselius
The Groningen ultraviolet experiment with the Netherlands astronomical satellite (ANS) **34**, 33
- Acket, G. A.**, J. J. Daniele, W. Nijman, R. P. Tijburg and P. J. de Waard
Semiconductor lasers for optical communication **36**,190
- , R. Tijburg and P. J. de Waard
The Gunn diode **32**,370
- Admiraal, D. J. H.**, B. L. Cardozo, G. Domburg and J. J. M. Neelen
Annoyance due to modulation noise and drop-outs in magnetic sound recording **37**, 29
- Adriaansz, M.**, see Vriens, L.
- Aitchison, C. S.**
Lumped components for microwave frequencies **32**,305
- Akkerman, H. J.**
Engineering technology (*35 years of IKO*) **39**,315
- Albrecht, C.** and J. Proper
A method of determining noise in X-ray films **31**,117
- Alcock, R. N.**, D. A. Lucas and R. P. Vincent
MADGE, a microwave aircraft digital guidance equipment, I. General principles and angle-measuring units **34**,225
- Aldefeld, B.**
Calculation and design of electro-mechanical devices **39**, 78
- Alphen, M. P. van**, R. E. J. van de Grift, J. M. Pieper and R. J. van de Plassche
The PM 2517 automatic digital multimeter **38**,181
- Appels, J. A.**, H. Kalter and E. Kooi
Some problems of MOS technology **31**,225
- , see Nielen, J. A. van
- Arink, G. J. A.**
The onboard computer of the Netherlands astronomical satellite (ANS) **34**, 1
- Ass, H. M. J. M. van**, P. Geittner, R. G. Gossink, D. Kùppers and P. J. W. Severin
The manufacture of glass fibres for optical communication **36**,182
- Asselman, G. A. A.** and D. B. Green
Heat pipes,
I. Operation and characteristics **33**,104
II. Applications **33**,138
- and A. J. van Mensvoort
A refrigerator-freezer with heat pipe **40**,350
- Auphan, M.** and G. Dale
A transducer for medical echography **38**,195
- Avoort, P. M. van den**, see Aagaard, E. A.
- Bacchi, H.** and A. Moreau
Real-time orthogonal transformation of colour-television pictures **38**,119
- Baig, W. G.**
An edge-enhancing double-focus camera for image processing **37**,180
- Bakker, P.**
Linear electron accelerators (*35 years of IKO*) **39**,325
- Barth, P. J.**, see Voorman, J. O.
- Bastings, L. C.**, see Bruninx, E.
- Baudet, P.**, M. Binet and D. Boccon-Gibod
Low-noise microwave GaAs field-effect transistor **39**,269
- Beasley, J. P.** and D. G. Squire
Electron-beam pattern generator **37**,334
- Beek, L. K. H. van**
The PD photographic process **33**, 1
- Beekmans, N. M.** and L. Heyne
An instrument for monitoring low oxygen pressures **31**,112
- Beenakker, C. I. M.**, P. W. J. M. Boumans and P. J. Rommers
A microwave-induced plasma as an excitation source for atomic emission spectrometry **39**, 65
- Beer, A. F.**
A MOS transistor store with discretionary wiring **31**,286
- Behr, J.-P.**, P. Pernards, B. Schendel and J. Schwandt
Modelling and simulation as an aid in designing a computer **39**,134
- Belevitch, V.**
The lateral skin effect in a flat conductor **32**,221
- Belouet, C.**
DKDP crystals for use in the TITUS tube **39**, 50
- Berg, J. F. H. van de**, T. E. G. Daenen, G. Krijl and R. E. van de Leest
The electrodeposition of aluminium **39**, 87
- Berz, F.**, see Murphy, N. St. J.
- Bethe, K.** and D. Schön
Thin-film strain-gauge transducers **39**, 94
- Beun, M.**
A flexible method for automatic reading of handwritten numerals,
I. General description of the recognition method **33**, 89
II. The thinning procedure and determination of the special points **33**,130
- Beuvens, H. J. H.** and J. H. Dettingmeyer
A non-destructive measurement of pressures in incandescent lamps **31**, 93
- Bhargava, R. N.**
The physics of radiative centres in GaP **32**,261
- Binet, M.**, see Baudet, P.
- Blasse, G.** and A. Brill
Characteristic luminescence,
I. The absorption and emission spectra of some important activators **31**,304
II. The efficiency of phosphors excited in the activator **31**,314
III. Energy transfer and efficiency **31**,324
- , see Brill, A.
- Bleeker, J. A. M.**, W. H. Diemer, A. P. Huben and H. Hui-zenga
Camera window for ultrasoft X-rays from celestial sources **40**, 12
- Bleekrode, R.**, M. Koedam and L. Rehder
Discharge lamps **35**,308
- Bloem, H.**, J. C. de Grijs and R. L. C. de Vaan
An evacuated tubular solar collector incorporating a heat pipe **40**,181
- Bloem, J.**, A. Bouwknecht and G. A. Wesseling
Amalgams for fluorescent lamps **38**, 83
- Bloemendal, W.** and C. Kramer
The Netherlands astronomical satellite (ANS) **33**,117
- Blom, D.**, H. W. Hanneman and J. O. Voorman
Some inductorless filters **33**,294
- Blume, P.**
Computer-aided design **36**,162
- Boccon-Gibod, D.**, see Baudet, P.
- Boer, F. J. de**, see Boumans, P. W. J. M.
- Boers, P. M.** and L. J. M. Bollen
A fast scanning microscope used in the development of avalanche photodiodes **35**, 23
- Bollen, L. J. M.**, J. J. Goedbloed and E. T. J. M. Smeets
The avalanche photodiode **36**,205
- , see Boers, P. M.
- Bootsma, G. A.**, see Meijer, F.
- Bos, J. G. G.**
Spiral-groove bearing systems with grease **35**,137
- Bosch, G.** and J. H. H. Janssen
Integrated circuit with Hall device for brushless d.c. motors **31**,366
- Bosma, H.** and W. G. Gelling
LSI — a revolution in electronics **37**,267
- Boucher, A.** and B. C. Easton
Epitaxial growth of gallium arsenide **32**,380

- Boudewijns, H. P. J.**, E. C. Dijkmans, P. W. Millenaar, N. A. M. Verhoeckx and C. H. J. Vos
Digital circuits in the video telephone 36,233
- Boulou, M.**, M. Furtado and G. Jacob
Light-emitting diodes based on GaN 37,237
—, see Schiller, C.
- Boumans, P. W. J. M.**
Multielement analysis by optical emission spectrometry
— rise or fall of an empire? 34,305
—, F. J. de Boer and J. W. de Ruiter
A stabilized r.f. argon-plasma torch for emission spectroscopy 33, 50
—, see Beenakker, C. I. M.
- Boutot, J. P.** and J. C. Delmotte
Two microchannel-plate photomultipliers with subnanosecond characteristics 38,240
- Bouwer, A. G.**, G. Bouwhuis, H. F. van Heek and S. Wittekoek
The Silicon Repeater 37,330
—, R. H. Bruel, H. F. van Heek, F. T. Klostermann and J. J. 't Mannetje
The Ophycograph 34,257
- Bouwhuis, G.** and P. Burgstede
The optical scanning system of the Philips 'VLP' record player 33,186
—, see Bouwer, A. G.
- Bouwkamp, C. J.**
Scattering characteristics of a cross-junction of oversized waveguides 32,165
- Bouwknegt, A.**, H. Nienhuis, D. J. Schipper and P. A. W. Tielmans
Electrodes in discharge lamps 35,356
—, see Bloem, J.
- Bowers, Brian**
The early history of the electric motor 35, 77
- Brandt, B. B. M.**, W. Steinmaier and A. J. Strachan
LOCMOS, a new technology for complementary MOS circuits 34, 19
- Breed, D. J.**, F. H. de Leeuw, W. T. Stacy and A. B. Voermans
Garnet films for fast magnetic bubbles 38,211
- Brehm, R.**, K. van Dun, J. Haisma and J. C. G. Teunissen
Transparent single-point turning of glass 39, 92
- Brice, J. C.**, M. J. Hight, O. F. Hill and P. A. C. Whiffin
Pulling large bismuth-silicon-oxide crystals 37,250
— and W. S. Metcalf
Quartz-crystal resonators using an unconventional cut 40, 1
- Bril, A.**, G. Blasse, A. H. Gomes de Mesquita and J. A. de Poorter
Fast phosphors for colour television 32,125
—, see Blasse, G.
—, see Sommerdijk, J. L.
- Brinkman, A. C.**, J. Heise and C. de Jager
Observation of cosmic X-ray sources with the Netherlands astronomical satellite (ANS) 34, 43
- Brinkman, G. A.**
Radioisotopes (*35 years of IKO*) 39,294
- Broek, C. A. M. van den** and A. L. Stuijts
Ferroxdure 37,157
- Broese van Groenou, A.** and J. D. B. Veldkamp
Grinding brittle materials 38,105
- Brongersma, H. H.**, F. Meijer and H. W. Werner
Surface analysis, methods of studying the outer atomic layers of solids 34,357
- Bronnes, R. L.**, R. C. Hughes and R. C. Sweet
Ceramic-to-metal bonding with sputtering as a metallization technique 35,209
- Brouha, M.** and A. G. Rijnbeek
Apparatus for solid-state research at very high pressures 36,245
- Brouwer, G.**, see Witmer, A. W.
- Brouwer, H. J.**, S. M. de Veer and H. Zeedijk
The SO₂ monitoring network in the Rhine estuary region 32, 33
- Bruel, R. H.**, see Bouwer, A. G.
- Bruninx, E.** and L. C. Bastings
Stoichiometric analysis of gallium selenide 34,350
- Bunge, E.**
Speaker recognition by computer 37,207
- Burgstede, P.**, see Bouwhuis, G.
- Burnett, D. J.**
INDA, a software tool for the production engineer 40,237
- Bussche, W. van den**, A. H. Hoogendijk and J. H. Wessels
Signal processing in the Philips 'VLP' system 33,181
- Butzelaar, P. F.** and L. P. J. Hoogeveen
A new method of measuring the oxygen demand of water 34,123
- Carasso, M. G.**, J. B. H. Peek and J. P. Sinjou
The Compact Disc Digital Audio System 40,151
- Cardozo, B. L.**, see Admiraal, D. J. H.
- Carl, K.**, J. A. M. Dikhoff and W. Eckenbach
The pigmentation of phosphors for colour television 40, 48
- Casimir, H. B. G.**
Theoretical physics and industrial research 32,149
Gilles Holst, pioneer of industrial research in the Netherlands 40,121
- Chalmerton, V.**
A channel-plate image intensifier for hard X-rays 37,124
- Christiaens, M.**, see Jager, F. de
- Christis, W. J.**
The optical sensors of the Netherlands astronomical satellite (ANS), III. The star sensor 34,218
- Clegg, J. B.** and E. J. Millett
The determination of carbon, oxygen and nitrogen in semiconductors by spark-source mass spectrography 34,344
- Coenders, J. W.**
Switching telephone and video-telephone signals 36,291
- Compaan, K.** and P. Kramer
The Philips 'VLP' system 33,178
- Corbett, B. D.**
MADGE, a microwave aircraft digital guidance equipment, II. The data link: data transmission and distance measurement 35,271
- Crucq, J.**
The reaction wheels of the Netherlands satellite ANS 34,106
- Daenen, T. E. G.**, see Berg, J. F. M. van de
- Dale, G.**, see Auphan, M.
- Daniele, J. J.**, see Acket, G. A.
- Daniels, A.** and F. K. du Pré
Miniature refrigerators for electronic devices 32, 49
- Daniels, J.**, K. H. Härdtl and R. Wernicke
The PTC effect of barium titanate 38, 73
- Dantzig, R. van**
BOL (*35 years of IKO*) 39,302
- Davies, R.**, B. H. Newton and J. G. Summers
The TRAPATT oscillator 40, 99
— and R. E. Pearson
Parametric amplifiers for a radio-astronomy interferometer 32, 20
- Davis, G. L.**
Transport of tungsten by the water cycle 36,133
- Day, P. E.**, see Janssen, P. J. M.
- Dekkers, N. H.** and H. de Lang
A detection method for producing phase and amplitude images simultaneously in a scanning transmission electron microscope 37, 1
- Delmotte, J. C.**, see Boutot, J. P.
- Denner, W.** and Heinz Schulz
Apparatus based on Philips PW 1100 diffractometer for crystal-structure research at high pressures 38,246
- Dettingmeijer, J. H.**, G. Dittmer, A. Klopfer and J. Schröder
Research on incandescent lamps, III. Regenerative chemical cycles in tungsten-halogen lamps 35,302
—, see Beuvens, H. J. H.
- Diemer, W. H.**, see Bleeker, J. A. M.
- Dijk, P. van**
The optical sensors of the Netherlands astronomical satellite (ANS), II. The horizon sensor 34,213
- Dijken, R. H.**
Designing a small d.c. motor 35, 96
- Dijkmans, E. C.**, see Boudewijns, H. P. J.
- Dikhoff, J. A. M.**, see Carl, K.
- Dimigen, H.** and H. Lüthje
An investigation of ion etching 35,199
- Dinklo, J. A.** and E. B. de Vries
The microprogram control of the Philips P1000 family of computers 34,132
- Dittmer, G.**, see Dettingmeijer, J. H.
- Dixon, Bernard**
Science and Education 38, 17
- Dolizy, P.**
Growth of alkali-antimonide films for photocathodes 40, 19
- Dollekamp, H.**, L. J. M. Esser en H. de Jong
P²CCD in 60 MHz oscilloscope with digital image storage 40, 55

- Dolphin, R. J.**, L. P. J. Hoogeveen and F. W. Willmott
An experimental system for the automatic determination
of organochlorine residues in milk 36,284
- Domburg, G.**, see Admiraal, D. J. H.
- Donjon, J.** and G. Marie
Polychrome data display using a single TITUS tube . . . 34,129
- Doorn, R. A. van** and N. A. M. Verhoeckx
An I²L digital modulation stage for data transmission . 37,291
- Döring, M.**
Ink-jet printing 40,192
- Dorleijn, J. W. F.**, see Miedema, A. R.
- Dötsch, H.**
The microwave generation and manipulation of magnetic
domains 37, 38
- Drift, A. van der**, W. G. Gelling and A. Rademakers
Integrated circuits with leads on flexible tape. 34, 85
- Drop, P. C.**, E. Fischer, F. Oostvogels and G. A. Wesselink
Metal-halide discharge lamps 35,347
- Druyvesteyn, W. F.**, F. A. Kuijpers, A. G. H. Verhulst and
C. H. M. Witmer
Single-mask bubble memory with rotating-field control . 36,149
- Duinen, R. J. van**, see Aalders, J. W. G.
- Dumont, F.**, J. P. Hazan and D. Rossier
The Phototitus optical converter 34,274
- Dun, K. van**, see Brehm, R.
- Easton, B. C.**, see Boucher, A.
- Eckenbach, W.**, see Carl, K.
- Eggermont, L. D. J.**, M. H. H. Höfelt and R. H. W. Salters
A delta-modulation to PCM converter 37,313
- Elst, J. H. R. M.** and D. K. Wielenga
The finite-element method and the ASKA program, ap-
plied in stress calculations for television picture tubes 37, 56
- Engel, F. L.**
The measurement of visual conspicuity 36, 71
- Engelen, G. A. J. van**, J. L. M. Hagen and W. A. L. Heijne-
mans
An equipment for measuring the magnetic fields of tele-
vision deflection coils 39,277
- Engelsen, D. den**, J. H. Th. Hengst and E. P. Honig
An automated Langmuir trough for building monomole-
cular layers 36, 44
- Engelsma, G.**
Phenol synthesis and photomorphogenesis 38, 89
- Enz, U.** and R. W. Teale
Photomagnetic effects 31, 33
- Esser, L. J. M.**, see Dollekamp, H.
- Essers, W. G.**, G. Jelmorini and G. W. Tichelaar
Plasma-MIG welding. 33, 21
- Finck, J. C. J.**, H. J. M. van der Laak and J. T. Schrama
A semiconductor laser for information read-out 39, 37
- Fischer, E.**, J. Fitzgerald, W. Lechner and W. Lems
Research on incandescent lamps, II. Transport and burn-
out in incandescent lamps 35,296
- , see Drop, P. C.
- Fischer, W. E.**
A data-base management system for CAD and CAM . . . 40,245
- Fitzgerald, J.** and H. Hörster
The temperature distribution in gas-filled incandescent
lamps 32,206
- , see Fischer, E.
- Flinn, I.**, see Murphy, N. St. J.
- Foederer, A. F.**, J. L. M. Hagen and A. G. van Nie
An instrument for measuring the curvature of reflecting
surfaces 40,338
- Franken, A. J. J.**, G. D. Khoe, J. Renkens and C. J. G.
Verwer
Experimental semi-automatic machine for hot splicing
glass fibres for optical communication 38,158
- Frankort, F. J. M.**
Vibration patterns and radiation behaviour of loud-
speaker cones 36, 1
- Fransen, J. J. B.** and J. H. N. van Vucht
An easily controlled alkali-vapour dispenser 36, 16
- Franssen, N. V.** and C. J. van der Peet
Digital tone generation for a transposing keyboard in-
strument 31,354
- French, R. C.** and P. J. Mabej
Error control in mobile-radio data communication . . . 39,172
- Frens, G.**, H. F. Huisman, J. K. Vondeling and K. M. van
der Waarde
Suspension technology 36,264
- Fuller, K. L.**
AVOID, a short range high-definition radar 32, 13
— and A. J. Lambell
Traffic-flow analysis by radar 31, 17
- Funk, W.**
Thick-film technology 35,144
- Furtado, M.**, see Boulou, M.
- Geittner, P.**, see Ass, H. M. J. M. van
- Gelling, W. G.**, see Bosma, H.
- , see Drift, A. van der
- Gerwen, P. J. van**, W. A. M. Snijders and N. A. M. Ver-
hoeckx
An integrated echo canceller for baseband data transmis-
sion 39,102
- Gestel, W. J. van**, F. W. Gorter and K. E. Kuijk
Read-out of a magnetic tape by the magnetoresistance
effect 37, 42
- Gibson, R. W.**
PINPOINT — a radio system for locating and monitor-
ing vehicles 35, 15
- Gieles, A. C. M.** and G. H. J. Somers
Miniature pressure transducers with a silicon diaphragm 33, 14
- Gielis, G. C. M.**, J. B. H. Peek and J. M. Schmidt
Station and programme identification in FM sound broad-
casting 39,216
- Gijsbers, T. G.**
COLATH, a numerically controlled lathe for very high
precision 39,229
- Goddard, N. E.**
Solid-state microwave electronics 32,282
- Goedbloed, J. J.**, see Bollen, L. J. M.
- Goedhart, D.**, R. J. van de Plassche and E. F. Stikvoort
Digital-to-analog conversion in playing a Compact Disc 40,174
- Gomes de Mesquita, A. H.**, see Bril, A.
- Gool, G. H. van**, see Witmer, A. W.
- Gorter, F. W.**, see Gestel, W. J. van
- Gossel, D.**
The frequency-analog signal as a basis for measurement
and control 34,288
- Gossink, R. G.**, see Ass, H. M. J. M. van
- Greebe, C. A. A. J.**
Electromagnetic, elastic and electro-elastic waves . . . 33,309
- Greeffkes, J. A.** and K. Riemens
Code modulation with digitally controlled companding
for speech transmission 31,335
- Green, D. B.**, see Asselman, G. A. A.
- Grift, R. E. J. van de**, see Alphen, M. P. van
- Grijs, J. C. de**, see Bloem, H.
- Groh, G.**
The challenge of picture processing 38,291
- Groot, J. de** and A. Mircea
Computer calculations of the Gunn effect 32,385
- Groot, J. J. de**, J. A. J. M. van Vliet and J. H. Waszink
The high-pressure sodium lamp 35,334
- Guétin, P.** and G. Schröder
Tunnelling in metal-semiconductor contacts under pres-
sure 32,211
- Guildford, L. H.**
The Dot Scan CCTV, a flexible system for real-time image-
processing experiments 38,310
- Haas, L. A. de** and S. S. Wadman
The Waalre complex of Philips Research Laboratories . 31,153
- Haeringen, W. van**, see Lang, H. de
- Hagen, J. L. M.**, see Engelen, G. A. J. van
- , see Foederer, A. F.
- Haisma, J.**, see Brehm, R.
- Hanenberg, J. G. van den** and J. Vredenburg
An experimental assembly robot 40, 33
- Hanneman, H. W.**, see Blom, D.
- Hansen, P.**
Magnetic anisotropy and magnetostriction of Ru and Ir
ions in yttrium iron garnet 35,225
— and J.-P. Krumme
The compensation wall 34, 96
- Härdtl, K. H.**
A simplified method for the isostatic hot pressing of
ceramics 35, 65
—, see Daniels, J.
- Harrop, P.**, P. Lesartre and T. H. A. M. Vlek
Low-noise 12 GHz front-end designs for direct satellite
television reception 39,257

- Hart, C. M. and A. Slob**
Integrated injection logic (1^2L) 33, 76
- Hart, P. A. H. and F. M. Klaassen**
The MOS transistor as a small-signal amplifier 31,216
- Hart, J. 't, S. G. Nooteboom, L. L. M. Vogten and L. F. Willems**
Manipulation of speech sounds 40,134
- Havas, P. G.**
Measurement of wire-diameter variations 31,111
Sorting of quartz tubing for high-pressure mercury-vapour lamps 32, 57
- Haverkorn van Rijsewijk, H. C., P. E. J. Legierse and G. E. Thomas**
Manufacture of LaserVision video discs by a photopolymerization process 40,287
- Hazan, J.-P. and L. Jacomme**
Characterizing optical fibres; a test bench for pulse dispersion 36,211
—, see Dumont, F.
- Heek, H. F. van, see Bouwer, A. G.**
- Heemskerk, J. P. J. and K. A. Schouhamer Immink**
Compact Disc: system aspects and modulation 40,157
- Heide, H. van der**
Stabilization by oscillation 34, 61
- Heijnemans, W. A. L., J. A. M. Nieuwendijk and N. G. Vink**
The deflection coils of the 30AX colour-picture system . 39,154
—, see Engelen, G. A. J. van
- Heise, J., see Brinkman, A. C.**
- Heitmann, H., B. Hill, J.-P. Krumme and K. Witter**
MOPS, a magneto-optic storage wafer of the discrete-bit type 37,197
- Hengst, J. H. Th., see Engelsen, D. den**
- Heusden, S. van**
Air-pollution monitors based on chemiluminescence . . 34, 73
— and L. G. J. Mans
Measurement of ozone in an aircraft 38,131
- Heuven, J. H. C. van**
P-I-N switching diodes in phase-shifters for electronically scanned aerial arrays 32,405
— and A. G. van Nie
Microwave integrated circuits 32,292
- Heyne, L., see Beekmans, N. M.**
- Heyns, H., H. L. Peek and J. G. van Santen**
Image sensor with resistive electrodes 37,303
- Hight, M. J., see Brice, J. C.**
- Hill, B., see Heitmann, H.**
- Hill, O. F., see Brice, J. C.**
- Hily, C., J. J. Hunzinger, M. Jatteau and J. Ott**
Real-time macro- and microthermography 37,241
- Hoek, W. J. van den and W. A. Klessens**
Carbon-filament lamps with a chemical transport cycle . 35,316
- Hoeve, H., J. Timmermans and L. B. Vries**
Error correction and concealment in the Compact Disc system 40,166
- Höfelt, M. H. H., see Eggermont, L. D. J.**
- Hofker, W. K.**
Geiger-Müller counters (*35 years of IKO*) 39,296
Semiconductor detectors (*35 years of IKO*) 39,298
Electronics for nuclear physics (*35 years of IKO*) . . . 39,312
Ion-beam technology (*35 years of IKO*) 39,320
— and J. Politiek
Ion implantation in semiconductors 39, 1
- Hofmeester, J. H. M. and J. P. Koutstaal**
Moving-coil motors 33,244
- Holster, P. L., C. J. Th. Potters and H. F. G. Smulders**
A water-pressure operated control system for dish-washers 33, 29
- Honds, L. and K. H. Meyer**
A linear d.c. motor with permanent magnets 40,329
- Honig, E. P., see Engelsen, D. den**
- Hoogendijk, A. H., see Bussche, W. van den**
- Hoogeveen, L. P. J., see Butzelaar, P. F.**
—, see Dolphin, R. J.
- Hoppe, W. J. J. van, G. D. Khoe, G. Kuyt and H. F. G. Smulders**
Very smooth fracture faces for optical glass fibres . . 37, 89
- Horn, B. L. ten**
Forty years of workshop technology 31,127
- Hornstra, J. and H. Vossers**
The Philips PW 1100 single-crystal diffractometer . . . 33, 61
- Hörster, H., E. Kauer and W. Lechner**
The burn-out mechanism of incandescent lamps 32,155
—, see Fitzgerald, J.
- Howden, H.**
Production of optical correction plates for projection television 39, 15
- Hoyer, A. and M. Schindwein**
Digital image enhancement 38,298
— and W. Spiesberger
Computerized mammogram processing 38,347
- Huben, A. P., see Bleeker, J. A. M.**
- Hughes, R. C., see Bronnes, R. L.**
- Huisman, H. F., see Frens, G.**
- Huizenga, H., see Bleeker, J. A. M.**
- Hunzinger, J. J., see Hily, C.**
- Immink, K. A. Schouhamer, see Heemskerk, J. P. J.**
- Jacob, G., see Boulou, M.**
- Jacobs, C. A. J. and J. A. J. M. van Vliet**
A new generation of high-pressure sodium lamps 39,211
- Jacomme, L., see Hazan, J.-P.**
- Jager, C. de, see Brinkman, A. C.**
- Jager, F. de and M. Christiaens**
A fast automatic equalizer for data links 37, 10
- Jagt, J. C. and P. W. Whipps**
Negative electron resists for VLSI 39,346
- Jansen, J. A. J., see Witmer, A. W.**
- Janssen, J. H. H., see Bosch, G.**
- Janssen, P. J. M. and P. E. Day**
Control mechanisms in the Philips 'VLP' record player 33,190
- Jatteau, M., see Hily, C.**
- Jelmorini, G., see Essers, W. G.**
- Jeu, W. H. de and J. van der Veen**
Liquid crystals for numerical displays 37,131
- Jong, H. de, see Dollekamp, H.**
- Jonker, G. H. and A. L. Stuijts**
Controlling the properties of electroceramic materials through their microstructure 32, 79
- Joseph, R. D.**
MOS transistors for power amplification in the HF band 31,251
- Kalis, H. and J. Lemmrich**
Frequency-analog speed control 33,260
- Kalter, H. and E. P. G. T. van de Ven**
Plasma etching in IC technology 38,200
—, see Appels, J. A.
- Kamerbeek, E. M. H.**
Electric motors 33,215
Torque measurements on induction motors using Hall generators or measuring windings 34,153
Scaling laws for electric motors 35,116
- Kaps, G.**
'Controlled cascading', a new open-loop control principle for adjustable frequency dividers 32,103
- Kasperkovitz, W. D.**
Frequency-dividers for ultra-high frequencies 38, 54
- Kauer, E., see Hörster, H.**
- Kelly, J. J. and G. J. Koel**
Galvanic effects in the wet-chemical etching of metal films 38,149
- Kessel, Th. J. van and R. J. van de Plassche**
Integrated linear basic circuits 32, 1
- Kessels, J. L. W. and A. J. Martin**
Parallel programs 40,254
- Keve, E. T.**
Pyroelectric materials based on triglycine sulphate (TGS) for infrared detection 35,247
- Khoe, G. D. and L. J. Meuleman**
Light modulation and injection in optical-fibre transmission systems with semiconductor lasers 36,201
—, see Franken, A. J. J.
—, see Hoppe, W. J. J. van
- Kilian, R. and M. Liehr**
Experimental etching equipment 38, 51
- Klaassen, F. M.**
The development of field-effect transistor electronics . . 33,203
—, see Hart, P. A. H.
- Klein Wassink, R. J.**
The thermal behaviour of electronic components during soldering 38,135
— and H. J. Vledder
The attachment of leadless components to printed boards 40,342
- Klerk, M.**
The electron microprobe 34,370

- Klessens, W. A., see Hoek, W. J. van den
 Klinck, M.
 Control of the surface-water purification plant for the Amsterdam Water-Supply Authority 36,273
 Kloosterboer, J. G., G. J. M. Lippits and H. C. Meinders
 Photopolymerizable lacquers for LaserVision video discs 40,298
 Klopfer, A., see Dettingmeijer, J. H.
 Klostermann, F. T., see Bouwer, A. G.
 Klotz, E., R. Linde, U. Tiemens and H. Weiss
 Flashing tomosynthesis 38,338
 Knippenberg, W. F.
 Inorganic chemical analysis 34,298
 — and B. Lersmacher
 Carbon foam 36, 93
 —, B. Lersmacher and H. Lydtin
 Products of pyrolytic graphite 37,189
 Knowles, J. E.
 Measuring the magnetic properties of a microscopic particle 39, 48
 Kock, A. J. R. de
 Vacancy clusters in dislocation-free silicon and germanium 34,244
 Koedam, M., see Bleekrode, R.
 Koel, G. J., see Kelly, J. J.
 Kooi, E., see Appels, J. A.
 Koster, W. G., see Vredendregt, J.
 Köstlin, H.
 Double-glazed windows with very good thermal insulation 34,242
 —, see Schaper, H.
 Koutstaal, J. P., see Hofmeester, J. H. M.
 Kramer, C., see Bloemendal, W.
 Kramer, P., see Compaan, K.
 Kreuwels, W. G. J.
 Structural testing of digital circuits 35,261
 Krijl, G., see Berg, J. F. M. van de
 Kroon, D. J. and M. Q. Mengarelli
 Monitoring the quality of surface water 34,113
 Krul, L. G. and P. Reijnierse
 Transmission of simple pictures 32, 42
 Krumme, J.-P., see Hansen, P.
 —, see Heitmann, H.
 Kruseman Aretz, F. E. J.
 Abstraction 40,225
 Kuijk, K. E., see Gestel, W. J. van
 Kuijpers, F. A., see Druyvesteyn, W. F.
 Küppers, D., see Ass, H. M. J. M. van
 Kurz, H.
 Lithium niobate as a material for holographic information storage 37,109
 Kuus, G.
 A getter for metal-iodide high-pressure mercury-vapour lamps 35,354
 Kuypers, W. and J. C. Tiemeijer
 The Philips PSEM 500 scanning electron microscope 35,153
 Kuyt, G., see Hoppe, W. J. J. van
 Laak, H. J. M. van der, see Finck, J. C. J.
 Lakerveld, H. G.
 High-speed solid-rotor induction motors 34,170
 Lambell, A. J., see Fuller, K. L.
 Lang, H. de, D. Polder and W. van Haeringen
 Optical polarization effects in a gas laser 32,190
 —, see Dekkers, N. H.
 Lechner, W., see Fischer, E.
 —, see Hörster, H.
 Leest, R. E. van de, see Berg, J. F. M. van de
 Leeuw, F. H. de, see Breed, D. J.
 Legierse, P. E. J., see Haverkorn van Rijsewijk, H. C.
 Lely, P. van der and G. Missriegler
 Audio tape cassettes 31, 77
 Lemke, M. and W. Schilz
 Microwave integrated circuits on a ferrite substrate 32,315
 Lemmrich, J., see Kalis, H.
 Lems, W., see Fischer, E.
 Lens, G. A., see Zaengel, T.
 Lersmacher, B., see Knippenberg, W. F.
 Lesartre, P., see Harrop, P.
 Liehr, M., see Kilian, R.
 Linde, R., see Klotz, E.
 Lippits, G. J. M., see Kloosterboer, J. G.
 Lucas, D. A., see Alcock, R. N.
 Luijckx, G.
 The cyclotron (*35 years of IKO*) 39,290
 Lüthje, H., see Dimigen, H.
 Lydtin, H., see Knippenberg, W. F.
 Maaren, A. C. van, O. Schob and W. Westerveld
 Boron filament: a light, stiff and strong material 35,125
 Mabey, P. J., see French, R. C.
 Magarshack, J.
 Gunn-effect oscillators and amplifiers 32,397
 Mannetje, J. J. 't, see Bouwer, A. G.
 Mans, L. G. J., see Heusden, S. van
 Marie, G., see Donjon, J.
 Martiensen, W.
 Science and Society 38, 25
 Martin, A. J.
 Distributed computations on arrays of processors 40,270
 —, see Kessels, J. L. W.
 Meijer, A.
 An analogue computer for simulating one-dimensional aerial arrays 31, 2
 Meijer, F. and G. A. Bootsma
 Investigation of the chemical behaviour of clean silicon and germanium surfaces 32,131
 —, see Brongersma, H. H.
 Meijer, R. J.
 Prospects of the Stirling engine for vehicular propulsion 31,169
 Meinders, H. C., see Kloosterboer, J. G.
 Melis, J. H. A.
 O-BUS: a system for flexible public transport by means of on-call buses 40,231
 Memming, R.
 Energy production by photoelectrochemical processes 38,160
 Mengarelli, M. Q., see Kroon, D. J.
 Mensvoort, A. J. van, see Asselman, G. A. A.
 Metcalf, W. S., see Brice, J. C.
 Meuleman, L. J., see Khoe, G. D.
 Meyer, K. H., see Honds, L.
 Meyer, W. and W. Schilz
 Microwave measurement of moisture content in process materials 40,112
 Michel, C.
 Observations of domains in ferroelectrics and ferromagnetics with a scanning electron microscope 36, 18
 Miedema, A. R.
 A simple model for alloys,
 I. Rules for the alloying behaviour of transition metals 33,149
 II. The influence of ionicity on the stability and other physical properties of alloys 33,196
 The heat of formation of alloys 36,217
 The atom as a metallurgical building block 38,257
 — and J. W. F. Dorleijn
 Electrical conduction in ferromagnetic metals 35, 29
 Millenaar, P. W., see Boudewijns, H. P. J.
 Millett, E. J., see Clegg, J. B.
 Mircea, A., see Groot, J. de
 Missriegler, G., see Lely, P. van der
 Mitchell, R. F.
 Acoustic surface-wave filters 32,179
 Moerkens, J. C.
 An electronic starter for long fluorescent lamps 31, 54
 Moreau, A., see Bacchi, H.
 Moutaán, K.
 IMPATT-diode oscillators 32,345
 Optical communication systems with glass-fibre cables 36,178
 Muijderman, E. A., G. Remmers and L. P. M. Tielemans
 Grease-lubricated spiral-groove bearings 39,184
 Murphy, N. St. J., F. Berz and I. Flinn
 Carrier mobility in MOS transistors 31,237
 Neelen, J. J. M., see Admiraal, D. J. H.
 Newton, B. H., see Davies, R.
 Nicia, A. J. A. and C. J. T. Potters
 Components for glass-fibre circuits 40, 46
 Nie, A. G. van
 A method of measuring mechanical stresses in passivation layers 39,130
 —, see Foederer, A. F.
 —, see Heuven, J. H. C. van
 Nie, C. P. van
 Improved ear-lobe clip for physiological transducers 33,102

- Nielen, J. A. van**
Operation and d.c. behaviour of MOS transistors . . . 31,209
—, M. J. J. Theunissen and J. A. Appels
MOS transistors in thin monocrystalline silicon layers . . . 31,271
- Nienhuis, H.**, see Bouwknegt, A.
- Nienhuis, R. J.**
Integrated audio amplifiers with high input impedance and low noise . . . 31,245
A MOS tetrode for the UHF band with a channel 1.5 μm long . . . 31,259
- Niessen, C.**
Computer-aided design of LSI circuits . . . 37,278
- Nieuwendijk, J. A. M.**, see Heijnemans, W. A. L.
- Nieuwland, J. M. van**, A. Petterson and C. Weber
The design and construction of a non-rectangular reverberation chamber . . . 37,176
- Nijman, W.**, see Acket, G. A.
- Nobel, D. de** and M. T. Vlaardingerbroek
IMPATT diodes . . . 32,328
- Nooteboom, S. G.**, see Hart, J. 't
- O'Hanlon, H.**
Ferrite-cored kicker magnets . . . 31, 66
- Oostrom, A. van**
Field emission of electrons and ions . . . 33,277
- Oostvogels, F.**, see Drop, P. C.
- Opdorp, C. van**, see Werkhoven, C.
- Opstelten, J. J.**, D. Radielović and J. M. P. J. Versteegen
Optimum spectra for light sources . . . 35,361
- Ott, J.**, see Hily, C.
- Otterloo, P. van**
Attitude control for the Netherlands astronomical satellite (ANS) . . . 33,162
- Overgoor, B. J. M.**
An integrated chopper circuit with MOS transistors . . . 31,248
- Pannenberg, A. E.**
Science and World Problems . . . 38, 33
- Parker, D. W.**, R. G. Pratt, F. W. Smith and R. Stevens
Acoustic surface-wave bandpass filters . . . 36, 29
- Pearson, R. E.**, see Davies, R.
- Peek, H. L.**, see Heyns, H.
- Peek, J. B. H.**, see Carasso, M. G.
- , see Gielis, G. C. M.
- Peelen, J. G. J.**
Light transmission of sintered alumina . . . 36, 47
—, B. V. Rejda and J. P. W. Vermeiden
Sintered hydroxylapatite as a bioceramic . . . 37,234
- Peet, C. J. van der**, see Franssen, N. V.
- Périlhou, J.**
An experimental endoscope with miniature television camera . . . 35,166
- Pernards, P.**, see Behr, J.-P.
- Persoon, E. H. J.**
A system that can learn to recognize two-dimensional shapes . . . 38,356
- Peschmann, K. R.**
Medical electroradiography — its potential and limitations . . . 39, 19
- Petersen, A.**, P. Schnabel, H. Schweppe and R. Wernicke
A small analog memory based on ferroelectric hysteresis . . . 37, 51
- Petterson, A.**, see Nieuwland, J. M. van
- PHLIQA Project Group**
PHLIQA 1, a question-answering system for data-base consultation in natural English,
I. Organization and performance . . . 38,229
II. The artificial languages and translation operations . . . 38,269
- Pieper, J. M.**, see Alphen, M. P. van
- Pistorius, J. A.**, J. M. Robertson and W. T. Stacy
The perfection of garnet bubble materials . . . 35, 1
- Plassche, R. J. van de**, see Alphen, M. P. van
—, see Goedhart, D.
—, see Kessel, Th. J. van
- Polaert, R.** and J. Rodière
Investigation of microchannel plates by scanning electron microscopy . . . 34,270
- Polder, D.**, see Lang, H. de
- Pölitiek, J.**, see Hofker, W. K.
- Polman, J.**, H. van Tongeren and T. G. Verbeek
Low-pressure gas discharges . . . 35,321
- Poorter, J. A. de**, see Bril, A.
- Potters, C. J. Th.**, see Holster, P. L.
—, see Nicia, A. J. A.
- Pratt, R. G.**, see Parker, D. W.
- Pré, F. K. du**, see Daniels, A.
- Proper, J.**, see Albrecht, C.
- Rademakers, A.**, see Drift, A. van der
- Radielović, D.**, see Opstelten, J. J.
- Radziwill, W.**
Steady-state performance of a class of electronically commutated d.c. machines . . . 35,106
- Raes, R.** and J. Schellekens
A speed-controlled d.c. motor for a washing machine . . . 34,163
- Rathenau, G. W.**
Innovation in electronic devices . . . 32,117
- Rehder, L.**, see Bleekrode, R.
- Reijnen, P. J. L.**
Thermogravimetric analysis applied to ferrites . . . 31, 24
- Reijnierse, P.**, see Krul, L. G.
- Rejda, B. V.**, see Peelen, J. G. J.
- Remmers, G.**
Grease-lubricated helical-groove bearings of plastic . . . 34,103
—, see Muijderman, E. A.
- Renkens, J.**, see Franken, A. J. J.
- Rennicke, K.**
Supply-voltage speed control for capacitor motors . . . 34,180
- Riemens, K.**, see Greefkes, J. A.
- Rijckaert, A. M. A.**
Making the tracks on video tape visible with a magnetic fluid . . . 40,129
- Rijnbeek, A. G.**, see Brouha, M.
- Roberts, E. D.**
Electron resists for the manufacture of integrated circuits . . . 35, 41
- Robertson, J. M.**, see Pistorius, J. A.
- Rodière, J.**, see Polaert, R.
- Roeder, E.**
Extrusion of glass . . . 32, 96
- Rommers, P. J.**, see Beenakker, C. I. M.
- Roosmalen, J. H. T. van**
A new concept for television camera tubes . . . 39,201
- Röschmann, P.**
YIG filters . . . 32,322
- Rossier, D.**, see Dumont, F.
- Rothgordt, U.**
Electrostatic printing . . . 36, 57
- Ruiter, J. W. de**, see Boumans, P. W. J. M.
- Salter, R. H. W.**, see Eggermont, L. D. J.
- Sangster, F. L. J.**
The 'bucket-brigade delay line', a shift register for analogue signals . . . 31, 97
Integrated bucket-brigade delay line using MOS tetrodes . . . 31,266
- Santen, J. G. van**, see Heyns, H.
- Saraga, P.** and J. A. Weaver
An experiment in flexible automation . . . 38,329
- Scha, R. J. H.**
Software . . . 40,219
- Schaper, H.**, H. Köstlin and E. Schnedler
Electrochemiluminescence in electrolyte-free solutions . . . 40, 69
- Scheer, J. J.** and J. Visser
Application of cryopumps in industrial vacuum technology . . . 39,246
- Schellekens, J.**, see Raes, R.
- Schemmann, H.**
Stability of small single-phase synchronous motors . . . 33,235
- Schendel, B.**, see Behr, J.-P.
- Schiefer, G.**
Broadband circulators for VHF and UHF . . . 36,255
- Schiller, C.** and M. Boulou
Investigation of crystal defects by cathodoluminescence . . . 35,239
- Schilz, W.**, see Lemke, M.
—, see Meyer, W.
- Schipper, D. J.**, see Bouwknegt, A.
- Schindwein, M.**, see Hoyer, A.
- Schmidt, J. M.**, see Gielis, G. C. M.
- Schmidt, U. J.**
Electro-optic deflection of a laser beam . . . 36,117
- Schnabel, P.**, see Petersen, A.
- Schnedler, E.**, see Schaper, H.
- Schnell, A.**
Electromechanical transducers with no hysteresis . . . 40,358
- Schob, O.**, see Maaren, A. C. van
- Scholl, G. J.**
A universal instrument for digital picture processing . . . 38,326

- Schön, D., see Bethe, K.
 Schouhamer Immink, K. A., see Heemskerck, J. P. J.
 Schouten, J. F.
 The EVOLUON, a permanent Philips exhibition . . . 31,187
 Schrama, J. T., see Finck, J. C. J.
 Schröder, G., see Guétin, P.
 Schröder, J.
 Temperature profiling of tungsten filaments in incandescent lamps by a chemical transport reaction . . . 35,332
 —, see Dettingmeijer, J. H.
 Schulz, Heinz, see Denner, W.
 Schwandt, J., see Behr, J.-P.
 Schweppe, H., see Petersen, A.
 Scott, J. P.
 Electron-image projector . . . 37,347
 Severin, P. J. W., see Ass, H. M. J. M. van
 Shannon, J. M.
 Ion-implanted high-frequency MOS transistors . . . 31,267
 Sinjou, J. P., see Carasso, M. G.
 Sintzoff, M.
 Transformation methods for improving parallel programs 40,278
 Skoyles, D. R.
 A fast actuator for an anti-lock braking system . . . 36, 74
 Slob, A., see Hart, C. M.
 Sluijterman, L. A. Æ.
 Chromatofocusing, a new protein-separation method . . . 39,125
 Smeets, E. T. J. M., see Bollen, L. J. M.
 Smets, A. J.
 The optical sensors of the Netherlands astronomical satellite (ANS), I. The sun sensors . . . 34,208
 Smith, F. W., see Parker, D. W.
 Smulders, H. F. G., see Holster, P. L.
 —, see Hoppe, W. J. J. van
 Snijder, P. J., see Voorman, J. O.
 Snijders, W. A. M., see Gerwen, P. J. van
 Somers, G. H. J., see Gieles, A. C. M.
 Sommerdijk, J. L. and A. Bril
 Phosphors for the conversion of infrared radiation into visible light . . . 34, 24
 — and A. L. N. Stevels
 The behaviour of phosphors with aluminate host lattices 37,221
 Spiesberger, W. and M. Tasto
 The automatic measurement of medical X-ray photographs . . . 35,170
 —, see Hoyer, A.
 Squire, D. G., see Beasley, J. P.
 Staas, F. A.
 Continuous cooling in the millikelvin range . . . 36,104
 Stacy, W. T., see Breed, D. J.
 —, see Pistorius, J. A.
 Steen, L. M. van der
 Digital integrated circuits with MOS transistors . . . 31,277
 Steinmaier, W., see Brandt, B. B. M.
 Stevels, A. L. N., see Sommerdijk, J. L.
 Stevens, R., see Parker, D. W.
 Stikvoort, E. F., see Goedhart, D.
 Strachan, A. J., see Brandt, B. B. M.
 Straten, P. J. M. van der and G. Verspui
 Chemical vapour deposition of wear-resistant coatings on tool steel . . . 40,204
 Stuijts, A. L.
 Renaissance in ceramic technology . . . 31, 44
 —, see Broek, C. A. M. van den
 —, see Jonker, G. H.
 Summers, J. G., see Davies, R.
 Swansenburg, T. J. B.
 Self-organizing systems . . . 38,364
 Sweet, R. C., see Bronnes, R. L.
 Tasto, M., see Spiesberger, W.
 Teale, R. W., see Enz, U.
 Teunissen, J. C. G., see Brehm, R.
 Theunissen, M. J. J., see Nielen, J. A. van
 Thissen, F. L. A. M.
 An equipment for automatic optical inspection of connecting-lead patterns for integrated circuits . . . 37, 77
 Thomas, G. E., see Haverkorn van Rijsewijk, H. C.
 Tichelaar, G. W., see Essers, W. G.
 Tielemans, L. P. M., see Muijderman, E. A.
 Tielemans, P. A. W., see Bouwknegt, A.
 Tiemeijer, J. C., see Kuypers, W.
 Tiemens, U., see Klotz, E.
 Tijburg, R. P., see Acket, G. A.
 —, see Verplanke, J. C.
 Timmerman, J.
 Two electromagnetic vibrators . . . 33,249
 Timmermans, J., see Hoeve, H.
 Tongeren, H. van, see Polman, J.
 Tooren, A. van
 A simple and flexible automatic extractor . . . 35,196
 Troye, N. C. de
 Digital integrated circuits with low dissipation . . . 35,212
 Tummers, L. J.
 MOS transistors . . . 31,206
 Vaan, R. L. C. de, see Bloem, H.
 Veen, J. van der, see Jeu, W. H. de
 Veer, S. M. de, see Brouwer, H. J.
 Veldhuis, J.
 Computer-aided research on multiwire telephone cables 40, 85
 Veldkamp, J. D. B., see Broese van Groenou, A.
 Velzel, C. H. F.
 Holographic strain analysis . . . 35, 53
 Ven, E. P. G. T. van de, see Kalter, H.
 Venema, A.
 The flow of highly rarefied gases . . . 33, 43
 Verbeek, T. G., see Polman, J.
 Verbunt, J. P. M.
 'Droplet interferometry' for investigating smooth surfaces . . . 33, 74
 Verheijke, M. L.
 Neutron activation analysis . . . 34,330
 — and A. W. Witmer
 Line-intensity calculations for X-ray fluorescence analysis . . . 34,339
 Verhoeckx, N. A. M., see Boudewijns, H. P. J.
 —, see Doorn, R. A. van
 —, see Gerwen, P. J. van
 Verhulst, A. G. H., see Druyvesteyn, W. F.
 Vermeiden, J. P. W., see Peelen, J. G. J.
 Verplanke, J. C. and R. P. Tijburg
 Determination of zinc-diffusion profiles in gallium phosphide and gallium arsenophosphide with the aid of radioactive isotopes . . . 37,121
 Verspui, G., see Straten, P. J. M. van der
 Verstege, J. M. P. J., see Opstelten, J. J.
 Verweij, H.
 The fining of glass . . . 40,310
 Verwer, C. J. G., see Franken, A. J. J.
 Vincent, R. P., see Alcock, R. N.
 Vink, A. T., see Werkhoven, C.
 Vink, H. J.
 Crystal defects and their transport in solids; industrial applications . . . 35,181
 Vink, N. G., see Heijnemans, W. A. L.
 Visser, J., see Scheer, J. J.
 Vlaardingerbroek, M. T., see Nobel, D. de
 Vledder, H. J., see Klein Wassink, R. J.
 Vlek, T. H. A. M., see Harrop, P.
 Vliet, J. A. J. M. van, see Groot, J. J. de
 —, see Jacobs, C. A. J.
 Voermans, A. B., see Breed, D. J.
 Vogten, L. L. M., see Hart, J. 't
 Volger, J.
 Vortices . . . 32,247
 Cryogenics: a critical review . . . 37, 91
 Volman, H. J. W. M.
 The 'push-pull' spiral-groove bearing — a thrust bearing with self-adjusting internal preloading . . . 35, 11
 Vondeling, J. K., see Frens, G.
 Vonk, R.
 Magnetic deflection in television picture tubes . . . 32, 61
 Voorman, J. O., P. J. Snijder, J. S. Vromans and P. J. Barth
 An automatic equalizer for echo reduction in Teletext on a single chip . . . 40,319
 —, see Blom, D.
 Vos, C. H. J., see Boudewijns, H. P. J.
 Vos, J. A. de
 Megadoc, a modular system for electronic document handling . . . 39,329
 Vossers, H., see Hornstra, J.
 Vredenburg, J. and W. G. Koster
 Analysis and synthesis of handwriting . . . 32, 73
 —, see Hanenberg, J. G. van den

- Vriens, L. and M. Adriaansz
Investigation of gas discharges by light scattering . . . 35,344
- Vries, L. B., see Hoeve, H.
- Vries, E. B. de, see Dinklo, J. A.
- Vrijer, F. W. de
Modulation . . . 36,305
I. Modulation of a sinusoidal carrier . . . 36,309
II. Modulation of pulse trains . . . 36,329
III. Quantization and coding of analog signals . . . 36,337
IV. Transmission of digital signals . . . 36,343
V. Modulation in telecommunication . . . 36,353
—, see Aagaard, E. A.
- Vromans, J. S., see Voorman, J. O.
- Vucht, J. H. N. van
Intermetallic compounds; background and results of
twenty years of research . . . 36,136
—, see Fransen, J. J. B.
- Waalwijk, J. M. and N. Wiedenhof
The Institute for Nuclear Physics Research 'has finished
its work' (35 years of IKO) . . . 39,286
- Waard, P. J. de
Anomalous oscillations with an IMPATT diode . . . 32,361
—, see Acket, G. A.
- Waarde, K. M. van der, see Frens, G.
- Wadman, S. S., see Haas, L. A. de
- Wal, J. van der
A fast refractometer for evaporated thin films . . . 35,142
- Walther, K.
Magnetoacoustic effects in bismuth . . . 32,233
- Waszink, J. H., see Groot, J. J. de
- Waumans, B. L. A.
Software aspects of the PHIDIAS system . . . 40,262
- Weaver, J. A., see Saraga, P.
- Weber, C., see Nieuwland, J. M. van
- Weiss, H., see Klotz, E.
- Werkhoven, C., C. van Oudorp and A. T. Vink
Influence of crystal defects on the luminescence of GaP 38, 41
- Werner, H. W., see Brongersma, H. H.
- Wernicke, R., see Daniels, J.
—, see Petersen, A.
- Wesselink, G. A., see Bloem, J.
—, see Drop, P. C.
- Wesselius, P. R., see Aalders, J. W. G.
- Wessels, J. H., see Bussche, W. van den
- Westerveld, W., see Maaren, A. C. van
- Whiffin, P. A. C., see Brice, J. C.
- Whipps, P. W., see Jagt, J. C.
- Wiedenhof, N., see Waalwijk, J. M.
- Wielenga, D. K., see Elst, J. H. R. M.
- Wijers, J. L. C.
Three special applications of the Philips high-speed
spark-machining equipment . . . 40,199
- Willems, L. F., see Hart, J. 't
- Willmott, F. W., see Dolphin, R. J.
- Witmer, A. W., J. A. J. Jansen, G. H. van Gool and G.
Brouwer
A system for the automatic analysis of photographically
recorded emission spectra . . . 34,322
—, see Verheijke, M. L.
- Witmer, C. H. M., see Druyvesteyn, W. F.
- Wittekoek, S., see Bouwer, A. G.
- Witter, K., see Heitmann, H.
- Zaengel, T. and G. A. Lens
Fast pressure gauge for hot and corrosive gases . . . 39,344
- Zeedijk, H., see Brouwer, H. J.
- Zijlstra, H.
A vibrating-reed magnetometer for microscopic particles 31, 40
Materials research for permanent magnets . . . 34,193