

THE IMPACT OF LSI TECHNOLOGY ON COMPUTER SYSTEMS

Gerald B. HERZOG

 **IFIP CONGRESS 74**



NORTH-HOLLAND PUBLISHING COMPANY

THE IMPACT OF LSI TECHNOLOGY ON COMPUTER SYSTEMS

Gerald B. HERZOG

*RCA Laboratories
Princeton, New Jersey, USA*

(INVITED PAPER)

The long awaited universal logic array seems to have arrived. It is a complete microprogrammed computer on a chip. This has been made possible by high density arrays of MOS transistors. Advances in the fabrication techniques of MOS transistors promises to give bipolar speeds at MOS power densities. When these CPU's on a chip process data with stage delays of a few nanoseconds, the system designer must consider giving up designing unique hardware for new tasks and merely program a low cost, high speed microprocessor. Thus, new systems will be assemblies of microprocessors rather than of gates; just as present machines are built of integrated circuits rather than discrete components.

1. INTRODUCTION

It is hard to imagine an industry that has the same combination of high technology and rapid growth as does the semiconductor industry, though one cannot overlook the dynamic growth and sophistication of the computer business. In fact, the interdependent partnership of semiconductors and computers is probably unequalled in the history of technological development. Without the invention and development of the transistor, it is unlikely that the computer industry could have progressed to a fraction of its present size. Conversely, without the low cost computer power now available for artwork generation, and logic and circuit simulation, we would not be able to design in only a few months time, logic arrays containing thousands of transistors.

In just over a decade of time we have progressed from large, slow vacuum tube machines to actual computers that we can carry in our pocket, complete with a tape reader. I refer, of course, to the HP-65 that contains 50 000 transistors exclusive of memory, 30 000 bits of memory, and magnetic tape strips that hold 600 bits of program information. This has been made possible by batch fabrication techniques that have reduced the cost of a transistor by 3 orders of magnitude over the past ten years. Complex logic arrays containing several transistors per gate interconnected and coupled to adjoining gates now cost the user only a penny or less. Unfortunately, for this low price to apply, huge quantities of the same type of logic interconnection pattern must be needed by the user. A \$100 000 engineering effort may be required to develop a working array of a few thousand interconnected transistors despite extensive computer design automation aids. Therefore, if only one thousand pieces of such an array are needed, a nonrecurring cost of one hundred dollars per piece has been incurred while the production cost for this array may be less than ten dollars.

Ever since the first efforts on large scale integration (LSI) in the early sixties, there have been pleas by the semiconductor technologists for computer system architects to make computers more regular, so that a given array type could be used repeatedly in a large machine. There have been efforts in that direction, but progress seemed slow to the semiconductor technologist. Consequently, he turned his attention to an area that is inherently regular -- memory. Here the same array can be used over and over without limit. Furthermore, increased performance over core memories was easily achievable with semiconductors. It seemed ideal. A transistor flip-flop would be easy to address and drive, would give large output signals thereby reducing sense-amplifier problems, and would be less temperature

sensitive than cores. The extensive work on thin magnetic films in the early sixties seemed ample evidence that the computer manufacturer wanted faster and faster memories.

The first effort to make bipolar semiconductor memories produced some technical successes, but many business failures. The semiconductor specialist found that speed was desirable, but not nearly as important as low cost. Except for scratch pad memories, high speed bipolar memories were not cost effective. Packages containing 64 or 256 bits required too much labor on the part of the user to interconnect them. The package cost per bit of memory was too high, and the chip itself was expensive. The bipolar flip-flop occupied too much area, making the chip size large and the yield low. As a rule of thumb, cost goes up proportionally to the square of the chip area. This is due to two effects. First, cost goes up directly with chip area because there are fewer possible chips for a given wafer size, and the cost of processing a wafer is fixed. Costs also go up because yield of good chips is approximately inversely proportional to chip area. Thus, for main memories, chip size had to be reduced drastically to compete economically with cores.

This led to clever capacitance storage cell designs based on dynamic circuits that were being widely used in PMOS logic arrays. This produced a memory cell that was not only volatile in terms of loss of data when power was turned off but one that was constantly losing its stored data and had to be refreshed every millisecond or so. Obviously an impractical, foolish idea. No experienced computer manufacturer would use such a memory. Yet, the cost performance incentive proved strong, and over 20 billion bits of refresh memory will have been sold to independent computer manufacturers by August 1974. When IBM committed to such a memory technique it seemed it would soon be accepted by everyone, though today there is still hesitancy on the part of some. Part of this may be economic because core memory manufacturers continue to reduce prices and one still hears that "next year" semiconductor memories will be below cores in cost. To achieve this it will be necessary to produce chips with 4096 bits of memory which means approximately 5 000 or 15 000 transistors per chip, depending on whether a one transistor per bit or three transistor per bit memory cell is used.

The slower than hoped for acceptance of semiconductor memories led to further investigation of ways to make semiconductor logic arrays more regular so that they could be used repeatedly in large computers. Since there was little incentive for the

computer manufacturer to do this, the semiconductor manufacturer undertook the task. He quickly discovered the concept of microprogramming pioneered by Wilkes at Cambridge University many years ago. Here was a way to make logic more regular through the use of memory, and the designer quickly devised ways to make read-only memories (ROM). Now large, regular patterns of cells (2 000 to 8 000 bits) could be manufactured and easily tailored to a customer's specific needs by modifying only one mask out of a set of 4 or 5. Furthermore, such ROM's could be made a part of more complicated logic networks and tailored so that a single array could perform slightly different functions for different customers. Still, there was not a need for a large volume of such networks in computer main frames, or even in associated peripherals and terminals a few years ago. Consequently, the semiconductor industry created a new market. It produced the electronic calculator that now pervades every major retail store throughout the free world. I suspect there may have been more transistors produced for calculators in the past two years than were produced for computers in the previous ten years.

This experience with calculator design and the development of programmed logic arrays gave semiconductor manufacturers confidence in undertaking more ambitious projects. We now see the emergence of small processors complete on a single chip or small family of chips. The day of a computer on a chip seems to have arrived. It is bound to have a major impact on the digital control field and indirectly on the computer industry.

The manufacturing industry today cries out for more automation to increase efficiency, reduce labor content, and improve quality. Large computers are generally not cost effective in such applications. Even the minicomputer has become too expensive for many such applications.

The microprocessor, which is distinguished from the minicomputer by being a nearly complete CPU on only one or two chips of silicon, will be the low cost answer. Generally organized as an 8-bit machine these microprocessors are already finding widespread use in communications networks, remote intelligent terminals, point of sale terminals, process control including self diagnostics, apartment house heating and cooling systems, security systems, etc. They are invading amusement parlors where they control space war games, electronic tennis, and other interactive skill games. Soon they will be used to automatically calculate the player's score in bowling alleys in America. One experiment has shown that the gas mileage of a large car can be significantly increased by incorporating a microprocessor to control engine operation. Such an onboard computer could also provide antiskid braking, automatic speed control, car performance diagnostics, and other desired safety and convenience options. Its quite likely people will have computers in their car before they have them in their home.

The microprocessor will drastically change the area of logic design. Rather than design a unique logic network for a given control or data manipulation function and attempt to minimize the number of gates, the designer will write a program for a microprocessor. The goal will be to minimize the number of steps in his algorithm to conserve memory. Logic is free, but memory still costs money. The future logic designer will be more of a software man than a hardware man. The hardware design will have been done by the semiconductor manufacturer.

2. SEMICONDUCTOR DEVICE PROCESSING

The batch processing of silicon wafers by photolithography techniques perfected in the past ten years has made all this possible. Early transistors were made individually on a chip of germanium about one millimeter on a side. Manually placed dots of indium approximately 350 micrometers in diameter were alloyed into the germanium to form the emitter and collector of the first computer transistor. Today, photographic techniques define patterns in silicon with dimensions of only a few micrometers and create nearly 1 000 transistors in that same millimeter square chip. This photographic technique is used in the fabrication of both bipolar and MOS transistors. Bipolar transistors are more difficult to fabricate because the depth profile of doping impurities put into the silicon during the processing must be carefully controlled. The unipolar MOS transistor on the other hand, is only sensitive to the horizontal spacing between different doping regions, and hence subject primarily to the photolithography perfection.

Today the digital world is on the threshold of a major change. Bipolar transistors have dominated since the earliest days of the discrete transistor computer. Early integrated logic gates sought to emulate discrete component designs, hence the families of RTL, DCTL, and DTL. Then came the realization that integrated transistors were cheaper than integrated resistors, and TTL arrived on the scene.

The bipolar digital IC logic market is clearly dominated by TTL with 70 per cent of U.S. dollar volume in 1973. High speed ECL captured only 10 per cent. Because of its higher price, this is a much smaller per cent of the digital market in actual units for new applications. The remaining 20 per cent of the bipolar market is divided among many older types of logic. The surprising number is the 40 per cent of the total digital market captured by MOS devices. The use of MOS devices for complex digital arrays got off to a bad start. Overly ambitious programs were attempted in the early sixties that went beyond the semiconductor processor skill level. Failure to deliver working arrays created a serious credibility gap. Often the overly ambitious programs were attempted by small entrepreneurial companies with insufficient resources to solve their problems quickly enough. These small semiconductor houses sometimes folded with serious losses to the system companies that had been promised the moon. Understandably, this new technology got a bad reputation. By 1973 this learning period had become past history for most companies. All device types are now available from major semiconductor manufacturers, P-MOS, N-MOS, and C-MOS.

While MOS was initially thought to be good only for low performance equipments because it was believed that it was inherently slow, various studies have shown that the device is inherently fast but limited by its own parasitic capacitance [1]. This has led to efforts to reduce the parasitic capacitance, and MOS logic inverters that switch in 0.5 nanoseconds have been fabricated. It is quite likely that MOS arrays five years from now will be operating at data rates in excess of those economically possible with bipolar arrays. I say economically possible because the speed of bipolar transistors can always be made faster than MOS transistors in equivalent environments. Unfortunately, the costs are generally prohibitive.

A simplistic interpretation of how transistors operate will help us compare bipolar and MOS devices and the operation of different forms of logic gates. From this I hope to show that MOS devices will, in fact, be the most common form of logic in the future.

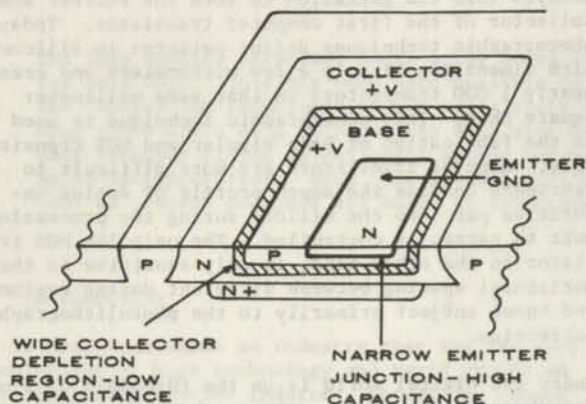


Fig. 1. Bipolar transistor.

The bipolar transistor of fig. 1 normally operates with a reverse bias on its collector and the emitter at some reference voltage, usually ac ground. If we think of the transistor as a charge control device, then a positive charge in the base region is necessary in an NPN transistor to change the field across the base-emitter junction capacitance, and cause current to flow from the emitter into the base region and across to the collector region. Since the base to emitter junction region is very narrow, the capacitance is very high and the bipolar transistor fundamentally has a low input impedance. If the transistor was made perfectly with perfect material of infinite carrier lifetime, the charge needed to turn the transistor on would be conserved, and the ratio of controlled current to controlling current, beta, would be infinite. In practice, some of the charge is lost and beta's range from in the tens to in the thousands. Typical beta values for switching transistors are under 100. Since the beta is finite, available output current of a transistor in a logic gate must be shared among the inputs of the gates it drives. To insure proper operation at even the slowest switching speeds, the number of logic loads on a gate must be restricted and the power supply voltages must be held within a given narrow range. Too low a voltage will reduce the fan-out or load driving capability. Too high a voltage will cause excess dissipation in the gate because of higher currents and the higher voltage.

High speed operation of a bipolar transistor is limited by the time it takes for carriers to get from the emitter to the collector and to charge the output load capacitance. The transit time from emitter to collector can be made very small in modern transistors. Since the inputs of the following gates are low impedance the circuit board wiring capacitance, etc. is not a major factor in determining gate switching speed. The saturation effect in a transistor can be a limiting factor, however. If during the switching operation, the collector voltage falls below the base voltage, so that the collector base junction is forward biased, then excess carriers appear in both the base and collector regions. The transistor cannot be completely turned off until all these carriers are removed. Such effects can be avoided by designing

circuits to limit the voltage swing in the collector, as in the case of the ECL gate. Alternatively, the base drive can be limited by feedback techniques as in the Schottky clamped forms of TTL. In both cases, the penalty is higher power dissipation. This power dissipation is the ultimate factor limiting the number of gates that can be integrated safely on a given chip.

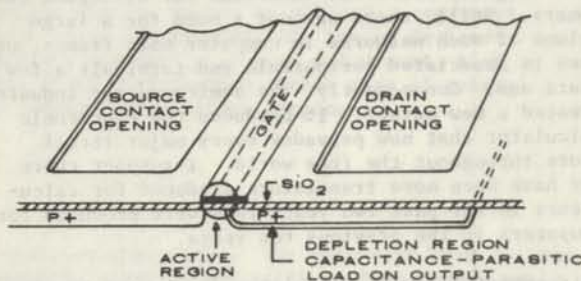


Fig. 2. MOS transistor.

The MOS transistor, fig. 2, is a charge control device where the controlling electrode is separated from the active region by a high quality dielectric, normally silicon dioxide. Because of the relatively large distance of the controlling element from the conducting region, compared to a bipolar, higher control voltages are needed and lower values of conductance are achieved for a given amount of controlling charge. Thus, to achieve conductance equal to that of a bipolar gate, larger devices occupying larger amounts of silicon area are necessary. However, since the real part of the input impedance of an MOS is nearly infinite, only modest amounts of conductance are necessary in an array environment to achieve large fan-outs.

Because the MOS is inherently a high impedance device, its speed is limited by its own output capacitance as well as the driven load capacitance. The transit time for the carriers to travel from the source to drain are comparable to those in a bipolar, i.e. less than one nanosecond. The MOS transistor is a good amplifier at hundreds of megahertz and in switching circuits it does not suffer from the storage effects of the bipolar. It is limited in speed only because of the way it is traditionally fabricated in silicon wafers. Although the only active region is the channel between the source and drain, the drain area and hence parasitic capacitance is determined by dimensions necessary to make electrical contact. In most instances, the parasitic capacitance is ten times higher than the capacitance associated with the active region. If all or most of the silicon could be removed, except in the channel region, and replaced by an insulator, then MOS integrated logic gates could approach the basic speed of the MOS device. This has been done and promises to be a way of making MOS devices that will replace TTL in the future.

As stated previously, MOS devices require relatively high voltages to turn them on compared to bipolars. The relationship between output current and control voltage is square law unlike the exponential relationship in bipolars. Therefore, a single MOS transistor makes a relatively poor logic gate. This was the cause of many of the problems with early P-MOS

logic arrays. Variations in processing effected the turn-on point such that the output of one array which was supposed to represent a "0" was regarded by another array as representing a "1" and usually left a device only partway turned on or off. Improved processing and better circuit design have virtually eliminated such problems in modern MOS logic networks.

One circuit approach which has distinct advantages is a complementary symmetry arrangement. Here two devices of opposite type are paired so that their combined switching characteristics are essentially the product of their individual characteristics. This gives a sharp transition between the "0" and "1" state resulting in good noise immunity and well defined "0" and "1" states. The complementary inverter, fig. 3a, approaches an ideal switch, fig. 3b.

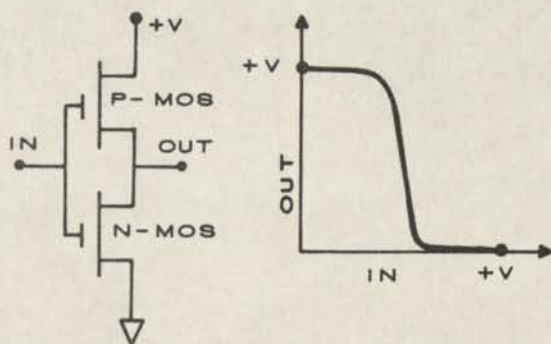


Fig. 3a. MOS complementary inverter.

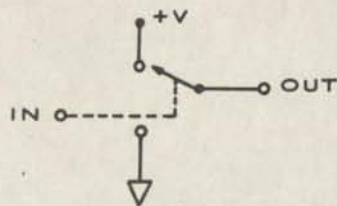


Fig. 3b. Ideal switch.

Since the "0" and "1" states are defined by ground and the supply voltage, +V, the logic gate works with a wide range of power supply voltages, typically from 3 to 15 volts. The logic decision point tracks at nearly the 50 per cent point of the supply voltage so that large changes in power supply voltage or ripple in the supply voltage do not cause logic errors. Furthermore, since no dc current flows into the MOS input, the only circuit losses are due to charging the output load capacitance. The circuit power is therefore proportional to CV^2f . Hence, at zero frequency or standby condition, the logic array dissipates essentially zero power. At long last there is an adaptive circuit that automatically adjusts its power requirements to the speed of the logic operation required.

When this form of logic is fabricated by processes that eliminate the undesired MOS parasitic capacitance, a superior logic family is achieved that will eventually replace TTL and even compete in many areas where ECL is now the only contender. It will compete with ECL on a systems basis because much more complexity can be achieved on a chip with this form of MOS circuitry compared to ECL. The economics of the situation will make it practical to use more gates to process in parallel and in pipeline fashion with MOS's but not with ECL.

At present, complementary symmetry logic gates are being fabricated in R&D activities that have propagation delays of 3 to 5 nanoseconds for gates with

fan-out of 3 and fan-in of 3. Inverters have only 2 nanoseconds delay and improvements are possible. These delays are achieved with power dissipation proportion to CV^2f which amounts to 5 milliwatts at 50 MHz for a 3 input and a 3 output logic gate. While this power approaches that of bipolar, the important fact to remember is that many logic gates do not switch on each logic cycle, hence, the system power is reduced in accordance to the activity of the logic gates. For example, a general purpose arithmetic unit with 560 transistors operating at a 50 MHz rate dissipates 100 milliwatts or about 1 milliwatt per gate.

The high speed is achieved by fabricating the complementary symmetry circuit devices in thin films of silicon deposited on sapphire substrates. In this way, no excess parasitic capacity is introduced into the structure of the MOS and its ultimate speed capability is approached. Details of this process are given in the references [2-6].

In summary, it seems clear that the era of the bipolar integrated circuit has reached its apogee and that large scale integration of superspeed MOS devices will be the technology of the future. Figure 4 indicates the area of performance that can be covered by silicon-on-sapphire C-MOS. The memory array of 256 bits shown in fig. 5 operates with less than 100 nanoseconds access time. It dissipates less than one microwatt per bit in standby and 0.2 milliwatts per bit when operating at a 5 MHz rate.

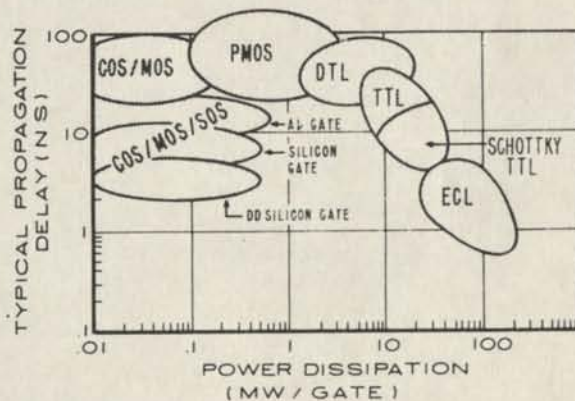


Fig. 4. Speed-power relationships for various logic forms.

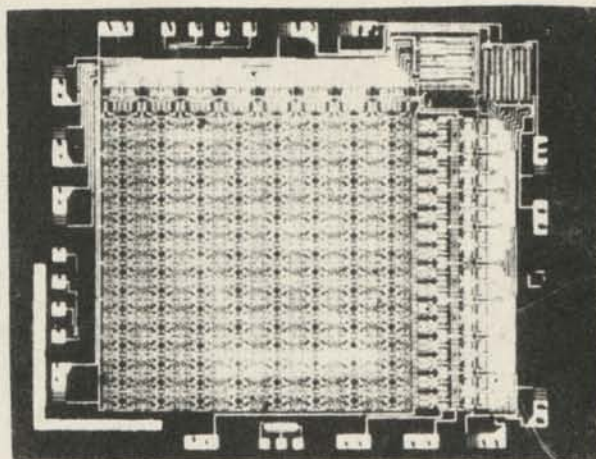
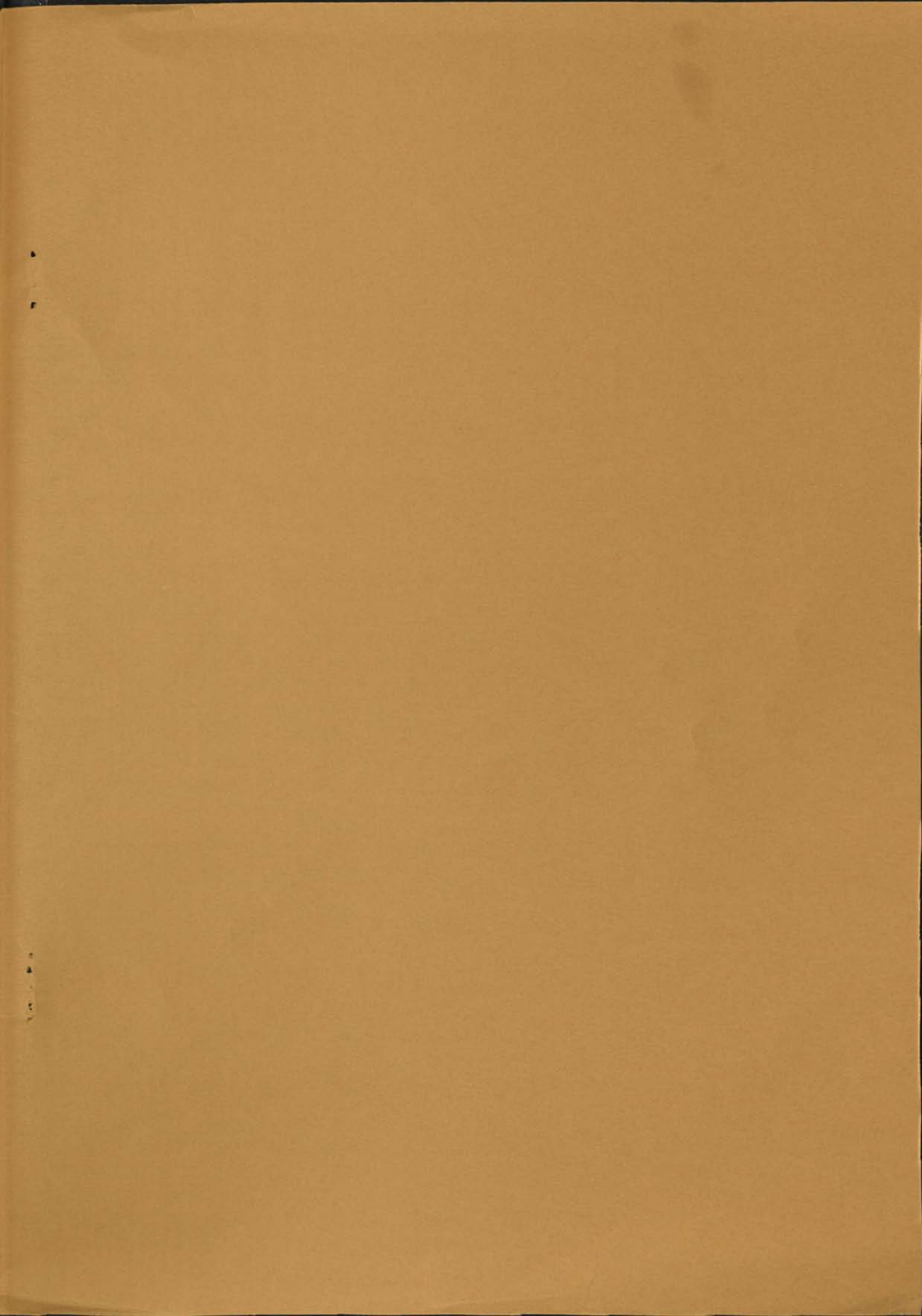


Fig. 5. 256-bit SOS memory



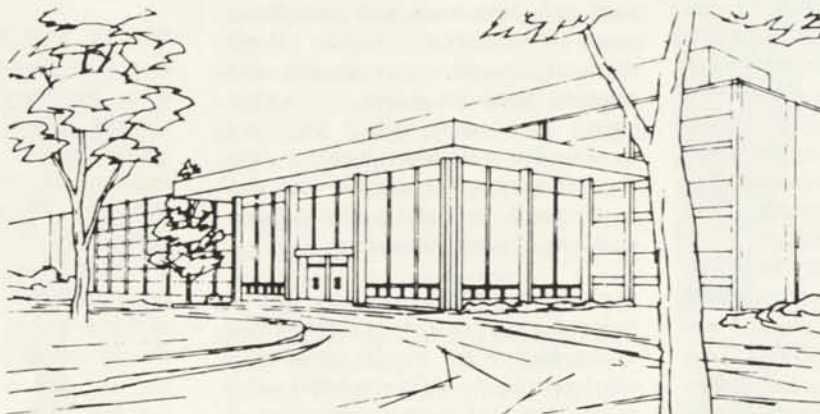
RCA

COS/MOS — from concept to manufactured product

G.B. Herzog

PE-658
ST 6501

RCA LABORATORIES PRINCETON, NEW JERSEY



DAVID SARNOFF RESEARCH CENTER

Copyright 1976 RCA Corporation
All Rights Reserved

PRINTED
IN
U.S.A.

COS/MOS product design

COS/MOS from concept to manufactured product

Giving birth to a new technology requires conviction, dedication, support and the many other demands of parenthood. The struggle to implement the CMOS circuit concept into a product line is a specific example of how a laboratory development was transferred to a product division and ultimately led to the establishment of a market.

G.B. Herzog

THE metal-oxide-semiconductor transistor is based on a relatively old idea. A patent was issued to Lillienfeld in 1930 on the basic concept. The Bell Telephone Laboratories' Research Staff was exploring this concept when the point-contact transistor action was discovered. However, it was not until the semiconductor industry had progressed to an understanding of silicon surfaces that the modern MOS transistor became possible. In 1960 Dr. Webster, as Director of RCA Laboratories' Device Research Laboratory, initiated a project on MOS transistors. By the end of 1960, Charlie Mueller and Karl Zaininger of the RCA Laboratories had shown control of conduction in an insulated-gate structure. A team ultimately consisting of Steve Hofstein, Fred Heiman, and Karl Zaininger began serious research efforts to understand the physics of MOS structures and to build devices with significant gain. Ground work for this effort had been initiated by T. Wallmark, who had been investigating silicon direct-coupled unipolar transistor logic (DCUTL). Tom Stanley of the Laboratories, had shown that DCUTL with its two-dimensional planar structure rather than the three-dimensional bipolar structure was ideally

suited to what today is known as large-scale integration. His analysis indicated that the devices could be scaled down yet retain their speed capability.

NMOS devices

Based on this work, some Government contracts were obtained which had the objective of fabricating large arrays of MOS transistors. Unfortunately, the importance of impurities in the silicon dioxide insulator was not known at that time; and although arrays of individual transistors were made with remarkably good functional yield, their characteristics were unpredictable and unstable. Since the n-type device gave higher gain and higher frequency response due to the higher mobility of the electron carriers, the research work was concentrated on solving the stability problems of n-type transistors.

Because the MOS transistor had a square-law type of transfer characteristic as opposed to the bipolar exponential characteristic, the MOS was better suited to the low distortion requirements of radio and tv amplifier service. Consequently, the first serious production ef-

forts were on transistors for high frequency and very high frequency amplifiers and mixers. These were depletion-mode MOS devices; i.e., normally conducting. Since they could be self-biased in a manner similar to vacuum-tube design, a certain amount of variability or drift in their characteristics was acceptable. Interestingly, these n-type transistors tended to be depletion-mode devices because of the impurities in the oxide. As the production equipment improved and cleaner oxides were grown, it became harder to achieve the depletion-mode conduction specified by the data sheets. Today, essentially the same devices are being produced, but the "impurity" is added by the controlled process of ion implantation.

MOS for large arrays

While the depletion-mode devices were being produced for use in consumer products, the potential of the MOS device was being evaluated for use in other fields. While Stanley had predicted its usefulness in large arrays, the instabilities, the low gain, and the slow switching speed made it unattractive for use in general-purpose computers. The square-law characteristic that reduced cross modulation and made it useful for rf amplifiers, made it a poor switching device as compared to bipolar transistors. It was hard to define when the device was "on" and when it was "off". With characteristics that drifted, a circuit might well produce an output representing a "1" for the same input that previously had given a "0" output. In fact, some of the n-type devices were so bad that they would change state while their characteristics were being observed on a curve tracer.

With the hope that MOS arrays would eventually play a role in terminals and other peripheral computer equipment, research efforts were directed toward improving the speed, stability, and switching characteristics of the MOS transistors. Work concentrated on the n-type devices since they were about double the speed of p-type. The device research groups worked on the physics of the problem; and in a cooperative program, the application groups ran extensive tests on the devices in various ambients (vacuum, nitrogen, etc.) to determine what was causing the characteristics to shift. Sodium in the oxide was eventually

labeled the culprit; and ultra-clean processing, plus phosphorus gettering in the oxide, provided a solution.

Faster circuit speed and a more sharply defined switching characteristic were shown to be possible, both analytically and experimentally, when depletion-mode n-type transistors were used as loads for enhancement-mode n-type switching transistors. Unfortunately, a simple way was not available to accurately and selectively change the surface conduction to provide both depletion-type and enhancement-type devices on the same wafer. Today, the needed control is achieved with ion implantation.

Complementary symmetry—another approach

An alternative and preferable circuit approach was complementary symmetry. This concept, pioneered at RCA Laboratories in the early 1950's, had many boosters, including the author. Consequently, efforts were made to build n-type and p-type devices on the same wafer. While this concept was actually more complex in principle than depletion-mode load devices, the circuit form was more tolerant of differences in transistor characteristics resulting from variations in doping levels. Complementary-symmetry MOS circuits also had many other desirable characteristics which today are well recognized. Still, in the early 1960's it was not clear that the effort to develop such circuits would be worthwhile in terms of RCA's product needs. The Computer Division insisted on the highest speed ECL for main-frame logic and used discrete bipolar devices in its peripherals, etc. There did not seem to be a need for slow logic arrays.

While RCA was concentrating on solving the stability problems of NMOS devices in an attempt to achieve high-speed performance, a venture company was formed to exploit the virtues of arrays of PMOS transistors. Strangely enough, one of the entrepreneurs behind the company was a previous Government employee who had monitored the RCA contracts and must have taken to heart our descriptions of the virtues of MOS transistor arrays. So while we struggled to make n-type devices stable, and virtually ignored the basically more stable, but slower p-type devices, other companies began announcing products. As

pioneers in MOS research, we felt awed and frustrated as various small companies publicized their plans to produce electronic desk calculators containing hundreds of devices in just a few small PMOS chips. Stanley's prediction had come true, but we were not participating. Still, we worried about the soft switching characteristics of the p-type device, the instabilities, and the complex clocking schemes that people were depending on to increase speed.

As history shows, our concerns were well justified. Most of these early, overly ambitious programs ultimately failed, causing severe financial losses to one major calculator company and the complete collapse of the PMOS-array vendor. Nevertheless, it was clear that there would be applications appropriate to the speeds and integration complexity of MOS transistors when the technology was better in hand. How could the RCA Laboratories play a part? Since we were trying to be responsive to the wishes of the Computer Division, we didn't propose work on MOS logic. At that point in time, however, there was considerable interest in content-addressable memories. This was an area of interest to our Computer Division that required logic of modest performance, memory system capability beyond that easily possible with cores, and low standby power. Complementary-symmetry MOS structures were ideal. If only we knew how to make them.

Fortunately, in 1960 Paul Weimer had started work on a different form of device. It was known as the thin-film transistors (TFT) because it used evaporated thin films of compound semiconductor material on a glass substrate. Most of the early TFT devices were so unstable they made silicon MOS devices look like Bureau of Standards references by comparison. They had the great virtue, however, of being easy to make in large arrays and could be made in either conductivity type. Since complementary-symmetry circuits tolerated large variations in device characteristics, I felt that we had a chance of making useful content-addressed memory arrays. I also felt that if real applications could be shown, a substantial device physics effort would be mounted to solve the stability problem. In fact, Paul and his co-workers achieved significant improvements in the stability of both n- and p-type devices. Although RCA dropped its TFT effort in

favor of a silicon-on-sapphire (SOS) thin-film program, other companies are currently pursuing TFT research for physically large-area array structures.

To more rapidly exploit the circuit advantages of complementary arrays for content-addressed memories, I requested a member of my circuit group to work with Paul Weimer's device-processing people. We provided some digital circuit help to Weimer's group while they in turn helped us set up equipment to make arrays of complementary TFTs.

Meanwhile, MOS technology skills were developing in the Somerville Electronic Components semiconductor product group. The Laboratories had been funding efforts to help them learn how to make complementary silicon devices,

Gerald B. Herzog, Staff Vice President, Technology Centers, RCA Laboratories, Princeton, N.J., received the BSEE and MSEE from the University of Minnesota in 1950 and 1951, respectively. He joined RCA Laboratories in 1951 and in 1952 helped design and construct the first completely transistorized television receiver. Subsequently he worked on special color reproducer systems, video tape recording systems, ultra-high-speed logic including microwave and tunnel diode circuits, and large scale integrated circuits, including complementary MOS and silicon-on-sapphire devices. At the RCA Laboratories he has served as Director of the Process Research Laboratory, Director of Digital Systems Research Laboratory, and Director of the Solid State Technology Center (with locations in Princeton and Somerville, New Jersey). Mr. Herzog has presented and published many technical papers on advanced semiconductor device applications and holds 23 U.S. Patents. He is a member of Sigma Xi, Eta Kappa Nu, Fellow of the IEEE, and a past Chairman of the ISSCC. He has received two RCA Achievement Awards, two David Sarnoff Outstanding Team Awards in Science, and the University of Minnesota Outstanding Achievement Award in 1972.



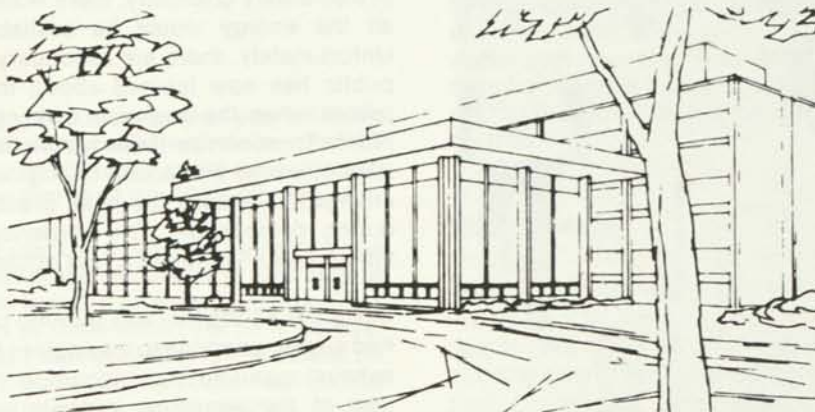
RCA

The potential for electronic automobile controls

G. B. Herzog

PE-677

RCA LABORATORIES PRINCETON, NEW JERSEY



DAVID SARNOFF RESEARCH CENTER

Copyright 1976 RCA Corporation
All Rights Reserved

PRINTED
IN
U.S.A.

The potential for electronic automobile controls

G. B. Herzog

Electronic control can improve gas mileage and reduce pollution in a number of ways—it can cut idling and pumping losses, calculate and control spark timing and transmission shift points accurately, distribute fuel to all cylinders evenly, and possibly even prevent uneconomical driving habits. Taken individually, the mileage gains are not necessarily cost-effective, but the savings can be significant if a number of these concepts are combined into one overall electronic control system.

Gerald B. Herzog, Staff Vice President, Technology Centers, RCA Laboratories, Princeton, N.J., received the BSEE and MSEE from the University of Minnesota in 1950 and 1951, respectively. He joined RCA Laboratories in 1951 and in 1952 helped design and construct the first completely transistorized television receiver. Subsequently he worked on special color reproducer systems, video tape recording systems, ultra-high-speed logic including microwave and tunnel diode circuits, and large scale integrated circuits, including complementary MOS and silicon-on-sapphire devices. At the RCA Laboratories he has served as Director of the Process Research Laboratory, Director of Digital Systems Research Laboratory, and Director of the Solid State Technology Center (with locations in Princeton and Somerville, New Jersey). Mr. Herzog has presented and published many technical papers on advanced semiconductor device applications and holds 23 U.S. Patents. He is a member of Sigma Xi, Eta Kappa Nu, Fellow of the IEEE, and a past Chairman of the ISSCC. He has received two RCA Achievement Awards, two David Sarnoff Outstanding Team Awards in Science, and the University of Minnesota Outstanding Achievement Award in 1972.

Final manuscript received June 10, 1976.



FOR MANY years, there have been great expectations for the expansion of electronics' role in the automobile. At the 1967 IEEE Automotive Conference, a paper that was given stated "However, the 1975 model U.S. automobile will be loaded with electronic devices making life simpler and safer for its driver." As with so many Utopian dreams, economics slowed down the realization of what was technologically feasible. Today, with the price of gasoline approximately twice what it was in 1967, some of these predictions have become more economically meaningful. If the price of gasoline goes to \$1.00 per gallon, quite sophisticated electronic controls will become economically justifiable. Government regulations on fuel efficiency and pollution levels for cars are also hastening the advent of electronic controls. In the following material, I have cited various references to illustrate what is technologically possible. It remains to be seen what electronic control systems can be justified at any given point in time. This information has been gathered in the course of the RCA Laboratories project on microprocessor controls for automobiles, which started late in 1973.

Pollution misconceptions

During the course of this project I discovered that misunderstandings were prevalent among electrical engineers, including myself, about the operation of gasoline engines, their efficiencies, and pollution levels. It seemed to many people early in the fuel-crisis/air-pollution era that a properly operating, efficient engine should generate minimum pollution, or conversely that an engine adjusted for minimum pollution should be efficient, contrary to what seemed to happen to cars in the 1973-1974 period. The simple theory was that if all the fuel burned to carbon dioxide and water vapor, as we learned in elementary chemistry, there would be no pollution and all the energy would be available to drive the car. Unfortunately, there are secondary reactions, which the public has now learned about, that produce nitrogen oxides when the engine is running at its most efficient point. To minimize these reactions, EGR (Exhaust Gas Recirculation) was added to engines and the timing was changed to reduce the peak pressure and temperature during combustion. The efficiency went down as the timing was changed to a less optimum range.

Since not all the fuel was burning completely, some cars had air pumps added to the engines to force air into the hot exhaust manifold and encourage the complete combustion of the remaining hydrocarbons and convert the carbon monoxide to carbon dioxide. Unfortunately, some

This 1951 Lincoln got 25 miles per gallon and won its class in the Mobil Economy Run that year, despite the fact that it weighed 5200 lb. A comparable car today might get half its mileage.



people, even ones with years of engineering experience, were so cynical of our engineering brethren in Detroit that they thought the air pump merely diluted the mixture coming out of the car's tailpipe to comply with pollution testing standards. In fact, the federal standards specify that during an EPA test run, all the exhaust emissions must be captured and the pollutants of each class weighed. The acceptable levels for a given year are specified in grams per mile for a specific driving pattern. These standards are given in Table I, and apply regardless of the size of the car.

Table I — Federal pollution standards.

Model Year	Allowed Emissions (grams per mile)		
	HC	CO	NO _x
1973	3.4	39.0	3.0
1975	1.5	15.0	3.1
1977	0.41	3.4	2.0
1978	0.41	3.4	0.4

These are the standards as originally set by the federal government. Because of the fuel crisis and difficulty of meeting the 1978 standards, the timetable has been extended repeatedly so that the most stringent requirements now fall in the 1980's. It is probable that the NO_x requirement of 0.4 gram per mile will be relaxed somewhat.

Fuel economy— past and present

A 1951 Lincoln weighing 5200 lb won that year's Mobile Economy Run (on a ton-mpg basis) with better than 25 mpg.² A comparable-sized car today might achieve 10 to 15 mpg or a little better than half the 1951 winning rate.

There are many things that can be done to increase fuel economy, and the mechanical engineers in Detroit and around the world are well aware of them. They have not suppressed these ideas as some would like us to believe, but have, in fact, made them available to the public on various occasions. Unfortunately, the public has not had the wisdom to accept fuel-economy devices at the necessary higher initial cost, or the economics has not

truly been favorable. Overdrive is one such example. I had one on my 1951 Mercury, loved it, and got great gas mileage on the highway.³ With gas at 25-30 cents per gallon, however, most people apparently did not want to pay extra for such a feature. It was also the beginning of the era of bigger and bigger engines, responding to customer demand for more performance and convenience. The stick-shift plus overdrive gave way to the automatic transmission, plus power steering, power brakes, air conditioning, etc. Hence, the 20-25 mile-per-gallon car of 1951 became the 10-13 mile-per-gallon car of today.

It is overly simplistic, however, to suggest that we can merely turn back the clock to achieve greater fuel economy. That great Model A Ford we keep hearing about had a low-compression engine and got rather poor mileage for such a light-weight car. Recently, as reported in the June 21, 1976 issue of *Automotive News*, Ford Motor Company ran the EPA mileage test on a 1936 Ford sedan that had spent the better part of its life in their museum. With 26,000 miles of use, it was in good condition, but gave only 15 mi/gal for city driving and 20 mi/gal highway, while generating 12.4 grams of hydrocarbons per mile and 86.8 grams of carbon monoxide per mile. Only the 2.8 grams of nitrogen oxides per mile came close to meeting today's pollution standards. For comparison, today's heavier Ford Granada, with similar size engine and three-speed transmission, gave 22 mi/gal city driving and 30 mi/gal highway. This indicates that good mechanical approaches can achieve significant results.

What further might be done? There are many possible improvements—some simple, some complex, and almost none economically justifiable by themselves. Taken together, however, several of them might be justified if the cost of fuel goes up and the cost of electronics continues to drop.

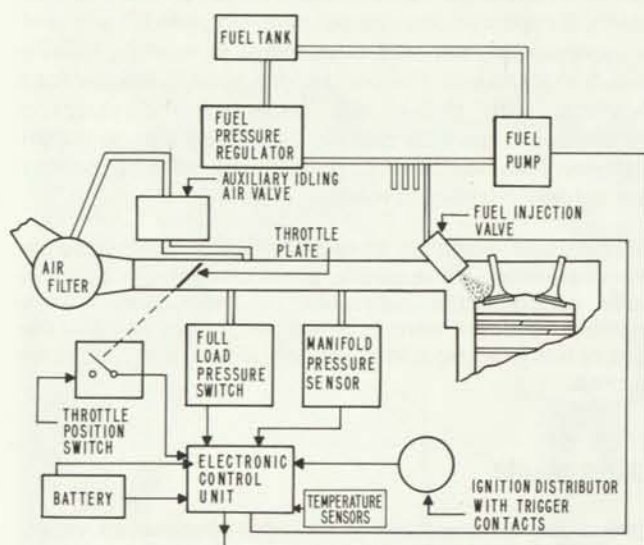
Idling waste

Idling is the most wasteful of all modes of operation, yet an engine may be idling as much as 40-60% of the time in city driving. We all recognize idling at stop signs as wasteful, but it accounts for perhaps only 10% of our city driving time. Other idling periods occur when we abruptly let up on the gas as we slow down for a corner, heavy traffic, or a

stop sign. In city driving we are almost always accelerating and then coasting. Besides wasting gas, this also increases pollution, because rapidly closing the throttle stops the air flow suddenly while the heavier fuel continues to flow for a moment, causing an overly rich mixture. A large V-8 engine burns nearly a gallon of gas per hour idling. If we spend 40% of the time with the engine doing no useful work, it is no wonder that a car gives poor gas mileage in city driving.

Two companies, one in Europe and one in Japan, have developed systems that shut off an engine at stop lights and restart it automatically to save gas. A hot restart of a V-8 engine takes less gas than idling for 24 seconds.⁴ One wonders, however, how long the starter would last with such an attachment. A more practical approach might be an electronic idle speed controller. On large cars with power steering, power brakes, and air conditioning, the idle speed is set high to insure that the engine doesn't stall when the air conditioning load comes on or when parking maneuvers impose a heavy load on the power steering pump. In most cars one can easily detect the drop in engine speed as the air conditioner clutch engages, though some cars now have electromechanical devices to adjust the idle setting to compensate for the added load. The idle speed is often set so high that the car will travel at about 10 mi/h on a level road, and faster if the choke is closed. That means that our brakes are fighting the engine as we try to stop at each stop sign.

An electronic control could maintain engine rpm at the lowest level consistent with pollution standards and still assure power for the accessories. This would save considerable gas in cars used primarily for city driving. But would the cost be justified by the fuel savings? Probably not. But if that same idle control were part of a more extensive control system, however, it would certainly be beneficial. The Bosch fuel injection system used in the VW



This early electronic fuel injection system was developed by Bosch and appeared on some Volkswagen engines in 1968. A set of contacts in the distributor triggered the injectors, and the control unit determined pulse duration after receiving information on engine speed, vacuum, and temperature. The control unit had 220 discrete components, but did not control spark advance.

Squareback and Fastback, circa 1968, shuts the fuel flow off when the throttle is closed and engine rpm exceeds the set idle speed, that is, when the car is pushing the engine. This not only saves fuel but gives improved engine braking, saving wear on the brakes.

Power loads and pumping losses

The most significant problem that the mechanical engineer in Detroit must face is that an automobile engine always operates in the dynamic mode, i.e., it is always changing its speed or power output. However, the testing facilities available to him are designed to test engines running at fixed loads and speeds. It is extremely difficult to measure the engine parameters during the transient as the throttle opens and closes, yet that is what is happening most in normal driving. We may think we drive at a constant speed on the highway, but we are a rather poor actuator in an elementary servo system and we are constantly speeding up or slowing down to hold an average speed. In addition, there is no such thing as a truly flat road, so the power output required is subject to changes even at a constant speed. It is claimed that the "cruise control" option available on some cars will save appreciable gas, especially in mountainous areas. Again, its cost cannot be defrayed by fuel savings, except perhaps for the traveling salesman or others doing a lot of highway driving. Such controls are purchased primarily because of their convenience value.

"A large car requires only about 20 hp to travel 60 mi/h on a level road."

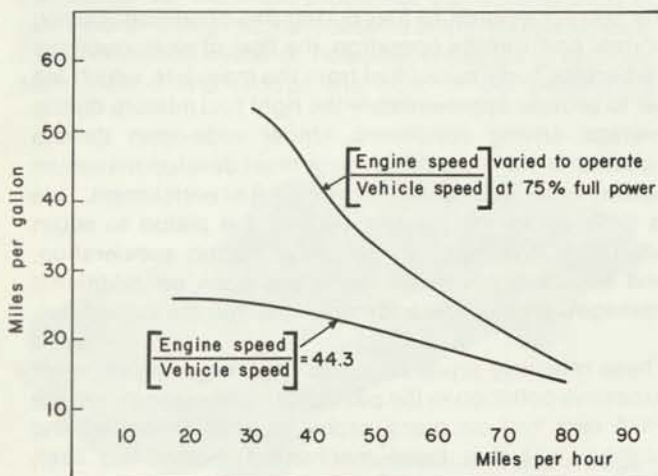
A large car requires only about 20 hp to travel at 60 mi/h on a level road. Yet such cars are often equipped with 200-hp engines because the purchaser demands rapid acceleration and power-robbing accessories. One study estimated that the normal driver uses maximum power only 1% of the time and uses less than 10% of engine power 43% of the time. When such an engine is running at modest speeds, at low power output, much of the engine's power is used to overcome pumping losses, i.e., working to suck air through the closed throttle structure. Even at 50 mi/h, much of the fuel energy input to a 200-hp engine is used to overcome the engine's losses, since only 20 hp is actually required. One reason a diesel engine is more efficient than a spark-ignited gasoline engine is that it operates at nearly wide-open throttle all the time, eliminating much of the pumping losses. Power output is controlled by the amount of fuel injected into each cylinder.

Pontiac engineers did some interesting experiments with gasoline engines several years ago. They modified a V-8 engine so that it could run on only four of its cylinders while the other four were vented to reduce the pumping losses. During idle and cruise, only four cylinders were in use. During periods requiring more power, such as during rapid acceleration, steep hills, etc., the other four

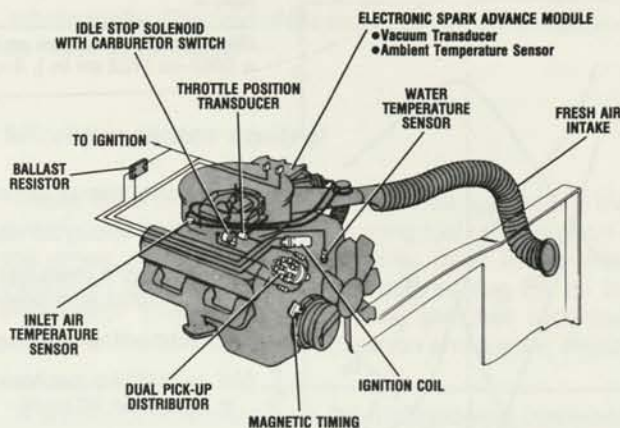
cylinders were brought into service. A 15% increase in gasoline mileage was claimed. With only mechanical linkages and electromechanical control mechanisms available, such a system might be a costly nightmare. But if fuel reaches \$1.00 per gallon and Federal legislation dictates a certain average gasoline mileage for all new cars, such a scheme might be revived for large cars. Hopefully, a microprocessor system will simplify its control. Alternatively, the diesel engine might make such complicated schemes unnecessary.

The electronic transmission

Since the spark-ignited gasoline engine is most efficient near wide-open throttle and at modest speeds, an infinitely variable gear ratio transmission is needed, but not available. The old stick-shift with three forward speeds plus overdrive has become the 4-speed and 5-speed transmission in today's small cars. The large-car purchaser still prefers the automatic transmission, however, even with the losses incurred in the torque converter. The point at which the automatic transmission shifts is a compromise based on fuel economy and performance. Performance takes priority because automotive magazines and the general public rate cars on the acceleration from 0-60 mi/h. When the shift points are controlled by mechanical governors, vacuum valves, and hydraulic pistons, the shift-point compromises made during design are literally "cast in iron". If electronic controls (which have been built and tested by the automotive industry) were used instead, the shift-point compromises could be selected by the driver. The driver could select an economy range for around-town shopping trips, modest performance for casual Sunday rides in the country, or super performance for the hot-rod set. Again, an electronic system just to control the transmission shifting would in itself probably not be justified by a 5-10% increase in gasoline mileage.⁵ However, if it were part of a larger control system it could be achieved at practically no increase in cost over today's hydraulic systems.



An ideal transmission could vary infinitely to operate an engine at the most efficient speed for each vehicle speed. This graph shows the difference such a transmission could make in mileage. The limitation of 400-rpm engine speed does not permit operation at 75% of full power below 30 mi/h. Taken from Kummer, *Technology Review*, Feb 1975. Source: Jandasek, *SAE Transactions*, 61, 95 (1953).



Chrysler Corporation's 'lean burn' ignition system can meet the 1976-1977 pollution standards without a catalytic converter or exhaust gas recirculation. It does this with a special carburetor designed for a lean 18:1 air/fuel ratio, magnetic timing, and electronic spark advance. The system is available only on Chrysler's 400-cubic-inch, four-barrel engine for 1976, but will be used on other engines in 1977.

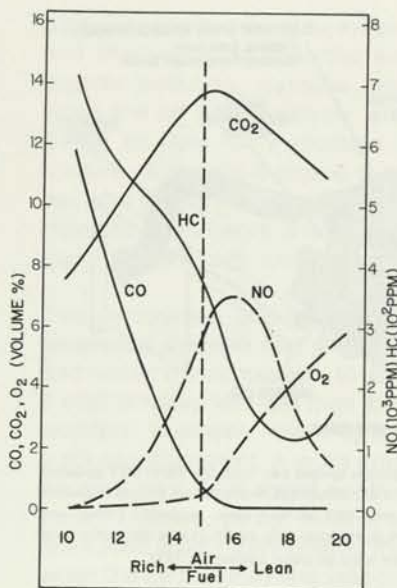
Ignition 'kludges'

The gasoline engine operates most efficiently and develops maximum power when the combustion produces maximum pressure in the cylinder just after the piston reaches the top of its compression stroke. Because the combustion rate is controlled by the fuel-mixture density, which is proportional to manifold pressure as set by the throttle position, and because the amount of time available for burning varies with the speed of the engine, the ignition point must be a variable. In the past, a simple vacuum diaphragm and governor mechanism on the engine's distributor, monitoring manifold pressure and engine rpm, controlled the point of ignition. This control

"...if a lean air-fuel mixture is used in conjunction with electronic control of the spark advance, neither a catalytic converter nor EGR is necessary to meet the 1976-1977 federal pollution standards."

system provided an approximation to the advance that was necessary to give best performance. Unfortunately, the condition of best engine performance also corresponds to a condition of maximum nitrogen oxide production under certain driving conditions. Therefore, over the past several years the distributor mechanism has become more complicated. Additional vacuum diaphragms have been added along with electronic controls used to switch the vacuum from one diaphragm to another, depending upon engine temperature, driving speed, etc. While these mechanical-pneumatic kludges have been satisfactory to meet the pollution standards of the 1974-1976 time frame, they will become increasingly inadequate as the allowable emissions are reduced.

The poor fuel economy achieved on the 1974 cars was due to the retarding of the spark-advance curve to meet the



Varying the air-fuel ratio in an engine varies its exhaust constituents. The vertical line marks the "ideal" stoichiometric air-fuel proportions for combustion; the vertical scales are different for some constituents to make the curves comparable. Taken from Kummer, *Technology Review*, Feb 1975.

Table II.

How thermodynamics and friction reduce mileage in the real-life automobile—a 2000-cc (122 cu in.), 4-cylinder engine with automatic transmission

	Thermal efficiency (per cent)	Miles per gallon	mpg reduction (per cent)
All chemical energy converted into mechanical energy	100	194	
Theoretical Otto cycle (ideal gases, perfect heat release)	57	116	
Real Otto cycle (heat loss dissociation, non-ideal gases, indicated efficiency)	33	68	
After subtracting pumping losses at road load, 40 mi/h	23	44.6	34
After subtracting mechanical losses (engine friction) at road load, 40 mi/h	17	33	26
After subtracting carburetor metering, choke, accelerator pump, fan, manifold distribution, and distributor retard losses	15.5	30	9
After subtracting automatic transmission losses	11.9	23	23
After subtracting power steering and generator losses	10.3	20	14
After subtracting air-conditioning losses (1.5 hp continuous—air-conditioning requires higher idle speed and is used only at certain times)	8.5	16.4	18

(As taken from Kummer, *Technology Review*, Feb 1975)

pollution standards of that year. When the 1976 standards took effect, the engines needed catalytic converters to further reduce the excess hydrocarbons and convert carbon monoxide to carbon dioxide. This allowed the engine designers to readjust the spark advance toward the optimum condition, resulting in improved fuel economy for the 1976 cars versus the 1974 cars. Chrysler Corporation's 1976 400-cubic-inch "lean burn engine" demonstrates that if a lean air-fuel mixture is used in conjunction with electronic control of the spark advance, neither a catalytic converter nor EGR is necessary to meet 1976-1977 federal pollution standards. Also, in an experimental project with a Mark IV, Ford Motor Company reported a 10-20 percent increase in mileage through the use of electronic spark advance and EGR control.⁶ This car did use a catalytic converter in addition, however.

Carburetion vs fuel injection

The gasoline engine operates best when the air-fuel mixture is near the stoichiometric ratio, that is, when there is just sufficient oxygen present to theoretically cause the fuel to burn completely. Unfortunately, this is also the

"What once was a fairly simple device has become another electromechanical kludge."

condition that generates the most nitrogen oxides. A lean mixture of 16 to 18:1 versus the stoichiometric ratio of 15:1 gives greater efficiency and less pollution, but less than maximum power from the engine. It is also more difficult to ignite than a rich mixture. Recent cars use so-called high-energy electronic ignition systems to ensure ignition of the leaner mixtures being used in the 1976 cars and to

prevent misfires that would put raw fuel into the catalytic converter and cause excessive heating.

Controlling the air-fuel ratio to the accuracy necessary to meet both pollution standards and fuel economy is quite difficult for the standard carburetor. What once was a fairly simple mechanical device has become another electromechanical-pneumatic kludge. The carburetor works on a fairly simple Venturi principle based on Bernoulli's laws of mass flow. Unfortunately, because of the wide dynamic range over which the engine operates, many approximations must be made. With the engine idling, the throttle is nearly fully closed, and very small amounts of air flow through the main body of the carburetor. Separate idle jets must be provided to supply the correct amount of fuel during this condition. During normal part-throttle operation, the flow of air through the carburetor body sucks fuel from the main jets, which are set to provide approximately the right fuel mixture during average driving conditions. Under wide-open throttle conditions, however, the engine must develop maximum power, and it is necessary to provide fuel enrichment. This is done by various means, including a piston to squirt additional fuel into the carburetor during acceleration, and so-called power valves, which open up additional passageways to allow additional fuel into the carburetor.

These relatively crude schemes have been the cause of excessive pollution in the past. One of the reasons that the 1974 cars had so many problems with hesitation and stumble was that these mechanical means had been refined to their practical limits. With the catalytic converter, though, somewhat larger amounts of excess fuel can be supplied to the engine at acceleration. The converter is then called upon to burn up the excess fuel to prevent it from being emitted as a pollutant. We are

"One of the reasons that the 1974 cars had so many problems with hesitation and stumble was that [with all the pollution controls added] these mechanical [carburetion] means had been refined to their practical limits."

currently hearing, however, about the intense heat that the catalytic converter gives off when it burns excess fuel. It is claimed that this can be a fire hazard near locations where fuel has been spilled or in high-underbrush areas.

Another problem with existing carburetor systems has to do with the fuel-mixture distribution. On the way from the carburetor to the various cylinders in a large engine, the fuel-air mixture changes composition. As the fuel travels to the furthestmost cylinders, some of the fuel droplets that are not fully vaporized strike the intake manifold walls and run along the walls as rivulets rather than as an air-fuel mixture. This causes some cylinders to run with a leaner mixture than others, so the carburetor must be set so that the cylinder with the leanest mixture still has a mixture rich enough for proper ignition.

An electronic fuel injection system can help overcome this problem by injecting fuel in the vicinity of the intake valve for each cylinder. This allows more precise control of the air-fuel ratio under all driving conditions, but not without a price penalty. The Cadillac Seville for 1976 has such a system. Fuel injection trades the problem of intake manifold design for the problem of exactly matching fuel injectors, which must be precisely made and are consequently expensive. In addition, some accurate way of measuring the air-mass flow must be devised to determine how much fuel to inject for a given operating condition. This has been a major problem in implementing fuel injection systems that must meet both performance requirements and pollution standards. As presently achieved with analog control circuitry, the cost is prohibitive for the mass-produced, standard-sized U.S. automobile. Using microprocessors in the next generation of fuel injection systems, however, promises to improve performance and reduce the cost of the electronics



Cadillac used electronically-controlled fuel injection to obtain 'improved driveability' on the 1976 Seville's 350-cubic-inch engine. Electronic fuel injection systems will probably remain luxury items, however, until microprocessor control and mass production bring costs down.

substantially. The cost of the fuel injectors can be reduced by mass production.

Microprocessor control

Finally, there is one major element in the automobile that can play a significant role in reducing fuel consumption—the driver of the car. Unnecessarily rapid acceleration wastes an enormous amount of gas. Letting up on the accelerator suddenly also wastes gas and increases pollution. A slow, steady acceleration and a slow, steady

"A microprocessor can be interposed between the driver and the engine to force better driving habits."

deceleration is the best mode of operation. A microprocessor can be interposed between the driver and the engine to force better driving habits. By putting the microprocessor in the driving loop, the car can be inhibited from unnecessarily rapid accelerations, except in emergency situations where, as with the kickdown feature of an automatic transmission, more rapid acceleration could be achieved upon driver command. For average driving the operator could select a microprocessor control mode designed to maximize the automobiles' fuel economy. This efficiency mode would suffice for a high percentage of driving conditions, but a high-performance mode could be selected by the driver when desired.

Conclusions

In summary, it seems certain that the microprocessor will play an important part in the car of our future; it is only a question of how far in the future. Other papers in this issue go into more detail on the contributions that RCA is making toward accomplishing that Utopian goal predicted in the 1967 paper referenced earlier. For those readers wishing more details on the what has been or could be tried to improve fuel economy, the February 1975 issue of the *MIT Technology Review* contains several papers that make extremely interesting reading.

References

1. Hugle, W.B.; "How microcircuits will be used in cars," *IEEE Automotive Conf. Record*, (Sep 21-22, 1967) pp. 1-9.
2. Cohn, C.E.; "Improved fuel economy for automobiles," *MIT Technology Review*, (Feb 1975) p. 52.
3. Cohn, C.E.; "Improved fuel economy for automobiles," *MIT Technology Review* (Feb 1975) p. 49. ("Overdrive, which was quite popular twenty years ago... can yield a 15 per cent improvement in fuel economy").
4. Cohn, C.E.; "Improved fuel economy for automobiles," *MIT Technology Review*, (Feb 1975) p. 49. (The companies are Toyota in Japan and E. Jucker Relaisbau of Zurich, Switzerland).
5. Kummer, Joseph T., Chemistry Department, Ford Motor Co.; "The automobile as an energy converter," *MIT Technology Review*, (Feb 1975) p. 27. ("If the transmission were more efficient in urban driving and if it could provide a better match between the engine speed and vehicle speed, there would be a large gain in fuel economy, perhaps in the ideal case as much as 50% for a vehicle with a large engine in urban driving.")
6. Moyer, D.F.; and Mangrulkar, S.M.; "Engine control by an on-board computer," *SAE Convergence '75, Automotive Electronics II*, (Feb 1975) pp. 75-77.