

Structured Hardware Design

New revised 6L 1A edition.

DJ Greaves

April 28, 1998

- 1 Building Blocks for Digital Electronics.
- 2 Clocking and Synchronous FSM Combination.
- 3 Technology, Speed, Power and Gate Delays.
- 4 Hardware and Software Design Partition.
- 5 Modular Design and Design Partition.
- 6 Further Circuit Structures.

Further copies of these notes can be printed from the Web page on the front cover. Copyright for these notes belongs to DJ Greaves. © 1998.

Preface

This is the new, revised, six lecture Structured Hardware Design course for 1A. The material covered is different from the older 1B course. Some of the material previously covered is now in the new 1B 8L ECAD course.

Bibliography (Book list)

Books related to the course, in approximate order of importance, are:

W.Ditch. *'Microelectronic Systems, A practical approach.'* Edward Arnold. The final chapters with details of the Z80 and 6502 are not relevant to this course.

Floyd. *'Digital Fundamentals'* Prentice Hall International.

T.J. Stoneham. *'Digital Logic Techniques'* Chapman and Hall. This is a basic book and relates more to the previous course on Digital Electronics.

Randy H Katz. *'Contemporary logic design.'* Benjamin Cummings ISBN 0 8053 2703 7

Texas Instruments. *'System 74 Series Logic Family.'* Texas Instruments.

Oldfield, Gray, Kean, Dorf. *'Field Programmable Gate Arrays.'* New York; Chichester: Wiley, c1995.

L.J. Herbst *'Integrated circuit engineering: establishing a foundation.'* Oxford University Press, c1996.

Glossary of jargon and acronyms

Here are some acronyms:

ALU	Arithmetic and logic unit.
ASIC	Application specific integrated circuit.
BICMOS	A process with both CMOS and bipolar transistors on one die.
CAD	Computer aided design.
CAE	Computer aided engineering.
CAM	Computer aided manufacture.
CLB	Configurable logic block.
CRC	Cyclic redundancy check (added to end of a data frame).
CMOS	Complementary metal oxide of silicon.
DRAM	Dynamic random access memory.
DMA	Direct memory access.
DSP	Digital signal processor/ing.
EDO	Extended Data Out for DRAMS.
ECL	Emitter coupled logic.
EMC	Electromagnetic compatibility.
EMI	Electromagnetic ingress or i(m)missions (sic).
FET	Field effect transistor.
FSM	Finite state machine.
FPGA	Field programmable gate array.
GaAs	Gallium Arsenide.
HDL	Hardware description language.
IEEE	Institute of electrical and electronic engineers.
IPR	Intellectual Property Rights.
JTAG	Joint technical advisory group
LSI	Large scale integration (say about 5000 gates, 10000 transistors).
LCA	Logic cell array.
MCM	Multi-chip module.
MSI	Medium scale integration (i.e. the larger ttl devices).
NRE	Non recurring engineering costs. i.e paid once per design.
OTP	One-time programmable
PAL	Programmable array logic.
PIO	Programmed Input and Output.
PCB	Printed circuit board.
PLL	Phase-locked loop.
PIN	Personal identification number.
RTL	Register transfer level.
SRAM	Static RAM.
SIMM	Single in-line memory module.
SSI	Small scale integration (i.e. about 4 gates on a chip).
TTL	Transistor-transistor logic.
UART	Universal asynchronous receiver and transmitter.
Verilog	An HDL language in wide use (not an acronym).
VHDL	VHSIC Hardware Description Language
VHSIC	Very High Speed Integrated Circuit
VNL	Verilog netlist.
VCO	Voltage controlled oscillator.
VLSI	Very large scale integration (big chips).
XNF	Xilinx netlist format.

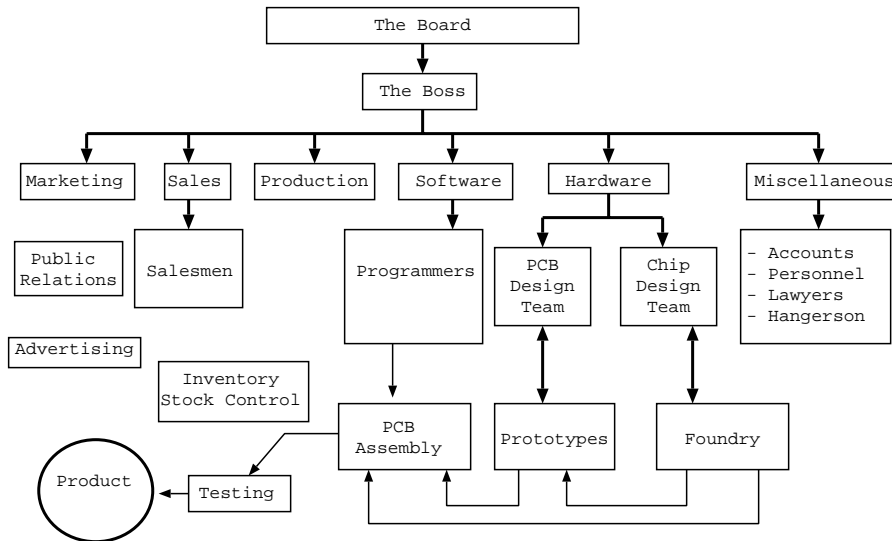


Figure 1: Structure of a small company with hardware products that use its own hardware, software and chips.

Introduction

A good hardware design works the first time. A structured approach to its design is required if the system is complicated. Simulation of the system behaviour before fabrication is also vital to help get it right. As Figure 1 shows, the final product relies on successful interworking between custom chips and the other parts on the circuit boards, together with the software.

Most designs today are held in a database which contains details of custom circuits, printed circuit boards, other chips, software, part numbers and so on. The database can help automate every part of the design process, except for the initial design decisions. These decisions define the basic structure of the product and they are what this course is about.

To create a new hardware design, as with software, it is important to use a top-down approach, starting with an accurate capture of the requirements specification and then decomposing the system into modules. As with software, the module must have well-defined functions and well-defined interfaces and the modules must be suitable for individual testing. However, two big differences are 1. that the modules in hardware are often made of different technology from one another, and 2. that the desirability of using pre-existing modules is much greater.

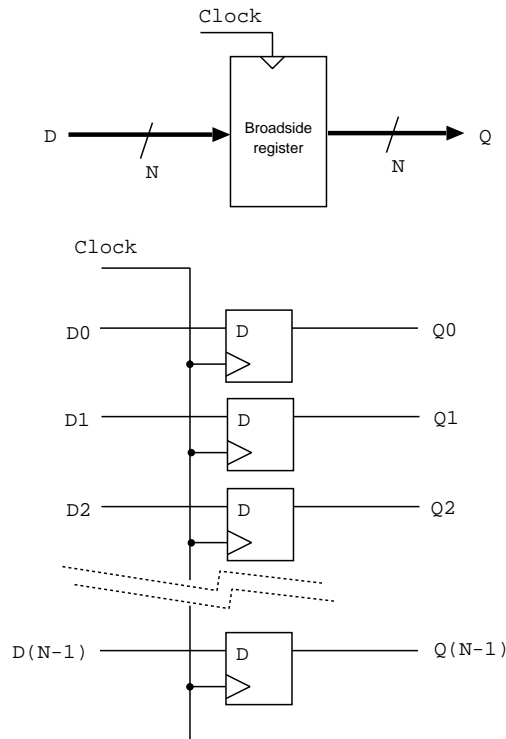


Figure 2: A broadside register: schematic symbol and internal structure.

1 Building Blocks in Digital Electronics.

This section introduces a number of building blocks that are frequently used in hardware designs. In times gone by, these were referred to MSI (medium scale integration) and LSI (large scale integration) subsystems. Today, many of these blocks rarely exist as separate components: instead they are placed inside custom (or semi-custom) application-specific, integrated circuits (ASICs).

1.1 Busses and Broadside Components.

A bus is a number of signals that operate in parallel to carry a binary number. In the schematics, a thicker line is used, and the number next to the slash is the number of signals in the bus.

Figure 2 demonstrates the concept of a broadside component, in this case, a broadside set of D-types. To form a broadside component, the building block is instantiated a number of times with the data inputs and outputs connected to the member lines of input and output busses. The control connections are paralleled and driven from a single external source. A broadside set of flip-flops is a *register*.

Figure 3 shows the schematic and a suitable circuit implementation for a broadside, two-input multiplexor. Most systems will regard a two input multiplexor as a fundamental building block and implement it at the transistor level, rather than from gates.

Figure 4 shows the schematic symbol for a dual-port register file. This device internally contains a number of broadside registers and multiplexors, but they are grouped into the file for convenient access. Such a file is a key part of all computers. The illustrated dual-port file has two reading ports and one writing port, so might be considered a tripple-ported file. On the positive edge of the clock, the data on the write port is stored in the addressed register. For reading, at any time, like an SRAM, the read address may be changed and the output data will change shortly afterwards.

Exercise: Sketch the circuitry that would turn a collection of standard broadside registers into a dual-port file. You will have to use one of the clock enabling techniques defined in section 2.

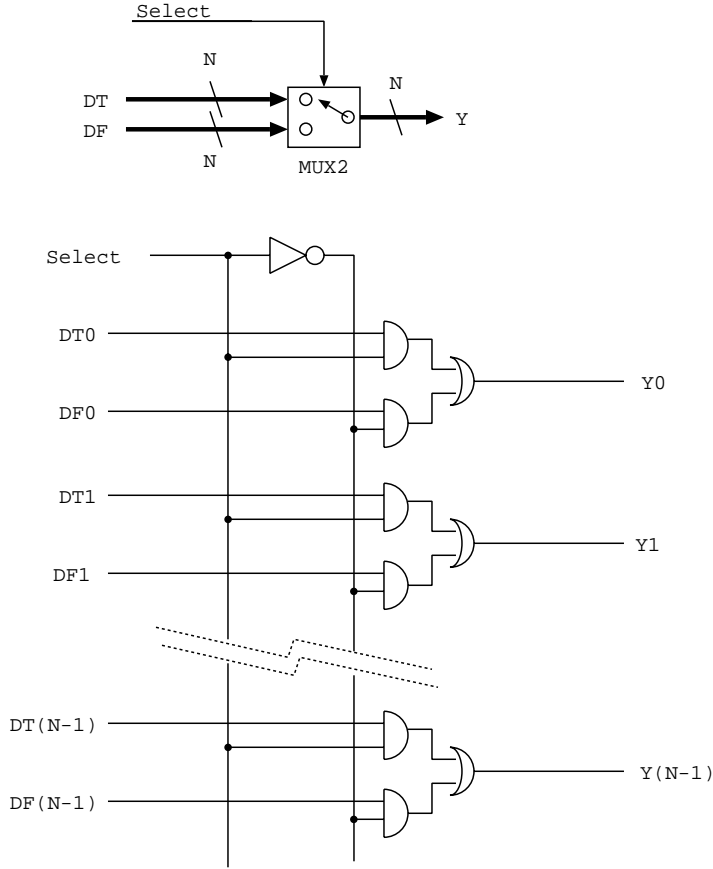


Figure 3: An N-bit broadside, two-to-one multiplexor.

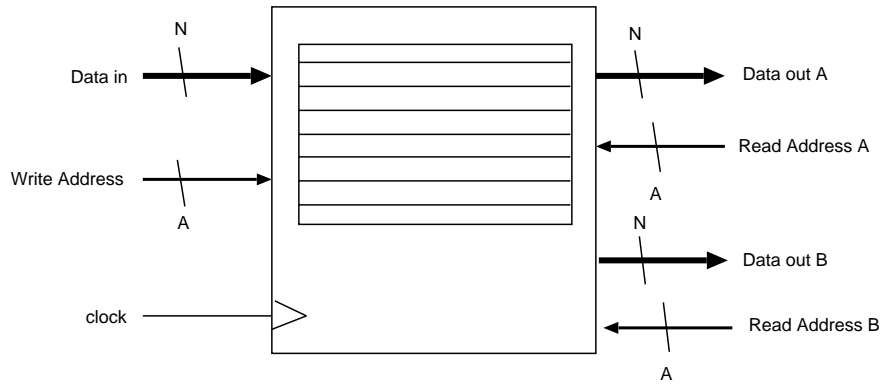
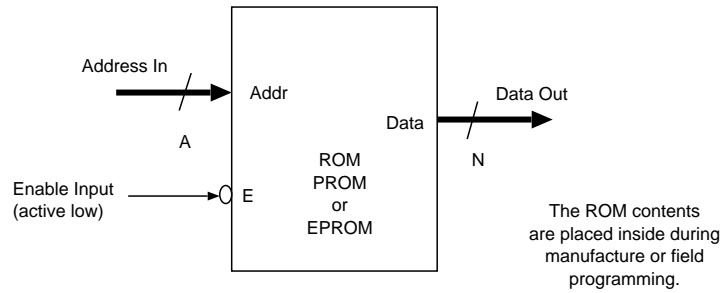


Figure 4: Dual port register file.



The ROM takes A address bits named A0 to A<A-1> and produces data words of N bits wide. For example, if A=5 and D=8 then the ROM contains 2*5 which is 32 locations of 8 bits each. The address lines are called A0, A1, A2, A3, A4 and the data lines D0, D1, ... D7

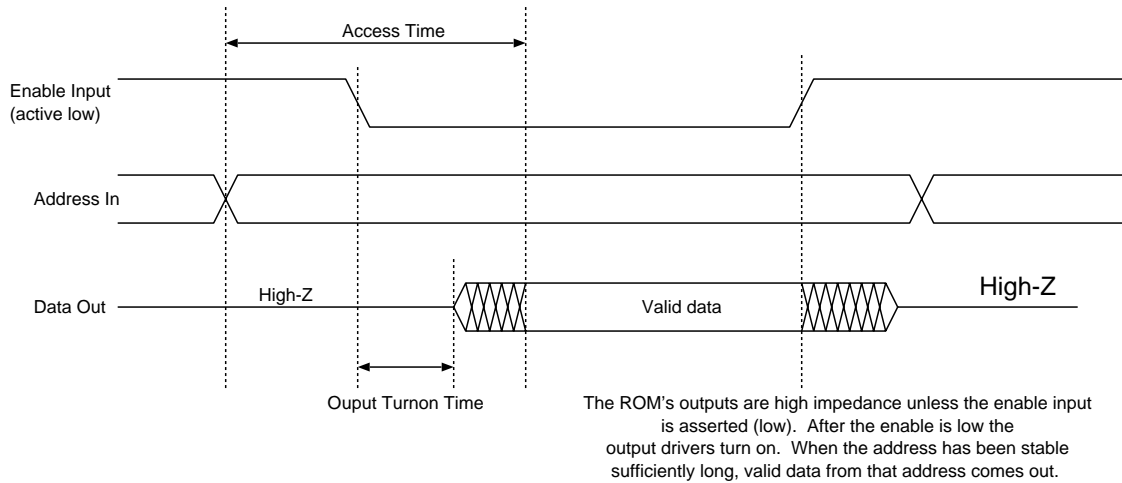


Figure 5: Read Only Memory (ROM).

1.2 Read Only Memories

Figure 5 shows the logic symbol for a read only memory (ROM). The contents are non-volatile and are placed in the ROM during manufacture using a fabrication mask or with a special programming step in the field. The techniques for storing field programmed information are:

- W-Sn fuse: One time programmable by melting selected fuses. This is an old technology which is not very reliable or dense. When used for ROM gives a PROM.
- Antifuse: One time programmable, high density process. Used on one contemporary family of FPGAs. Not used in ROMs today.
- Floating gate technology with ultraviolet erase. A static electrical charge is stored on the floating gate and is put there using a supervoltage that causes tunneling or other breakdowns in the insulator. Discharge is achieved using intense ultraviolet illumination to make the charge leak away. When used for a ROM, gives an EPROM.

A variant on EPROM is an OTP-EPROM (one-time-programmable) which the same device packaged in a windowless package at much lower cost.

- Electrically reprogrammable floating gate devices. Again charge is stored on the gate, but electrical techniques used for programming can be applied for erasing. Used widely in personal organisers and (Flash) memory cards and in most modern PAL devices. Often, erase of individual bits is impossible: only block erase is supported.

Mask programmed ROMs of several megabytes may cost sub 50 pence each in quantities of tens of thousands. Flash memory may be more than 100 times more expensive. The choice of which tech-

nology to use depends on the number of units that are going to be made and the expected number of reprogramming cycles required. There is always a testing problem with one time programmable devices!

1.3 RAM memories

Figure 6 shows a random access memory or RAM. It is as equally random access as a ROM, so is a misnomer. RAM is normally volatile, meaning it loses its contents when power is disconnected. Static RAM (SRAM) normally uses one flip-flop per bit stored. A transparent latch suffices, so it is not as complicated as having a D-type per bit. A possible internal structure for an SRAM is shown in Figure 7. However this would lead to a long, skinny chip. In practice, a square array of storage elements is used, meaning that the on-chip tri-state bus is much wider than the external data bus. The on-chip bus will have width approximately equal to the square root of the number of bits stored in the device. Modern SRAM devices have an access time of 20 nanoseconds and have capacities of about half a megabyte. Price is about six pounds per megabyte, but depends greatly on the access time required. Second-level cache chips for PCs are SRAMs and produced in high volume, so are currently a good way to purchase SRAM, even for other designs or uses.

A property of the SRAM, shared with the ROM, is that there is no clock and, during read, the output is essentially a combinatorial logic function of the address input value.

When writing, the address must not change and the data present on the data pins at the back edge of the (active low) write pulse is the data actually written.

Figure 8 shows the schematic symbol for a dynamic RAM. DRAMs store their data not in flip-flops, but in small capacitors. The value stored leaks away over time unless refreshed. It is also destroyed by readout: after readout, the device always writes back the value internally, removing the job from external circuitry, but meaning that the external logic must not start a new cycle immediately. Therefore the cycle time and access times of DRAM are different. Typically, today, DRAM cycle and access times are 60 and 120 nanoseconds and capacity is 2 to 4 megabyte.

All DRAMs are accessed by applying the address in two halves. First a row address is presented and RAS taken low, then a column address is applied, and CAS is taken low. These row and column addresses typically correspond to the actual physical topology of the rectangular array of bits on the device. In external view, some DRAMs are one bit wide and others are four.

DRAM is much cheaper than SRAM owing to the smaller silicon area needed per bit. Smaller packages are also possible, owing to the multiplexed address bus.

Refresh cycles do not transfer data in or out of the chip, but must be done sufficiently often anyway. A typical specification is that 512 refresh cycles are executed every 4 milliseconds. Since DRAM consumes most of its power when being refreshed, a designer who puts all of the refresh cycles in one batch at the end of every 4 millisecond period makes the design of power supplies much more difficult.

Modern EDO (extended data out) DRAMs guarantee that the data will be valid on the output for a period after CAS is deasserted, thereby easing system timings and helping the designer.

DRAMs are often packed on SIMMs (single in line memory modules). Capacities of SIMMs are around 16 to 32 Mbyte with prices of tens of pounds per SIMM.

High performance systems today use synchronous DRAMs (SDRAM) which include broadside registers on the address or data connects. This gives a pipeline delay (section 2.4) but allows faster operation.

Exercise: When would you use SRAM and when DRAM ?

Exercise: Sketch the circuitry that could be put around a bank of DRAM to make it appear like an SRAM with a single, non-multiplexed address bus and similar control signals. What shortcomings have you made ?

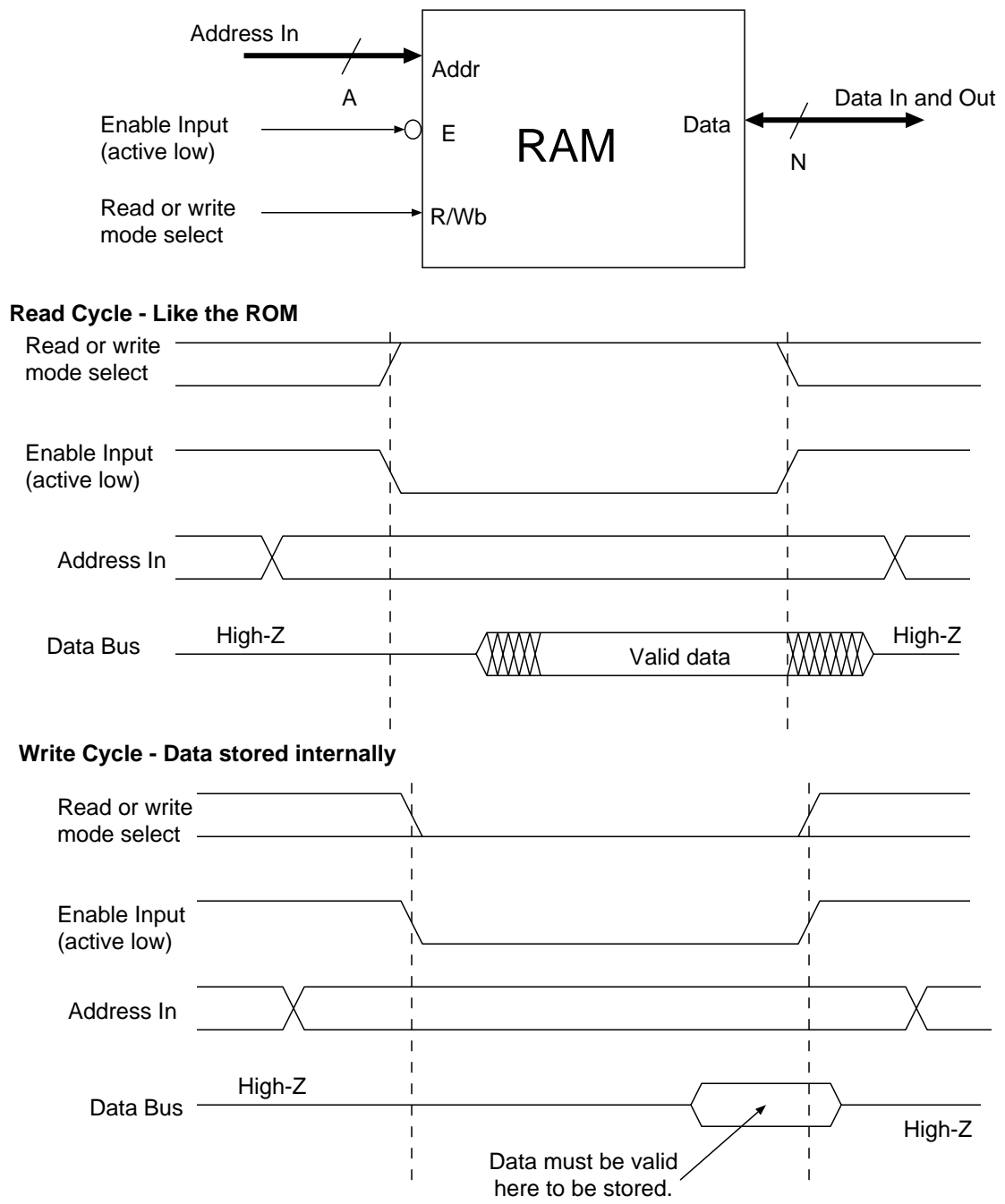


Figure 6: Read and Write Memory (RAM).

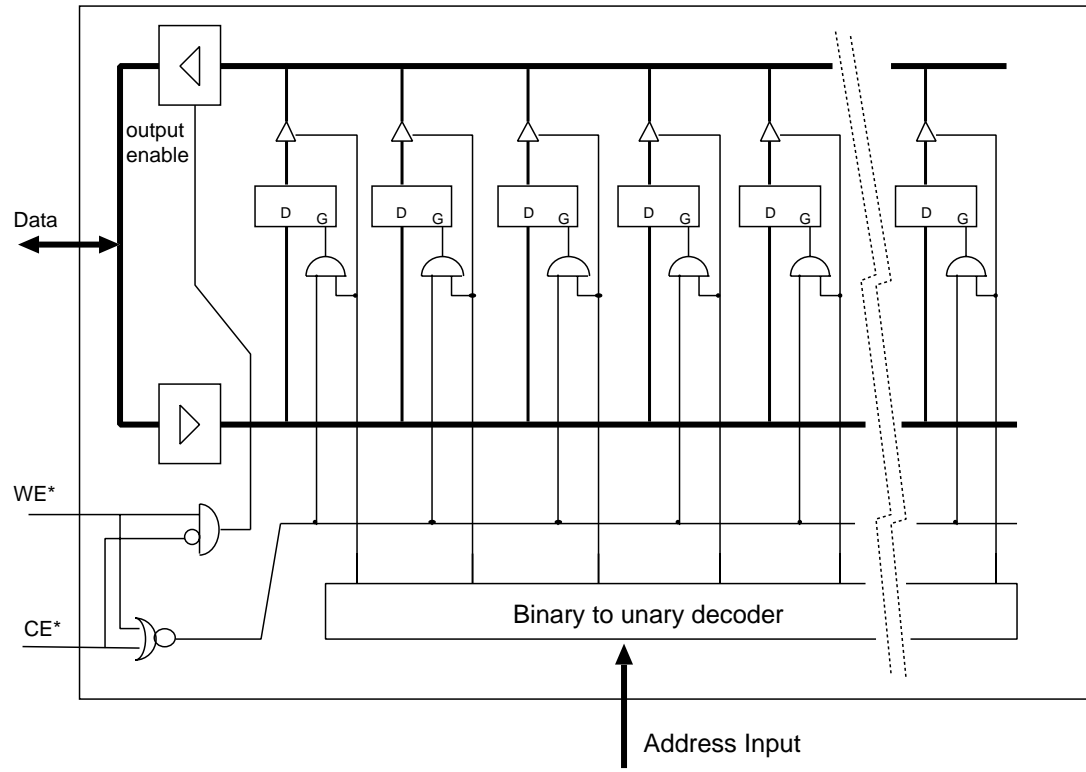
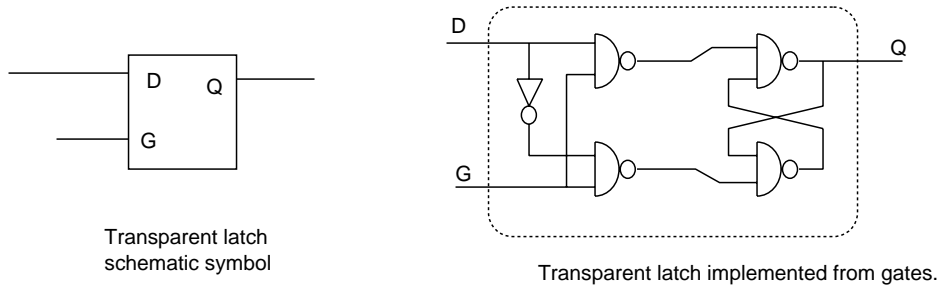
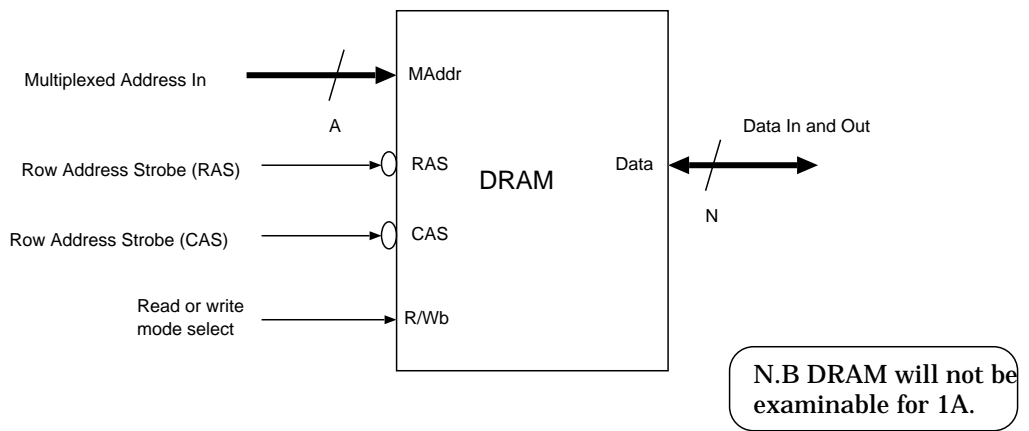
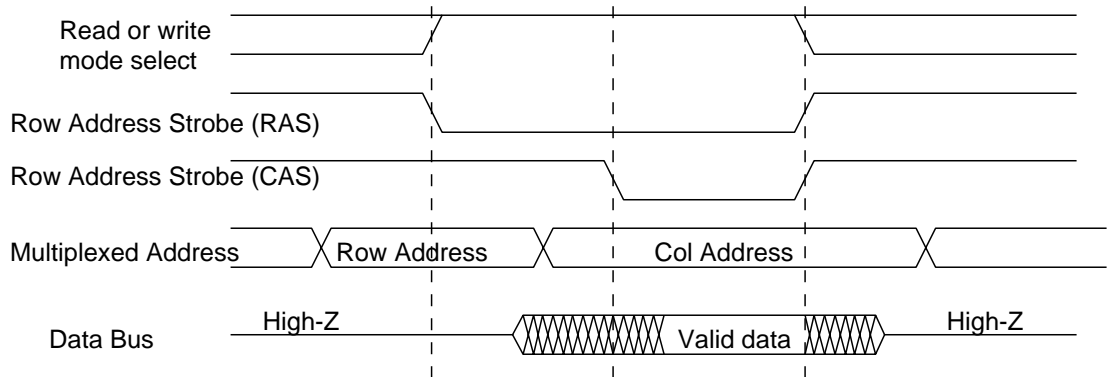


Figure 7: Schematic and possible implementation of a transparent latch and a naive implementation of an SRAM using broadside transparent latches.



Read Cycle (write is similar)

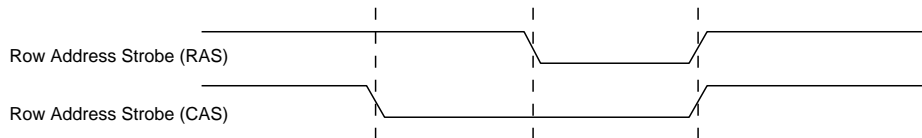


A DRAM has a multiplexed address bus and the address is presented in two halves, known as row and column addresses. So the capacity is $4^A \times D$. A 4 Mbit DRAM might have $A=10$ and $D=4$.

When a processor (or its cache) wishes to read many locations in sequence, only one row address needs to be given and multiple col addresses can be given quickly to access data in the same row. This is known as 'page mode' access.

EDO (extended data out) DRAM is now quite common. This guarantees data to be valid for an extended period after CAS, thus helping system timing design at high CAS rates.

Refresh Cycle - must happen sufficiently often!



No data enters or leaves the DRAM during refresh, so it 'eats memory bandwidth'. Typically 512 cycles of refresh must be done every 8 milliseconds.

Figure 8: Dynamic RAM (DRAM).

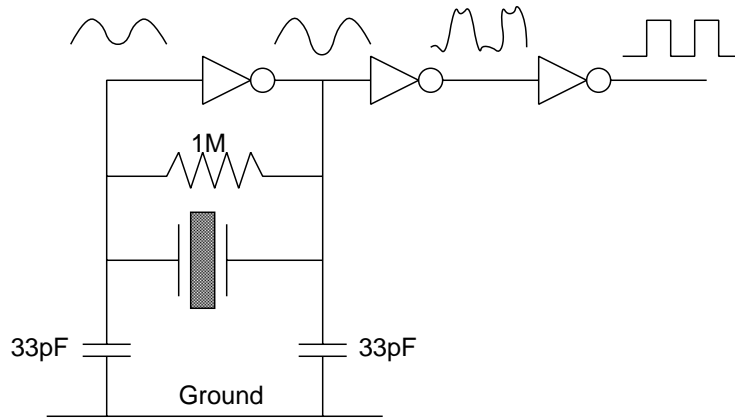


Figure 9: A crystal oscillator clock source.

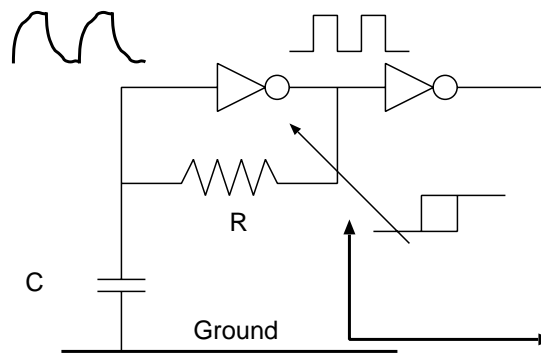


Figure 10: An RC oscillator clock source.

1.4 Miscellaneous Components

Figure 9 shows a circuit frequently used as a clock source. A carefully cut slither of quartz crystal is used as the timing reference. Quartz is pizeoelectric, so can be made to change shape by applying an AC signal across metal coatings applied to its sides. Also, as it changes shape, it generates an equivalent electrical signal. Since the speed of sound in quartz is very high, a slice about 100 microns thick will resonate at several MHz. The details of the circuit are beyond the scope of this course, but suffice to note that the sine waves generated by the crystal become good quality square waves, suitable for use as a clock, after passing through a few inverters. A crystal might cost one pound to add to a design, but is accurate to about 50 parts per million in frequency, so is highly useful.

Figure 10 shows a cheaper circuit that serves as a clock when accuracies of only a few percent are needed. The first inverted must be a special Schmidt inverter with abrupt switching and hysteresis for the circuit to work. The capacitor is the only component that is hard to make on an integrated circuit, so a single pin on the device serves to connect the external capacitor. The crystal oscillator always needs two pins.

Figure 11 shows the typical arrangement of clock distribution inside a large, modern IC. The clock is actually provided from a lower-frequency external reference and multiplied up internally with a phase-locked loop. Skew in the delivery to the various parts of the device is minimised by using a balanced clock distribution tree. Inverters are used instead of buffers to minimise pulse shrinkage (duty-cycle distortion).

Figure 12 shows a circuit that is often used to generate a reset signal when a system is switched on. When power is applied, the capacitor is initially discharged, so has zero volts across it. Therefore the output is a logic one. After a while, the resistor charges the capacitor such that the input to the inverter crosses the switching threshold of the Schmidt input and the output from the circuit

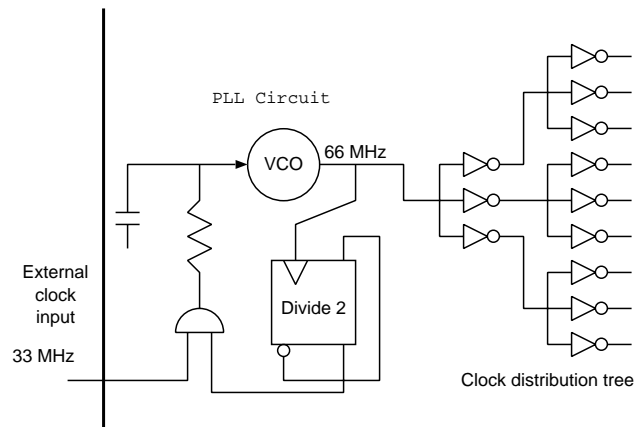


Figure 11: Typical clock multiplication and distribution system for a large IC.

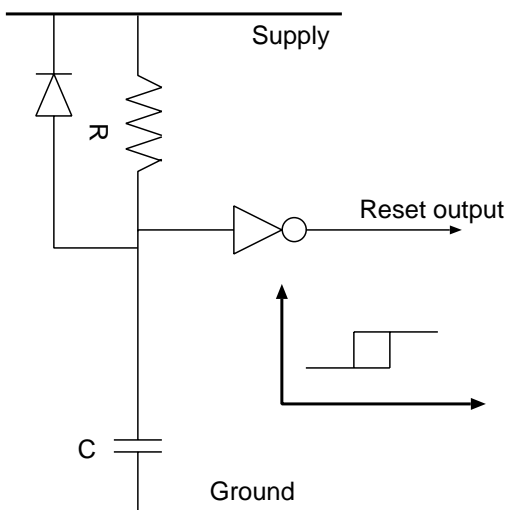


Figure 12: A circuit often used as a power up reset.

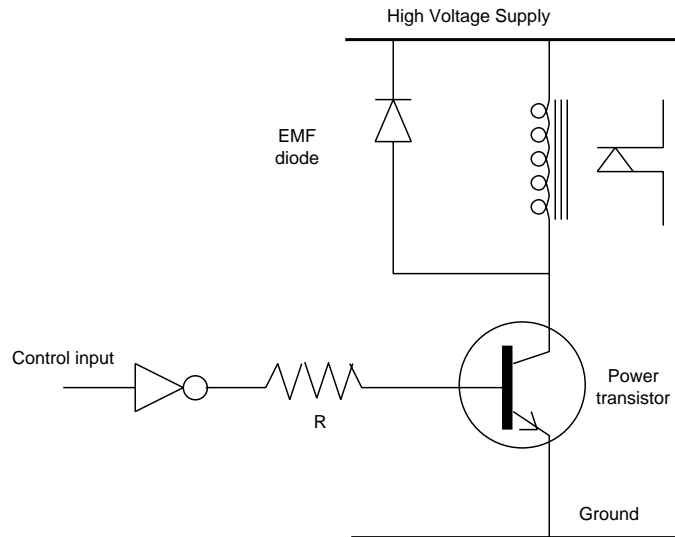


Figure 13: A typical structure used to drive a heavy current or high-voltage load.

is deasserted. Typical delay times used are about half a second, which is long enough for the power supply to become stable and for oscillators to settle down. The diode discharges the capacitor into the supply rail when power is removed. On most devices, a reset switch is also fitted, which discharges the capacitor when pressed, therefore starting the device as from power up.

Figure 13 shows a typical structure used to drive a high power load from a logic signal. The semiconductor processing used to fabricate logic devices may be too expensive to waste on large transistors or may not be able to tolerate the higher voltages required to control the load. Therefore, external driver transistors in individual packages or in power driver ICs are used. Typical applications are motor, loudspeaker, solenoid, relay, and print head driving.

Figure 14 shows the circuit typically used to get a clean signal out of a mechanical switch. In such a switch, when one contact hits another, it bounces off with a supersonic ping, which means that the circuit is made and broken a few thousand times for each switching operation executed. Using a two-pole switch and an RS latch formed from a pair of cross-coupled NAND gates, a clean output is generated. The system works on the basis that the *first* time the commuting contact hits the stator contact, the latch changes state, and stays there until the switch is moved all the way back to the other side, where the reverse happens.

1.5 ALU and ALU circuits

The arithmetic and logic unit (ALU) is a fundamental block of all computers: it is the basis of the integer execution unit. Figure 15 shows the normal schematic for an ALU. The ALU is combinatorial, but is shown connected to a flag register, which does have a clock input. For a detailed example, look at the 74181 in the System 74 databook.

The illustrated ALU has two N-bit inputs and a 4 bit function code input. The output is also N-bit. The function code determines what function of the two inputs is computed. Typical functions are add, add with carry, bitwise AND and OR, subtract and identity functions of the two inputs. The instruction set for the ARM microprocessor contains 16 data manipulation instructions which are a good example of the 16 most useful functions that an ALU can perform.

Like all combinatorial logic functions, the ALU is guaranteed to compute its output within a fixed time interval from the last input change. The value of this delay is typically dependent on the carry chain speed inside the ALU and it will set the maximum clock frequency of a simple microprocessor.

Question: Suggest why ALU's typically only support addition and subtraction as arithmetic operations and not multiplication or division ?

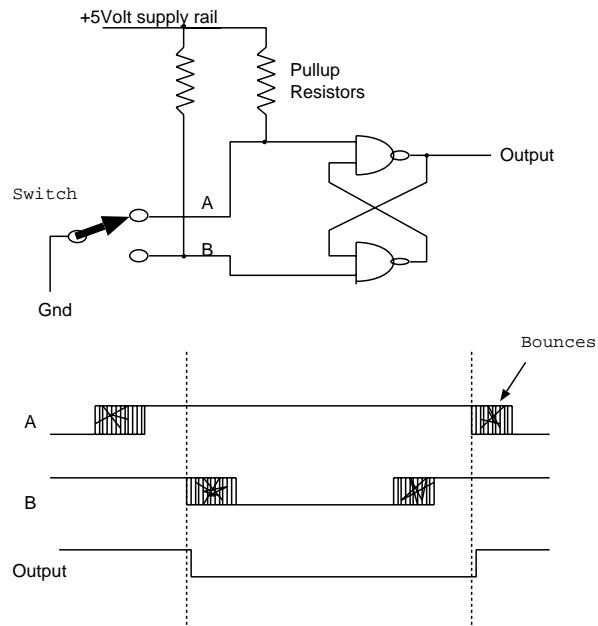


Figure 14: A debouncer circuit for a two-pole switch.

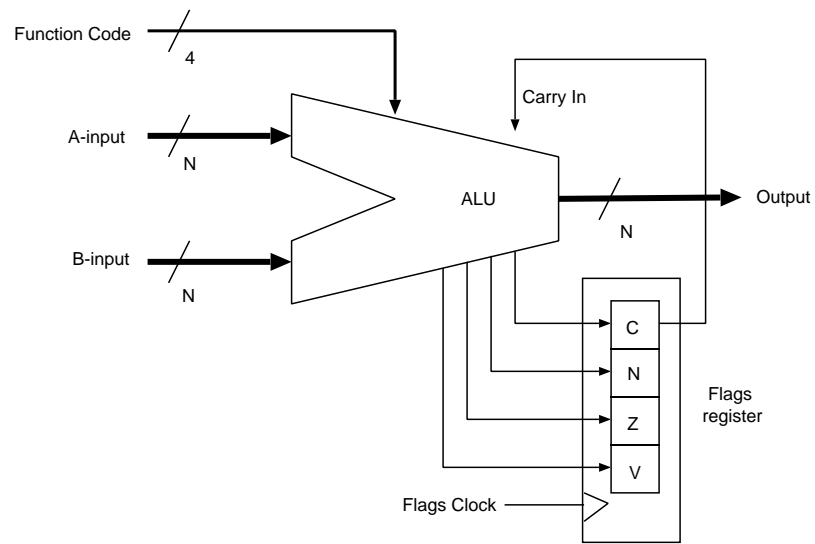


Figure 15: Schematic symbol for an ALU, connected to a flags register.

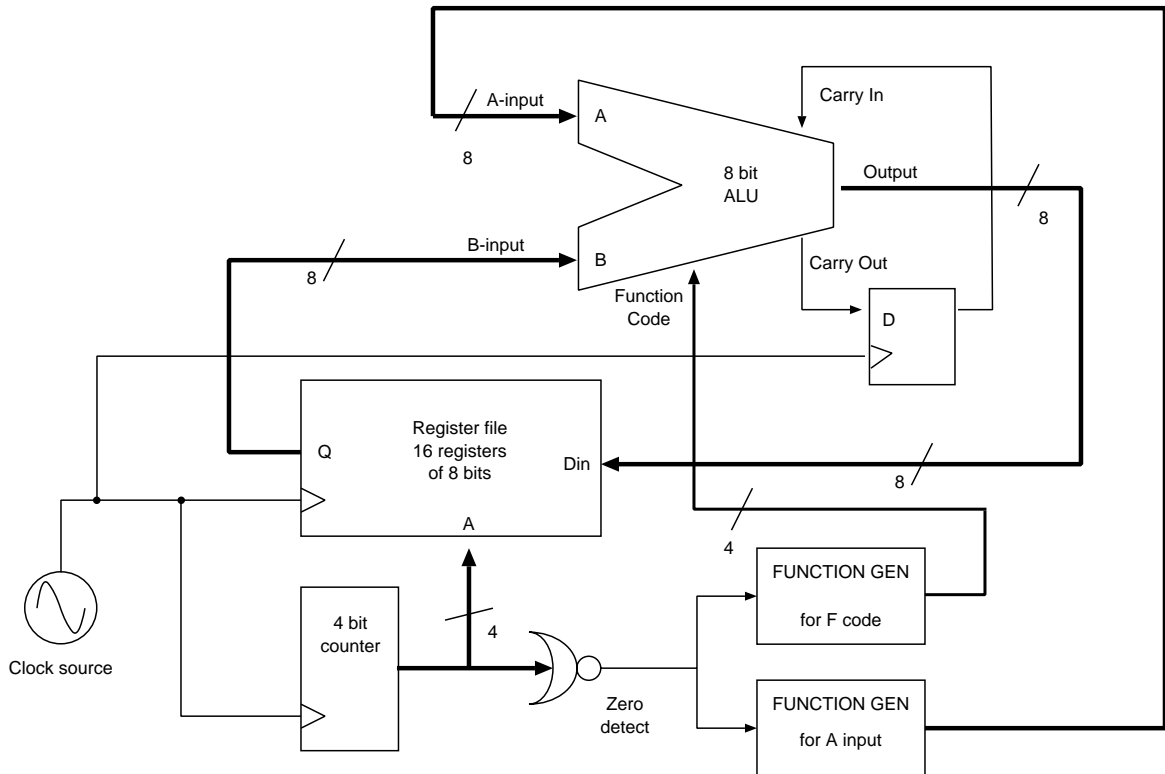


Figure 16: Example structure using an ALU and register file.

Figure 16 is a nonsense circuit designed to show the interaction of an ALU and a register file. (These two together form the heart of the execution unit of a computer.) In the circuit, the addressed register is both read and written at the same time. Even though the address is incremented each time, it is important to understand that the old value read is updated in the ALU and written back to the same location. The only possible interference between one register and the next is through the carry flip-flop, which will allow the output of one ALU operation to affect the behaviour of the next.

Two function generates are illustrated. These are just combinatorial blocks and could be considered as ROMs. To achieve some particular function these could be hardwired or programmed. In very early computers plugboards were used for such functions.

If the ALU function code generator is programmed to specify 'A+B' when the input is one and 'B+Cin' otherwise, and the A input function generator is programmed to produce 0 when the input is one and 0 otherwise, we will have programmed up a very large counter. Will it complete this universe lifetime ?

1.6 Microprocessor Circuits

Figure 17 shows the schematic symbol and internal block diagram for a microprocessor. The internal block diagram is beyond the scope of this course, since it is coupled with a subject known as programming, but the external view is fair game.

The main features of the microprocessor are that it has a clock and reset input and it has an address bus output and a bi-directional data bus. The microprocessor generates read and write cycles on its busses and these are essentially identical in form to the waves shown in figure 6 for the SRAM.

The microprocessor places addresses on the address bus and sets the read/write signal to show the direction of the intended transfer on the data bus. It then raises its operation request output (opreq) and either places the data to be written on the data bus (for a write) or expects the

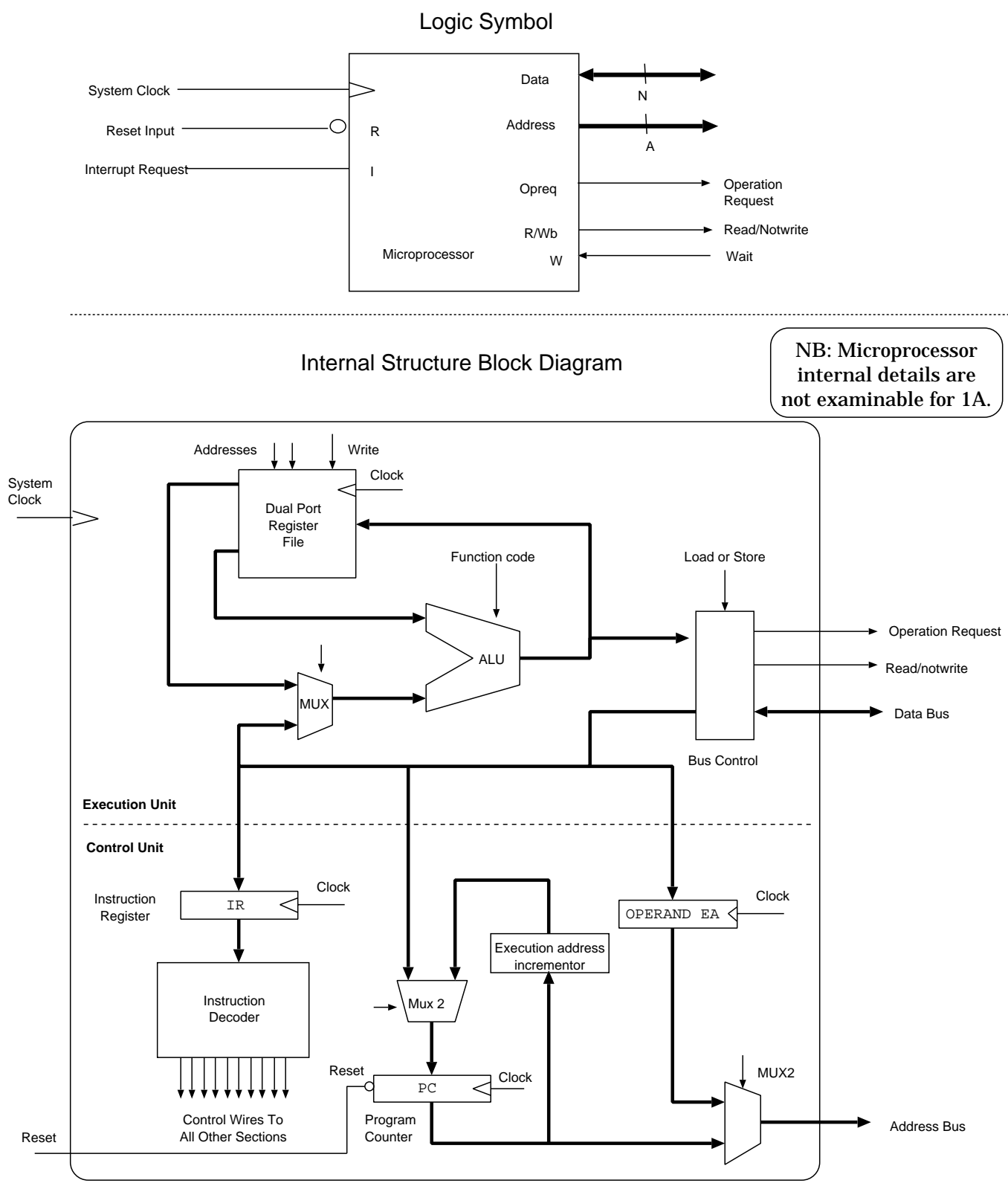


Figure 17: A microprocessor logic symbol and simplified internal structure.

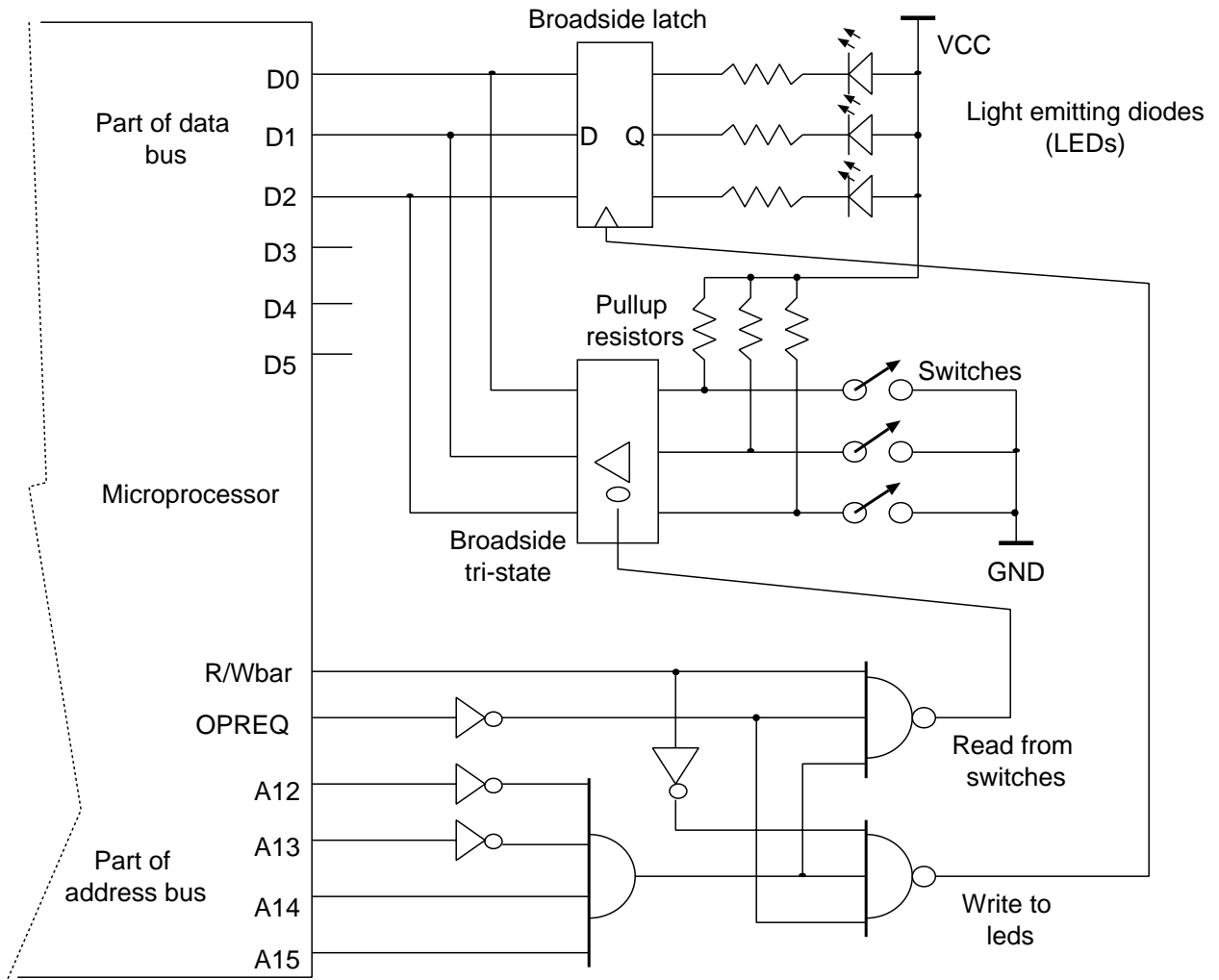


Figure 18: Example of memory address decode and simple LED and switch interfacing for programmed IO (PIO) to a microprocessor.

addressed device to drive the data bus with its information (for a read).

Figure 18 shows how to connect some LEDs and switches to a microprocessor for simple programmed input and output. Only the low order data bits are used, so when reading from the switches, random values may appear on the higher order bits. Conversely, when writing to the LED registers, the higher order bits are ignored.

Exercise: If the address bus is 16 bits, how many addresses are consumed by the switches and leds. If the data bus is 8 bits, what percentage of addressable bits have been wasted and is this a sensible design ?

Figure 19 is the circuit of small, microprocessor based computer. The microprocessor can address 2^{16} locations of eight bits (65536 bytes) and some of these are filled in with devices.

Exercise: What would happen if the microprocessor tried to write to the read only memory ?

Exercise: Figure 20 shows the signals involved in a parallel port, as frequently used to connect printers to computers. The basic operation is that the master places a byte of data on the data bus and then strobes the strobe signal low until acknowledge is asserted by the device. The master must not present a new byte to the printer while the printer has put busy high. Design the parallel port interface logic by drawing on the ideas shown in figure 18. To do this, it is sufficient if the busy signal appears to the processor like one of the switch inputs and the data output is rather like eight LEDs. Considering that the processor really does not care about the acknowledge signal

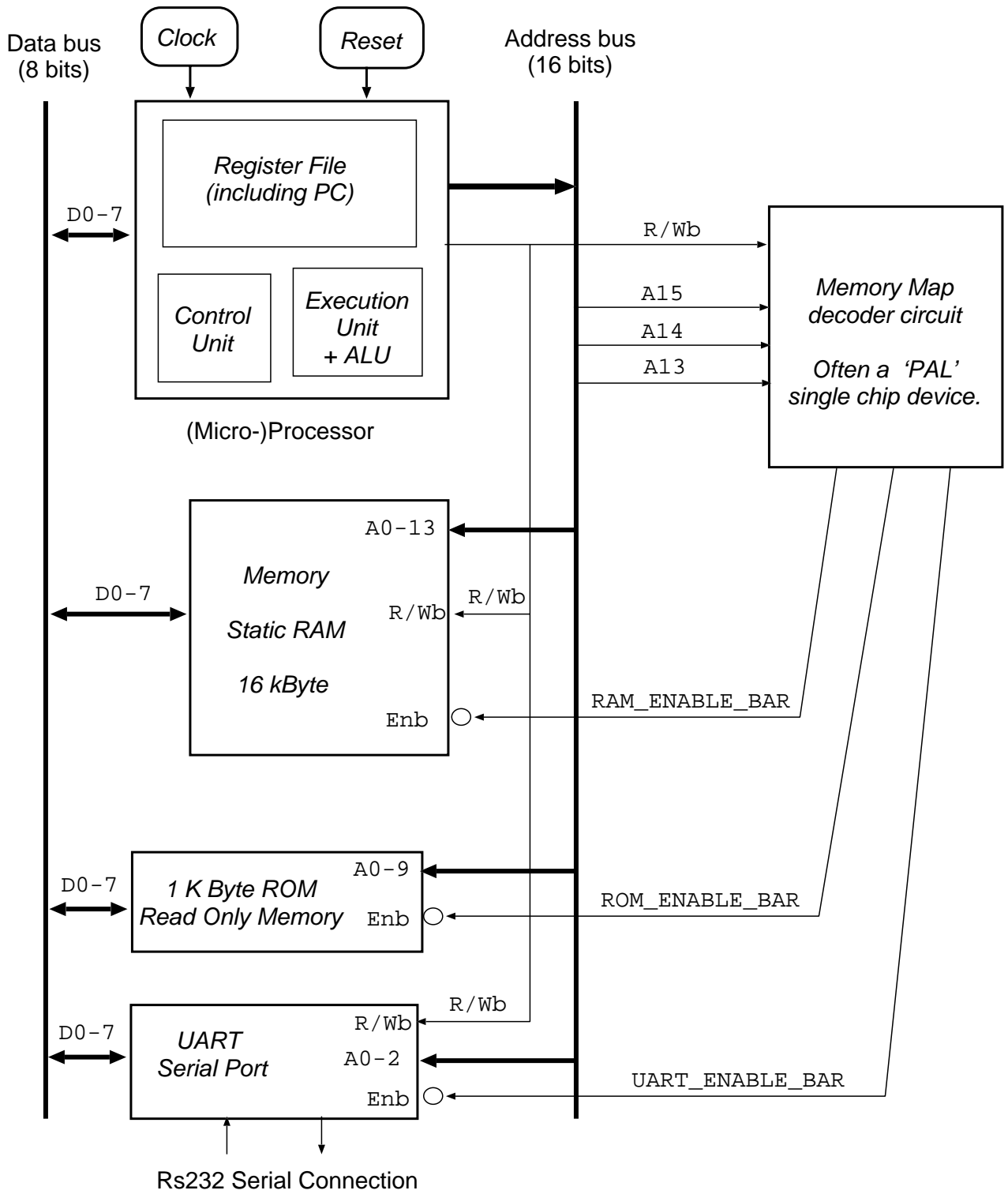


Figure 19: A small computer.

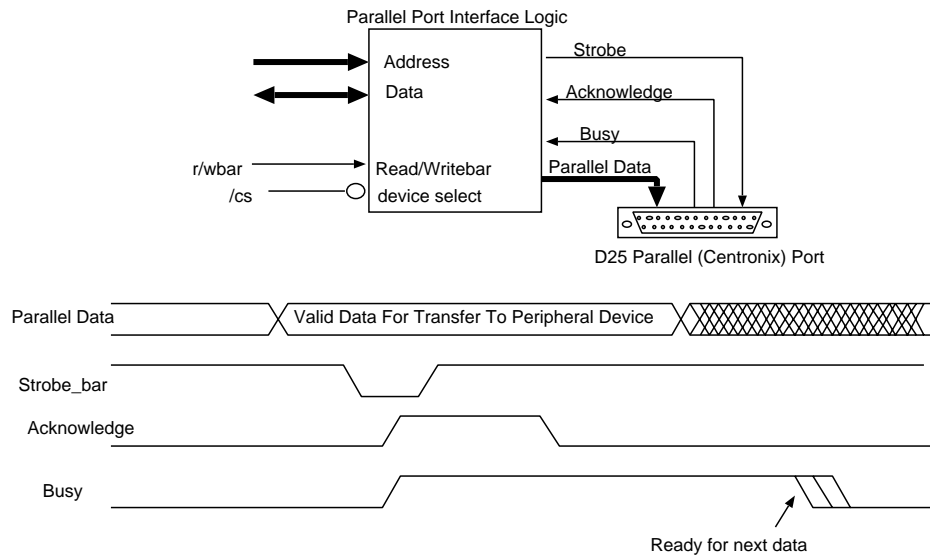


Figure 20: PC or Centronics Parallel Printer Port.

or when the strobe signal is deasserted, you might care to use an RS latch to help the processor generate the strobe signal.

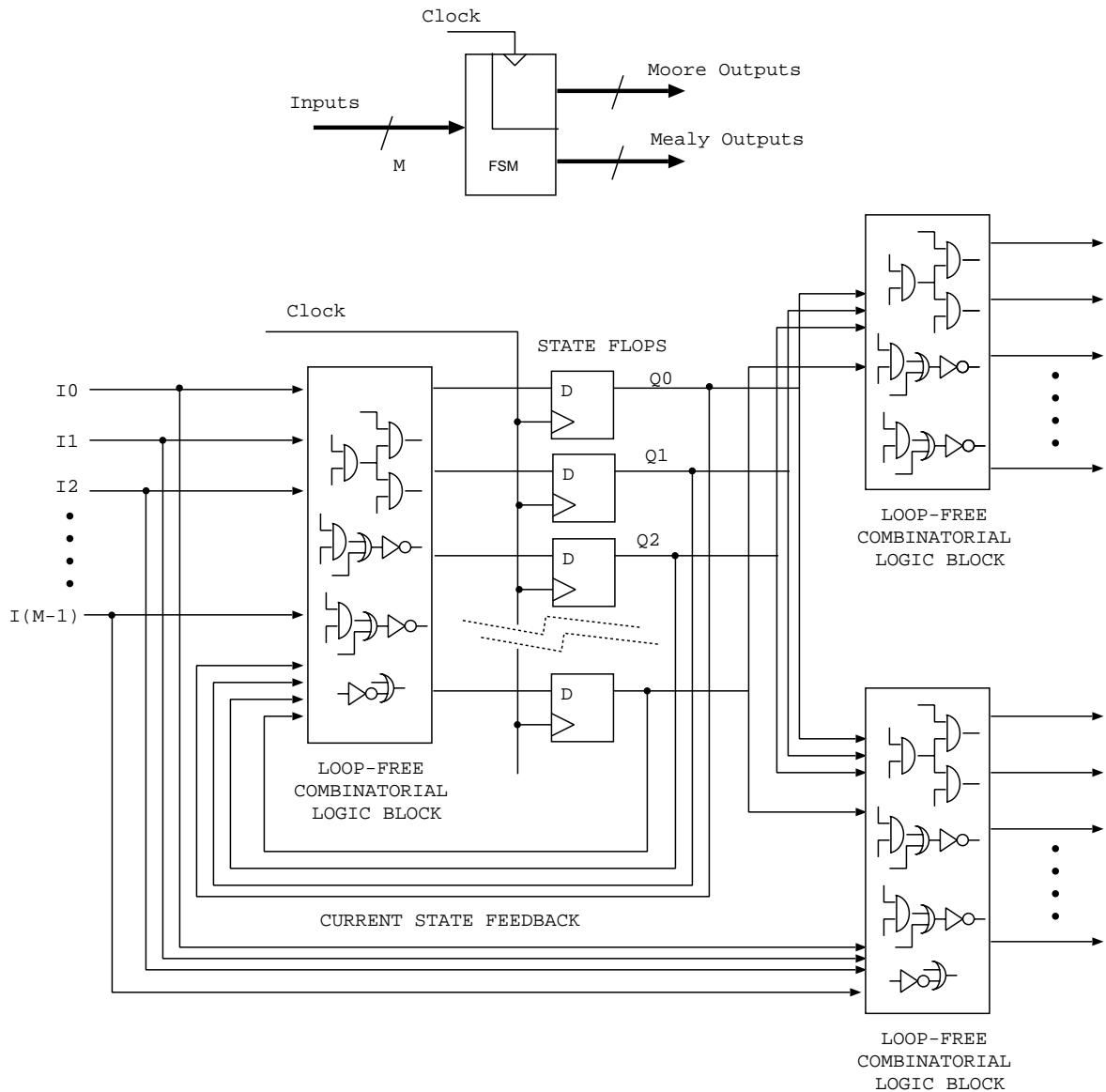


Figure 21: Canonical synchronous, finite-state machine.

2 Clocking and Synchronous FSM Combination.

This section looks at aspects of clock signals and considers the how, when and why of multiple clocks in a system. The decisions relating to clocks are fundamental to any hardware design and are normally taken very early in the process.

2.1 Finite State Machines.

Figure 21 shows the schematic symbol and canonical circuit for a finite-state machine. It is a mantra among mature hardware designers to consider everything as a collection of finite-state machines. Beginner hardware designers tend not to have this approach, and instead think nothing of treating the clock input to a flip-flop as suitable for wiring to anything convenient. They do not get far like this, since all CAD tools, design techniques and formal methods will stop providing any help and circuit bugs will creep up everywhere. (An exception is certain *licensed* asynchronous hardware designers, like Dr Simon Moore, who often don't even have clocks in their circuits.)

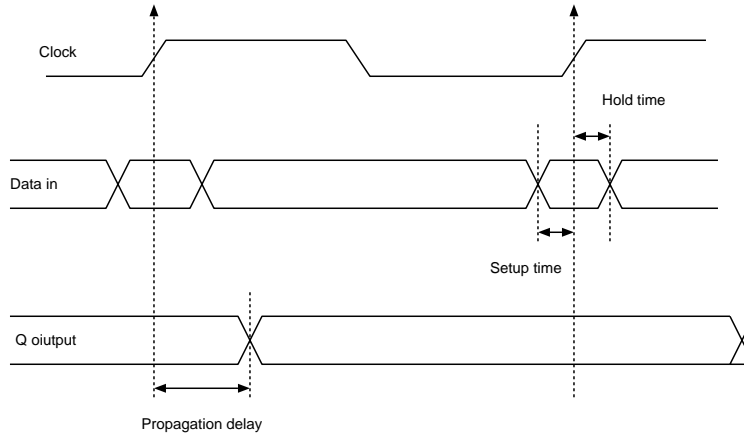
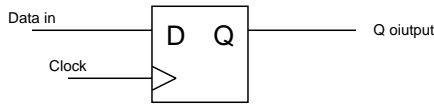


Figure 22: Timing specifications for a flip-flop.

The FSM has inputs and a clock. It has two types of outputs. Moore outputs depend only on the current state. Mealy outputs may also depend on the current inputs.

Figure 22 shows a D-type flip-flop and its three significant timing specifications. The flip-flop input must not change at the same time as the positive edge of the clock: indeed it must be stable for a period before (the set-up time) and remain stable a period after (the hold time). The output will change after a propagation delay time from the positive edge of the clock. The negative edge of the clock has no effect. The J and K inputs of JKs and the T inputs of Ts have the same set and hold requirements.

In a synchronous FSM, the maximum clock frequency is normally determined by the longest path of logic which ends at the input of a flip-flop. This path is known as the *critical path* and is illustrated in figure 23.

2.2 Speed Increase with Johnson Counter

A *Johnson* counter is based around a shift register. Figure 24 shows a Johnson counter that divides by five. The significant features of Johnson counters actually become clearer for larger counters, but they are twofold: 1. the division ratio is restricted to twice the number of flip-flops (e.g. 5 flip-flops can give a divide by 10 at maximum) and 2. the number of logic gates is always low, and so the counter is fast.

Exercise: Design a Johnson Counter that also has a control input that makes it divide by six or seven. Optimise your design for two features: 1. maximum clock speed, 2. that the ratio input can be changed at any time without causing a hazard. Note: clearly, metastability cannot be avoided, but by being hazard free, we mean that the state path taken will always be one or the other of the two desired paths, whatever happens. Why might the second feature be important in practice (hint: consider derived clocks) ?

2.3 Speed Increase with One Hot Coding

In *One Hot* coding, at all times, the flip-flops in a counter or FSM are arranged to be all zero except for one that is one. This clearly minimises the number of gates required to decode the current state and can greatly reduce the number of gates that determine the next state.

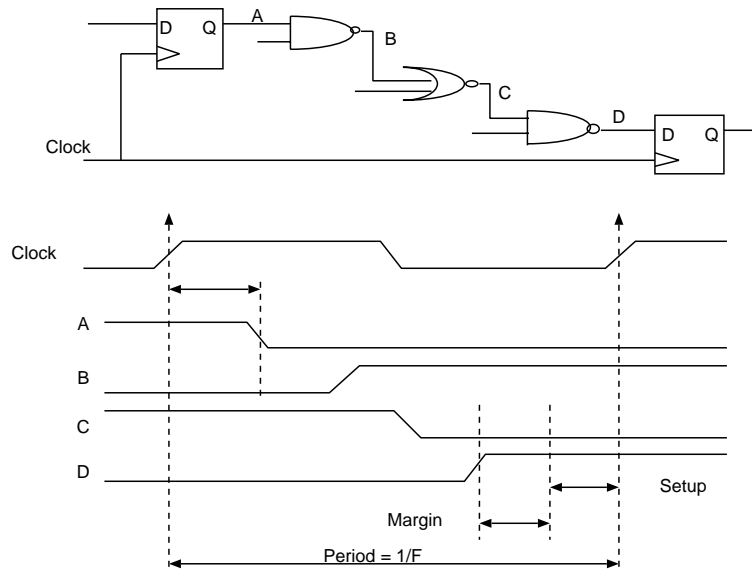


Figure 23: Typical nature of a critical path.

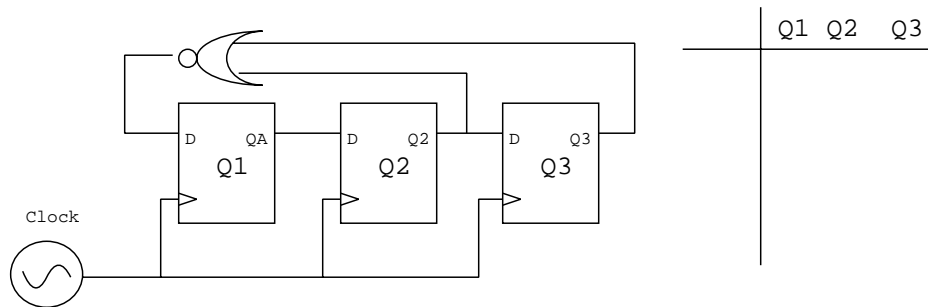


Figure 24: A counter that does not use binary.

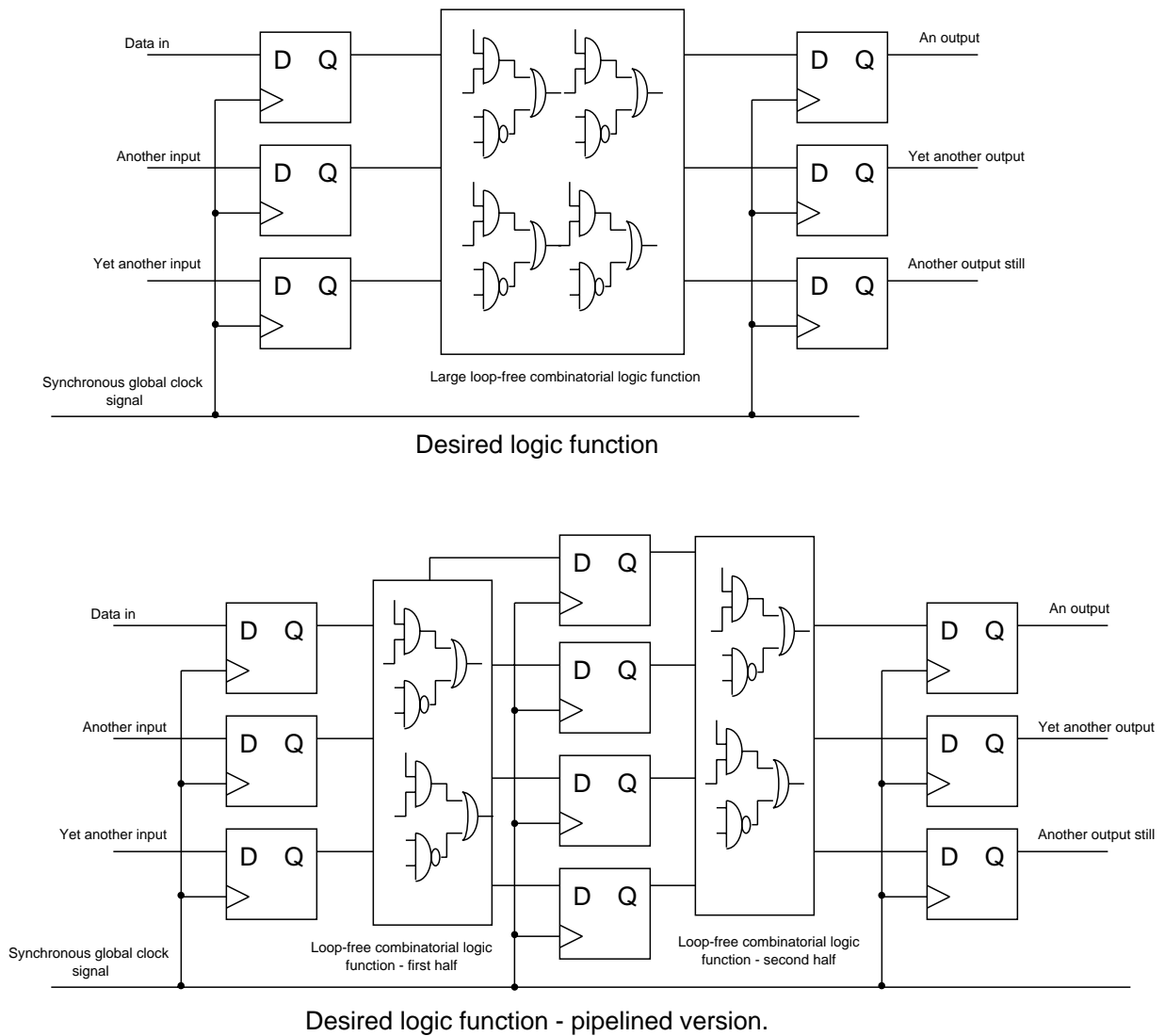


Figure 25: Example of introducing pipelining to increase throughput, but also delay.

2.4 Speed Increase using Pipelining

Figure 25 shows a circuit before and after the addition of a pipelining stage. The pipeline stage may have a fewer or greater number of flops as the input or output. The benefit of the pipeline stage is that a complex combinational logic function is split into two halves, thereby reducing the maximum length of unlocked logic leading to the inputs of the output flip-flops. The disadvantage of the pipeline stage is that data takes longer to be processed. This disadvantage is not often a problem: for instance, if the data has just come down a link in a network, then the additional delay will be just as if the link were a few meters longer. More than one stage of pipelining can be inserted if needed.

Some modern logic synthesiser tools will insert pipeline delays where possible, if helpful. Others will back-propagate them if they are instantiated at the outputs.

Exercise Sketch the circuit of a large adder with ripple carry and then modify it to include a pipeline stage. Note, the output signals should always reflect the result of a single addition in any one clock cycle. Explain or estimate the effect of the pipeline delay on the maximum clock frequency. If the adder is now used to form a counter where an input value N is added to an accumulator register each clock cycle *a)* what goes wrong, and *b)* how can this be fixed ?

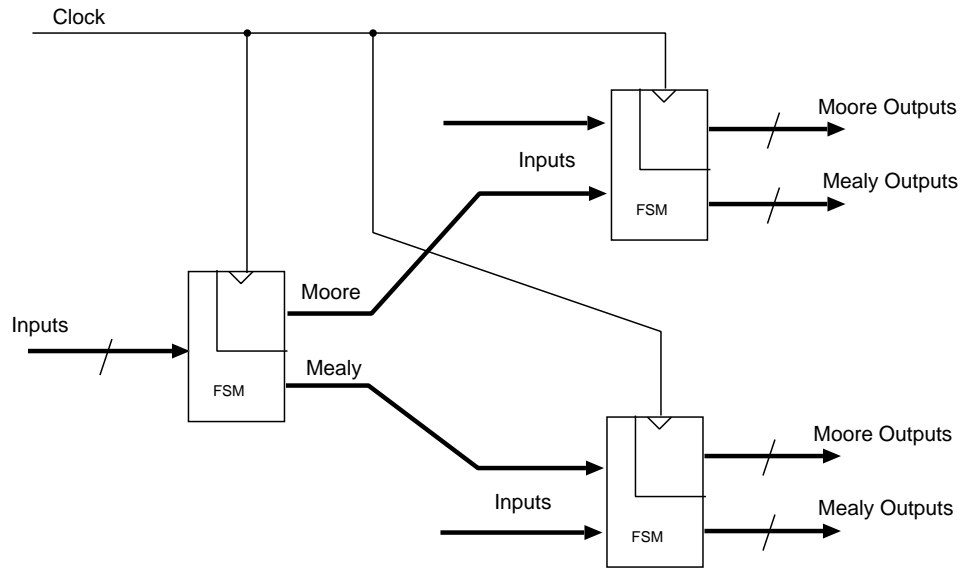


Figure 26: Cascading FSMs to produce larger circuits.

2.5 Derived Clocks.

Figure 26 shows a concatenation of finite-state machines running from the same clock. Of course, canonically, these together simply form one larger FSM. However, the decomposed view is vital in practice, since the density of wiring between FSMs will be lower than within the FSM and the functionality, authorship, history and even ownership of the IPR (intellectual property rights) of the FSMs will be different.

If we feed a Mealy output into another FSM, delays accumulate, reducing maximum clock frequency. This does not happen with Moore outputs. Moore outputs should be used if possible, but for some designs, a Mealy output is needed in order to reduce processing delay of the unit as a whole.

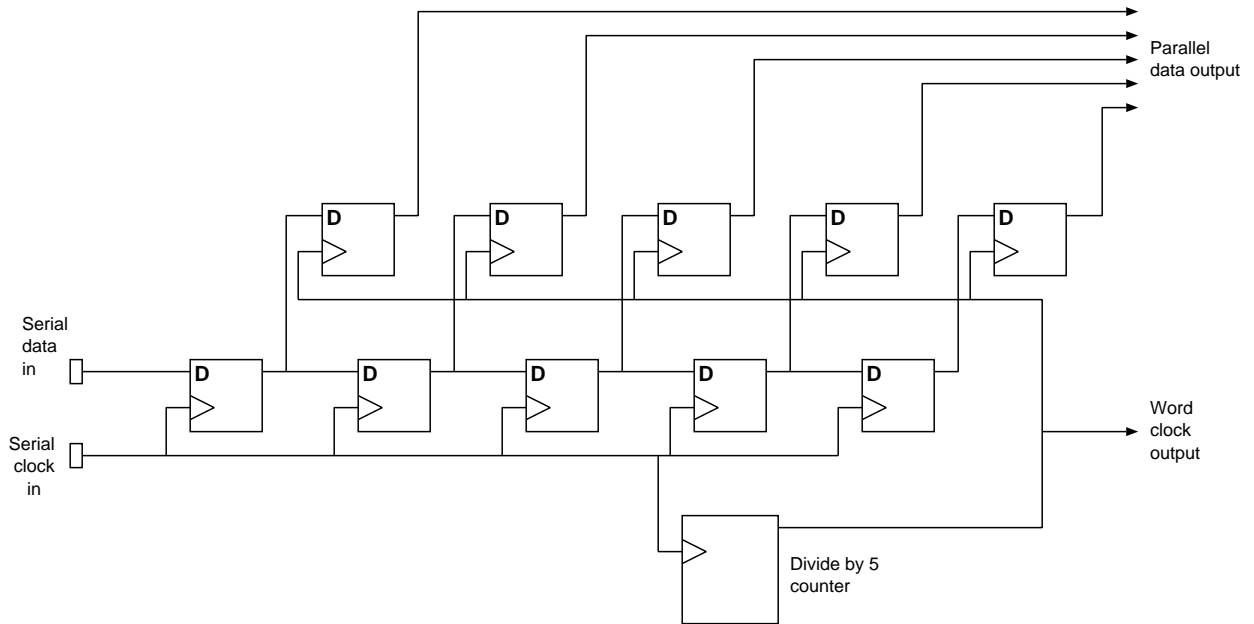


Figure 27: An example of a derived clock: a serial-to-parallel convertor.

Figure 27 shows a typical example of the use of a derived clock. Power consumption is increased in parts of a circuit where fast clocks are used, and so using lower frequency clocks where needed is an important design consideration.

The serial-to-parallel convertor can be redrawn as a pair of FSMs as shown in general in figure 32. It is then clear that the design rule about two Moore outputs from the master not changing at once cannot be applied. Therefore what can we do ? This is not easy. Real hardware often is not easy.

Exercise: Give the schematic for a reciprocal parallel-to-serial convertor. Does combining the two convertors result in the identity function ? If not, why not. Note: they can be combined in both orders.

2.6 Gated and guarded clocks

Sometimes we do not wish a flip-flop to change its value every clock cycle. Clearly, J-K device support this by putting both inputs low. For subsystem FSMs or broadside D-type registers, we need the concept of the clock enable.

Figure 28 shows a D-type with a clock enable input and the equivalent circuit. The complexity of the additional multiplexor is not present in a real implementation, since it can be adsorbed into the circuit which makes the D-type.

Figure 29 shows the safe way to gate a clock. The clock is OR-ed with the negated enable signal. In this clock gating method, the derived clock is kept normally high, except when a clock pulse is to be let through. Other arrangements for gating clocks tend to produce glitches in the derived clock.

Gating clocks should be avoided normally. If flip-flops have clock enable inputs these should be used instead. Providing your own multiplexors before the fops is also a good alternative. Clocks do have to be gated at times though, most commonly in generating a write strobe to an SRAM (figure 6).

Exercise: Design the logic to place in front of a J-K flip flop that turns it into a D-type with clock enable input.

Given that from Digital Electronics you know that a J-K can be equally as simple as a D-type to implement from transistors or gates, it is reasonable to assume that a clock enable D-type is also as

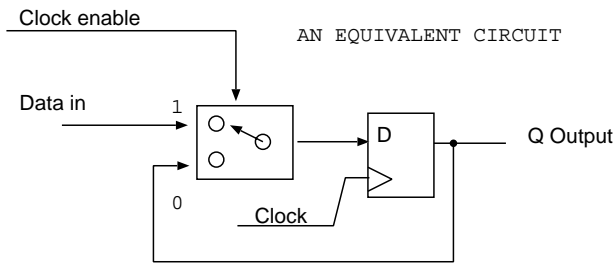
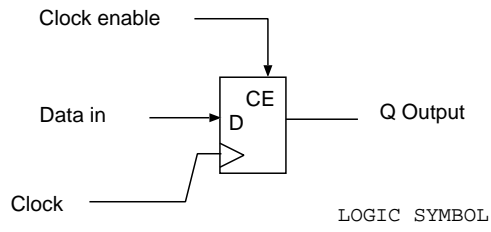


Figure 28: A D-type with clock enable and the equivalent circuit.

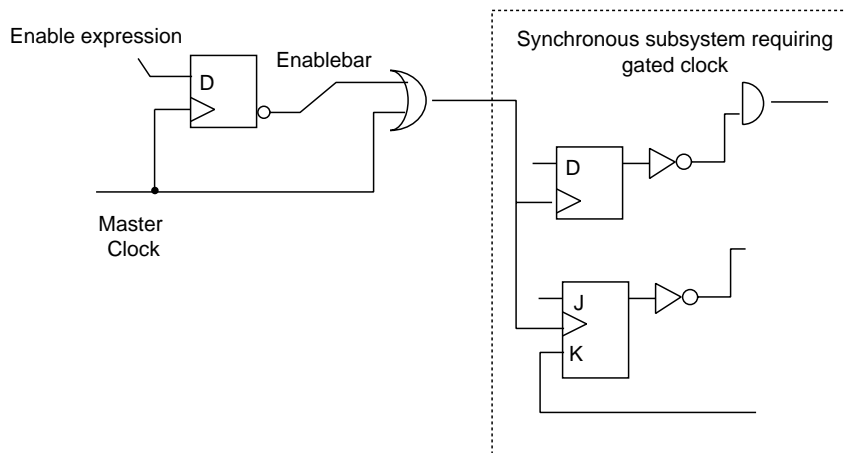
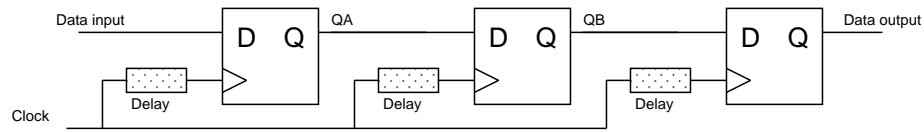
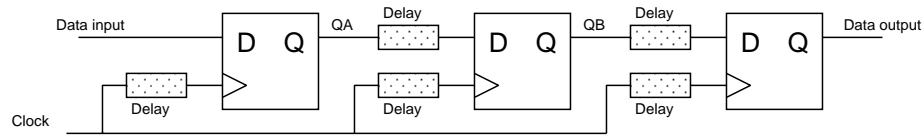


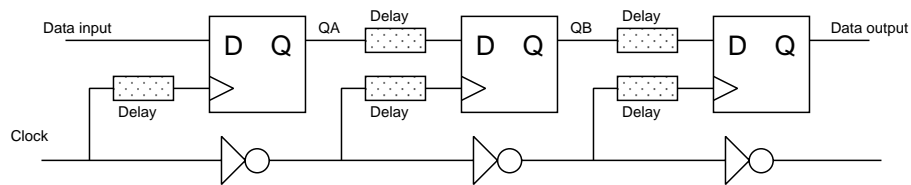
Figure 29: A 'safe' alternative to using clock enables ?



a) A three-stage shift register with some clock skew delays.



b) A three-stage shift register with data delays and clock skew delays.



c) A solution for serious skew and delay problems ?

Figure 30: Synchronous systems, but with clock skew.

simple. Therefore why should we ever use gated clocks instead of clock enable inputs ? The answer is that if we have a large number of flip-flops that are all to share the same clock-enable expression, we would have two penalties of using clock enable: 1. we have to route the clock enable signal to all of them, and 2. we have to make the flip-flops from logic which is capable of being clocked at every opportunity. We use clock-enable inputs where single or smaller numbers of flip-flops are to be enabled in unison, such as in a broadside register with clock enable. For single flip-flops, the clock enable is best considered as just one aspect of the FSM that the flop is part of.

2.7 Clock and Data Skews

Skew is when a signal that is intended to arrive at several places at the same time actually arrives at various times.

Figure 30 shows a shift register circuit with clock skew. Skew is the enemy of synchronous, finite-state machines, since in the worst scenario, the output of one flop may have changed, causing a change in the input to another flop, before that second flop has been clocked itself: this destroys the FSM behaviour.

If there is considerable delay in both the data and the clock, and these delays are random, then the shift register is not likely to work well. An example of such a system is where the shift register is actually a communications medium snaking around a number of circuit boards in a system. The delays are in the interconnecting cables. Typically, the logic at each station is not simply a D-type, but is logic which can read and change the data as it passes and the data may be a bus instead of a single wire. In this case, the solution is to invert the clock at each hop.

Exercise: Draw a timing diagram for the inverting-the-clock solution and work out the maximum frequency of operation. Use the following figures: conductor delay random between 4 and 9 nanoseconds, flip-flop set-up time 2 nanoseconds, hold time 1 nanoseconds and propagation time 5 nanoseconds. Note: in practice you would not run such a system close to the maximum you have

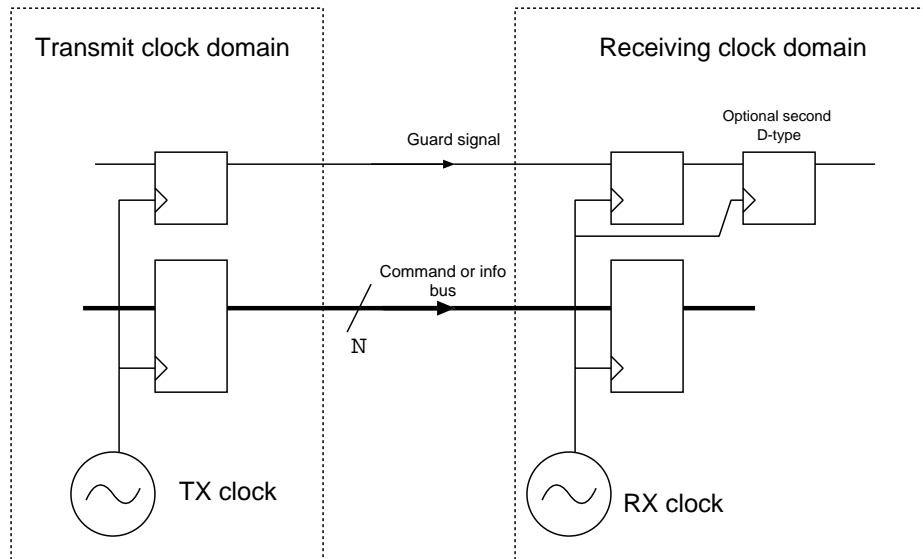


Figure 31: Use of a guard to cross an asynchronous boundary

just calculated.

2.8 Crossing Asynchronous Boundaries

Figure 31 shows two clock domains with signals crossing between them. There may be signals in the reverse direction between the clock domains, but these would not demonstrate anything further that is significant.

The first thing to note is that for a single signal between the domains, there is always the possibility of metastability and failure of the whole system. However, with careful flip-flop design, a designer can reduce this to less than the probability of winning the lottery every week for a year: so this is not a real concern. Since there is more than one clock in the world, we do have to build these circuits in reality.

The main thing to note is that there will be skew between the lines and in the characteristics of the individual flip-flops. Therefore it is impossible to be sure of capturing a consistent view of them at the receiver if they are all changed at once or changed freely by the source. Therefore a protocol is required. A suitable protocol is to denote one of the signals as a *guard* or *qualifier*. Only if the guard has its active value is it implied that the data on the remaining signals is valid. The source must execute or implement a protocol that simply delays the assertion of the guard with respect to the other wires such that the other wires are sure to be seen as settled if the guard is seen to be true.

A good way to help avoid metastability is to use two flip-flops in tandem (as shown for the guard). Even if the first oscillates a while after being violated, it is most unlikely that it is still oscillating a clock cycle later when the second flip-flop is used. The additional delay in the second flip-flop will result in the guard being presented to the receiver logic later than the signal lines, thus ensuring that they are stable.

Exercise: Is it possible to use changes of state of a guard signal to indicate that the data on the remaining wires is valid? If so, sketch an example circuit and explain the benefits. What assumptions regarding the relative clock frequencies in use have you made, if any?

2.9 FSM clocks derived from another FSM

Figure 32 shows a FSM whose outputs are used to clock another FSM. Examples of this are common, as in the serial to parallel convertor shown in figure 27 and the CFR two-chip decomposition

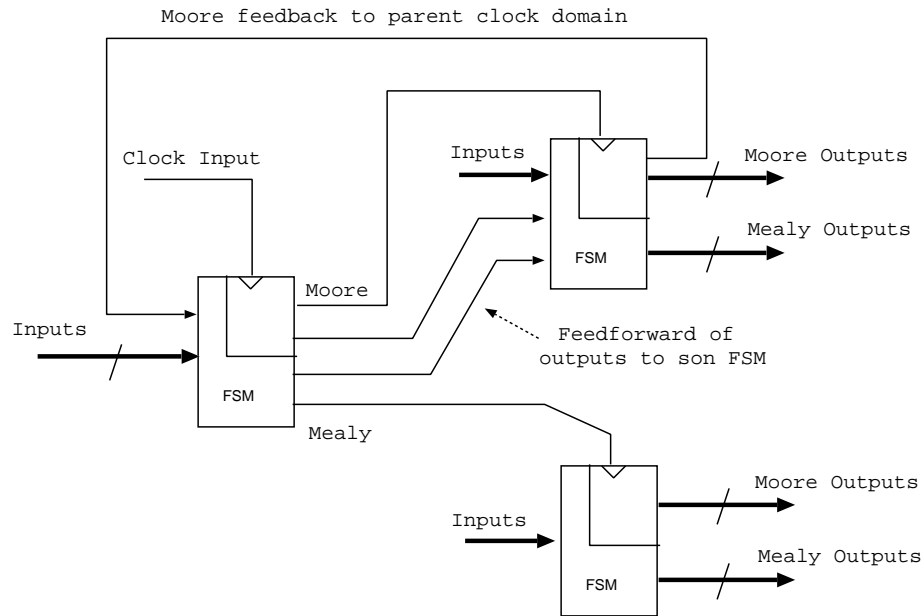


Figure 32: Example paths between FSMs with derived clocks

shown in section 4.4. Such structures must be used carefully. The Mealy output clocking a slave FSM is probably a bad idea. The Moore output used to clock a slave is normally fine, but the logic function generating it must be free from hazards, otherwise the clock will glitch. A function of one state variable is best, and state assignment should have considered that.

The non-clock signals between the master and slave FSM also need attention. The Moore output from the slave to the master will not cause hazards, but will tend to restrict the maximum frequency of safe clocking for the system as a whole. The Moore output from the master to the slave needs careful attention to hazards, since it must never change at the same time as the slave clock signal.

Exercise (long): Consider an FSM or system which generates a set of output signals that have to be carried via an intermediate subsystem which uses a separate (asynchronous), faster clock, and then delivered to a receiving system that must be clocked with a transported version of the original clock. An example of this would be where a company has fire alarm systems on two sites which are interconnected on a high speed data link owned by a telephone company, where the link is also carrying many other services. Sketch the basics of the circuits involved.

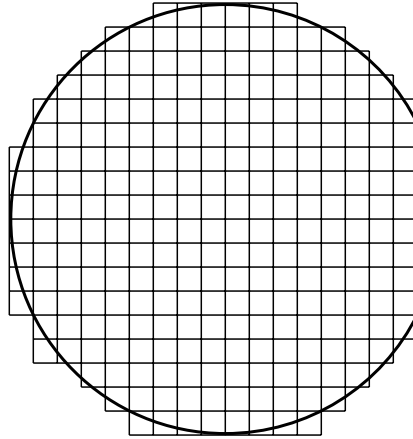


Figure 33: A wafer outline showing where it will be diced (Chips are not always square).

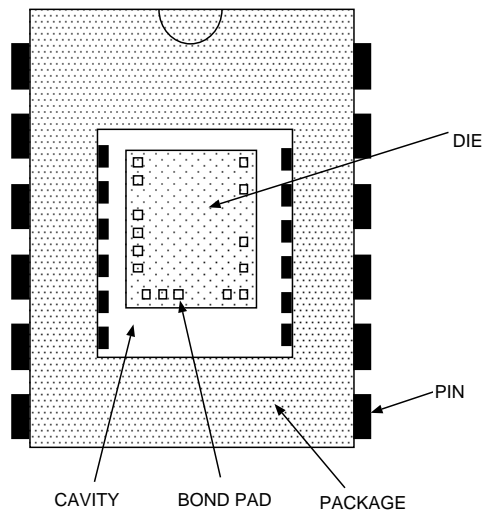


Figure 34: A chip in its package, ready for bond wires.

3 Technology

The cost of a system depends on the cost of its parts. There are many ways to design a system and typically we want to minimise cost. Therefore, before you can design a system, you need to be familiar with a vocabulary of technologies and pre-existing parts. You also need to know how to produce new integrated circuits for your application design. Application-specific integrated circuits (ASICs) and make up a large fraction of the IC market by revenue.

3.1 Basics of ICs.

Note: this section is a briefest of introductions to VLSI.

Chips are made of Si or GaAs with layers of Al on top for signal wiring. Chips are made on a rectangular grid on a wafer of the substrate material. Wafers are about 0.3 mm thick and from 3 to 8 inches in diameter. Most of the wafer depth is only present for handling reasons and the active portion of the chip lies in the top few microns of the wafer or is built up above the top surface.

Wafers are processed using ion implantation and photo-resistive etching at elevated temperatures in the presence of dopant and etchant gasses. Each photo-resist is first exposed to an ultraviolet or X-ray source under a mask generated by the CAD tools. Between 7 and 25 masks might be used, depending on the technology. The masks are stepped over the wafer to form the rectangular

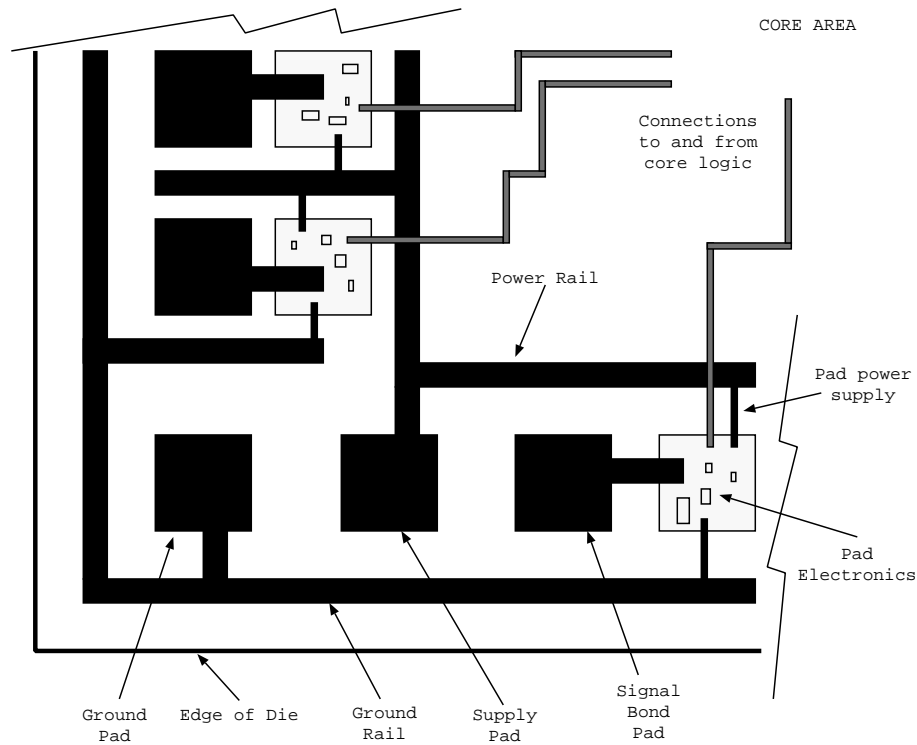


Figure 35: The corner of a typical chip showing IO and power pads.

array of chips. Masked areas are left with an intact layer of resist and then remain unaffected by the next processing step, hence successive masks build up complex three dimensional structures forming the active components and conductors. Typical chips are today made using 0.25 to 0.5 micron feature sizes. The smaller a transistor the faster it operates.

A wafer contains many chips which are diced after wafer testing. The chip is then packaged to ease handling and to dissipate heat.

A typical chip is 4 to 100 square millimetres.

Wires are connected to chips at bond pads. Bond pads are about 100 microns square and 150 microns apart or half this size (linear dimension) in modern devices. The corner of a typical device is shown in Figure 35. Most chips have the shown power and ground ‘rings’ around the outside with the bond pads in between. These rings provide low resistance power supplies to the pad logic. The pad logic is buffer circuitry to provide high-power drive to off-chip loads and to protect the internal circuitry from external static discharges.

Figure 34 shows the chip carrier cavity containing a die. Connections are made from the conductors in the package to the bond pads around the chip’s perimeter using cold-welded gold wires.

3.2 Will an ASIC be core or pad bound ?

A chip may be either:

- **Pad bound** if it has many more pins than natural perimeter, or
- **Core bound** if it has sufficient active core area to dominate.

A pad bound chip is wasteful of area since we have to use a bigger die to get the pads on, but the body of the die is not fully occupied with gates.

Generally we end up with separate fabrication processes for each of the following subsystems:

- memory,

- random logic,
- analogue circuits and analogue to/from digital conversion,
- high speed (e.g. clock rate greater than about 100 MHz).

Also, we must group together functions which can be implemented in the same technology, so they can most easily be put on one chip. This is hard if the functions do not otherwise have much to do with each other: we might prefer to put the two sections of logic further apart.

3.3 Chip cost versus area

The cost of a chip divides into two parts: NRE and per-device cost.

NRE — non-recurring expenditure. NRE is paid once, regardless of how many devices are made.

It consists of:

- design cost (labour, computer time, management overheads),
- mask making cost.

Per device — recurring expenditure consists of:

- cost of basic silicon wafer,
- cost of processing a wafer,
- cost of testing devices,
- cost of device packages and packaging,
- cost of further testing, possibly under stressful conditions.

The per device cost is influenced by the yield — the fraction of working dice.

The fraction of wafers where at least some of the die work is the ‘wafer yield’ and is typically close to 100 percent for a mature fabrication process.

The fraction of die which work on a wafer (often simply the ‘yield’) depends on wafer impurity density and device size. Die yield goes down with chip area.

The fraction of devices which pass wafer probe (i.e. before the wafer is diced) and fail post packaging tests is very low. However, full testing of analog sections or other lengthy operations are typically skipped at the wafer probe stage.

3.4 Example of CMOS die cost

A six inch diameter wafer has area $(3.14r^2) = 18000 \text{ mm}^2$.

A chip has area A , which can be anything between 2 to 200 mm^2 (including scoring lines).

Dies per wafer is $18000/A$

Processed wafer cost might be 5000 pounds. (Actually it’s probably about one third this cost to the silicon foundry, but there are overheads and profit to add on.)

Probability of working = wafer yield \times die yield (assume wafer yield is 1.0 or else included in the wafer cost).

Assume 99.5 percent of square millimetres are defect free. Die yield is then

$$P(\text{All } A \text{ squares work}) = 0.995^A$$

cost of working dice is

$$\frac{5000}{\frac{18000}{A} 0.995^A} \text{ pounds each.}$$

Here is a program to estimate the cost of a working die given a six inch wafer with a processing cost of 5000 pounds and a probability of a square millimetre being defect free of 99.55 percent. Its output is shown in table 1.

Area	Wafer dies	Working dies	Cost per working die
2	9000	8910	0.56
3	6000	5910	0.85
4	4500	4411	1.13
6	3000	2911	1.72
9	2000	1912	2.62
13	1385	1297	3.85
19	947	861	5.81
28	643	559	8.95
42	429	347	14.40
63	286	208	24.00
94	191	120	41.83
141	128	63	79.41
211	85	30	168.78
316	57	12	427.85
474	38	4	1416.89

Table 1: Output from the die cost program.

```

#include <math.h>
display(double area)
{
    double dies = 18000.0 / area;
    double yield = dies * pow(0.995, area);
    double cost = 5000.0 / yield;
    printf("%.0f %.0f %.0f %.2f\n", area, dies, yield, cost);
}
int main()
{
    double area = 2.0; /* area in square mm */
    while (area < 500.0)
    {
        display(area);
        area = trunc(area * 1.5);
    }
}

```

Large dies sometimes use redundant design such that one, two or three faults can be tolerated. Laser or ion beam trimming is used to patch out faulty sections (of say a RAM array) and substitute previously redundant sections.

Exercise: Given that the cost of patching a die doubles its processing cost and is 90 percent successful and requires 10 percent additional area in the design, for what size and type of device is it suitable? *It is necessary to make various assumptions to answer this.*

3.5 Classes of Integrated Circuits

Figure 36 presents a taxonomy of chip design approaches. The top-level division is between standard parts, ASICs and field programmable parts.

3.5.1 Standard Parts

A *standard part* is essentially any chip that a chip manufacturer is prepared to sell to someone else along with a datasheet. The design may actually previously have been an ASIC for a specific customer which is now on general release. However, most standard parts are general purpose logic, memory and microprocessor devices. These are normally full custom designs designed in-house by the chip manufacturer to make the most of in-house fabrication line, perhaps using optimisations not made available to others who use the line as a foundry. Other standard parts include graphics controllers, LAN controllers, bus interface devices, and miscellaneous useful chips.

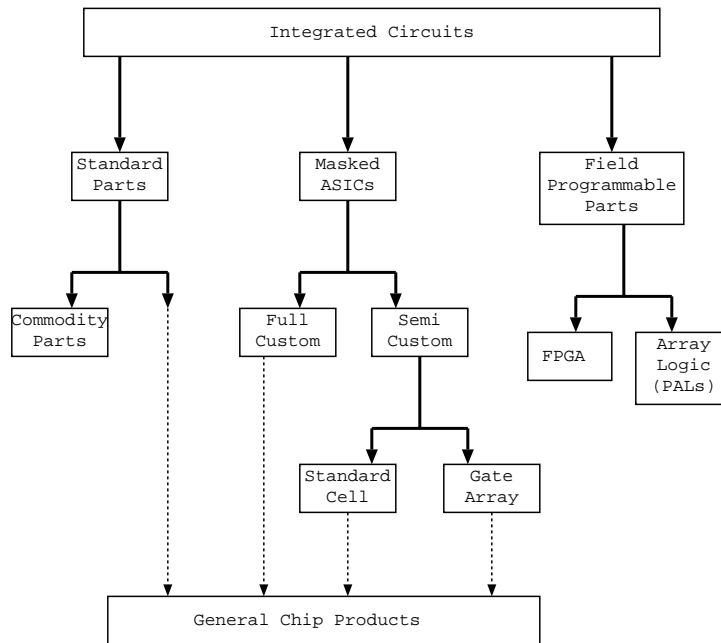


Figure 36: A taxonomy of integrated circuits.

3.5.2 Masked ASICs.

A masked ASIC (application specific integrated circuit) is a device manufactured for a customer involving a set of masks where at least some of the masks are used only for that device. These devices include full-custom and semi-custom ASICs.

A full-custom chip (or part of a chip) has had detailed manual design effort expended on its circuits and the position of each transistor and section of interconnect. This allows an optimum of speed and density and power consumption.

Full-custom design is used for devices which will be produced in very large quantities: e.g. millions of parts where the design cost is justified. Full-custom design is also used when required for performance reasons. Microprocessors, memories and digital signal processing devices are primary users of full-custom design.

In semi-custom design, each cell has a fixed design and is repeated each time it is used, both within a chip and across many devices which have used the library. This simplifies design, but drive power of the cell is not optimised for each instance.

Semi-custom is achieved using a library of logic cells. Figure 45 shows a cell from the data book for a standard cell library. This device has twice the ‘normal’ drive power, which indicates one of the compromises implicit in standard cell over full-custom, which is that the size (driving power) of transistors used in a cell is not tuned on a per-instance basis.

There are two types of semi-custom devices:

- standard cell
- gate array.

In standard cell designs, cells from the library can freely be placed anywhere on the device and the number of IO pads and the size of the die can be freely chosen. Clearly this requires that all of the masks used for a device are unique to that device and cannot be used again. (Mask making is one of the largest costs in chip design).

In gate array designs, the silicon vendor offers a range of chip sizes. Each size of chip has a fixed layout and the location of each transistor, resistor and IO pad is common to every design that uses that size. A gate array is configured for a particular design by wiring up the transistors (and other

components) in the desired way. Many transistors will be unused. The wiring up is, as usual, done with the top two or three layers of metal wiring. Therefore only two or three masks need be made to make a new design.

Silicon vendors will typically have stocks of wafers which have had the bottom 15 to 20 process steps made and are only awaiting final metalisation to be turned into usable devices. Hence design time is low. Risk is also low since the finished chip has only a small NRE cost and gate array technology is intrinsically conservative and hence reliable.

The disadvantage of gate arrays is their intrinsic low density of active silicon.

Standard cell designs use a set of well-proven logic cells on the chip, much in the way that previous generations of standard logic have been used as board-level products, such as Texas Instruments' System 74.

A variation on the semi-custom approach is to include full-custom macrocells such as processor cores in fixed positions on the wafer.

Exercise: Ignoring the title, determine from the data sheet whether the cell in figure 45 is for standard cell or gate array use? What other information about the cell is needed to prepare an audit of resources used in a design which has used this cell? The audit referred to is typically a report generated by the CAD tools which gives summary information.

3.5.3 Field programmable logic

About 25 to 40 percent of chip sale revenue now comes from field programmable logic devices. These are chips which can be programmed electronically on the user's site to provide the desired function. The Xilinx FPGA parts used in the Part 1B E+A classes are one of the most important examples of field programmable logic.

Field programmable devices may be volatile (need programming every time after power up), reprogrammable or one-time programmable. This relates to how the programming information is stored inside the devices, which can be in RAM cells or in any of the ways mentioned for ROM devices in section 1.2.

3.5.4 FPGAs

An FPGA (field programmable gate array) consists of an array of configurable logic blocks (CLBs), as shown in Figure 37. Not shown is that the device also contains a good deal of hidden logic used just for programming it. Some pins are also dedicated to programming. Such FPGA devices have been popular since about 1990.

Each CLB (Figure 38) typically contains one or two flip-flops, and has a few (five shown) general purpose inputs, some special purpose inputs (only a clock is shown) and two outputs. The illustrated CLB is of the look-up table type, where the logic inputs index a small section of pre-configured RAM memory which implements the desired logic function. For five inputs and one output, a 32 by 1 SRAM is needed. Some FPGA families now give the designer write access to this SRAM, thereby greatly increasing the amount of storage available to the designer. However, it is still an expensive way to buy memory.

All CLBs within a FPGA generally have the same structure, but FPGAs are available with lower and higher functionality CLBs. The best size of CLB is not yet clear. Modern designs of FPGA have a hierarchy of CLB interconnection patterns, giving CLB clusters within clusters.

An FPGA is very like a mask-programmed gate array to use. The design flow and CAD tools are virtually identical. The expenditure before the designer has the first device in her hands might be 1000 times lower. The cost of further devices is at least 10 times higher than mask-programmed devices, owing to the programming cost and wasted die area devoted to the programming activities. Clearly there is a crossover production volume point!

FPGAs also tend to be quite slow, owing to the signals passing through hidden logic used only for configuration.

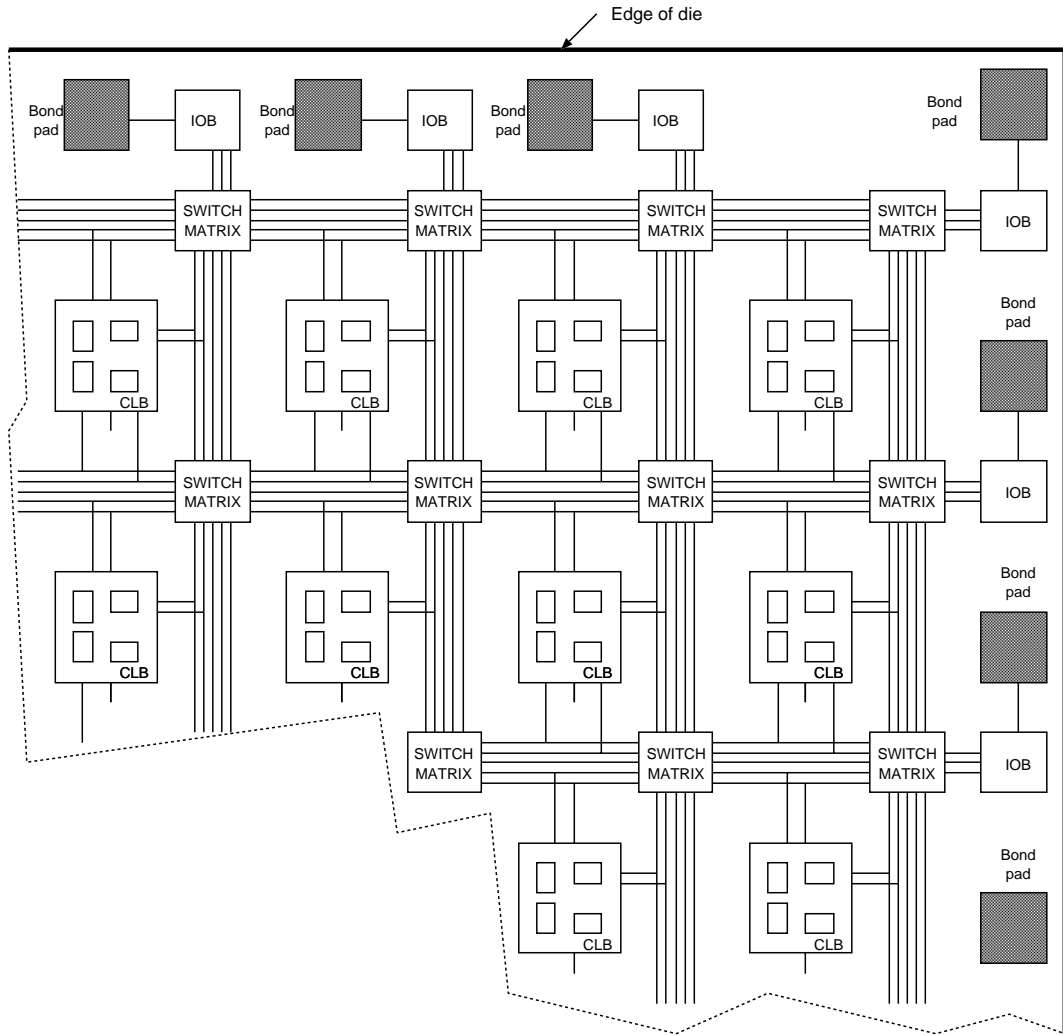


Figure 37: Field programmable gate array structure, showing IO blocks around the edge, inter-connection matrix blocks and configurable logic blocks.

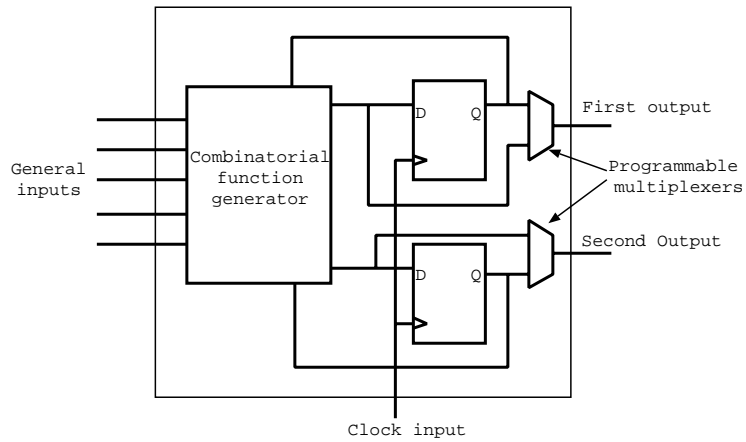


Figure 38: A configurable logic block for a look-up-table based FPGA.

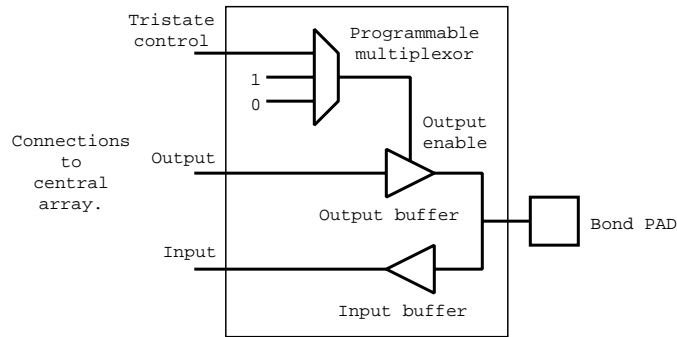


Figure 39: A simple IO block for an FPGA.

Often a designer or company will build prototypes and early production units using FPGAs and then use a drop-in mask-programmed equivalent once the design is mature and sales volumes pick up.

Exercise: Attempt to sketch the logic needed inside one of the programmable switch matrix boxes. You should find this is a vast amount of logic and therefore understand why FPGAs are so expensive for their functionality. (If you know how to use pass transistors, this will help your design).

3.5.5 PALs

A PAL is programmable array logic device. Figure 40 shows a typical device. Such devices have been popular since about 1985. The illustrated device has 8 product terms per logic function, and so can support functions of medium complexity. Such devices are very widely used and can feature high speed operation with clock rates of above 100 MHz. They are really just highly structured gate arrays. Every logic function must be multiplied out into sum-of-products form and hence is achieved in just two gate delays.

Programmable *macrocells* (Figure 41) enable the output functions to be either registered or combinatorial. Small devices (e.g. with up to 10 macrocells) offer one clock input; larger devices with up to about 100 macrocells are also available, and generally offer several clock options. Often some macrocells are not actually associated with a pin, providing a so called *buried state* flip-flop.

Figure 42 shows part of a PAL description for the PAL shown in figure 40, as entered by a designer in a typical PAL language, and part of the fuse map that would be generated by the PAL compiler. Each product line has seven groups of four fuses and produces the logical AND of all of the signals with intact fuses. An 'x' denotes an intact fuse and all of the fuses are left intact on an unused product lines in order to prevent the line ever generating a logical one (a gets ANDed with a bar etc.). The fuse map is loaded into a programming machine (in a file format known as JEDEC), an unused PAL is placed in the machine's socket and the machine programs the fuses in the PAL accordingly.

Exercise: How large a binary counter can the illustrated device implement? (Hint: are there sufficient product terms for the most-significant bit?)

PALs achieve their speed by being highly structured. Their applicability is restricted to small finite state machines and other *glue logic* applications.

3.6 Delay-power product

Aside from the selection of design methodologies just discussed, the designer must choose the basic circuit technology to use:

- silicon or GaAs
- bipolar or unipolar

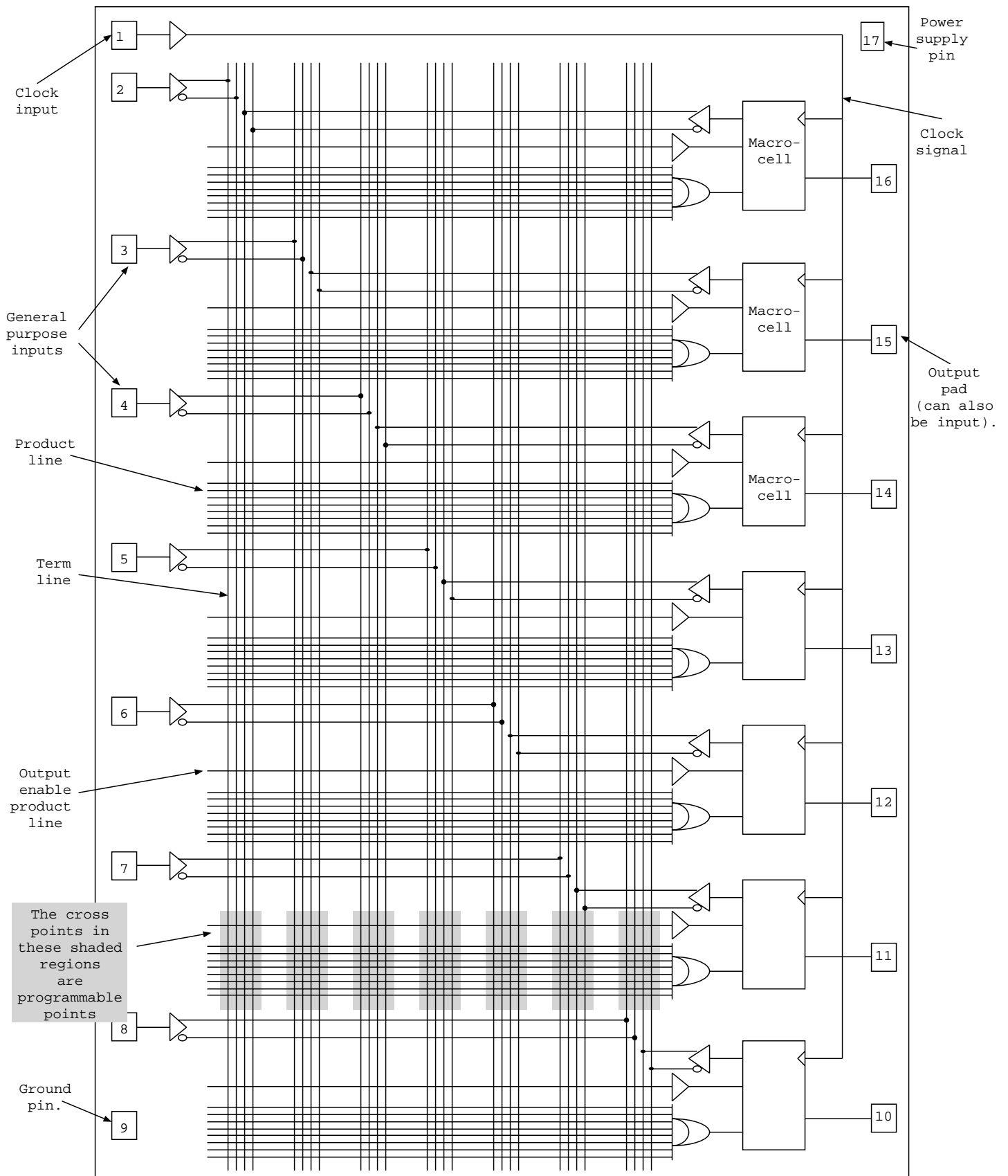


Figure 40: A typical PAL with 7 inputs and 7 I/Os.

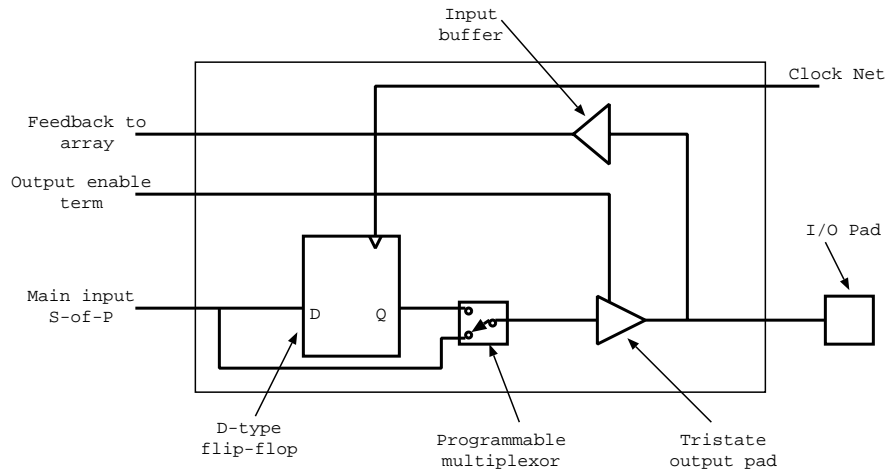


Figure 41: Contents of the PAL macrocell.

```

pin 16 = o1;
pin 2 = a;
pin 3 = b;
pin 4 = c

o1.oe = ~a;
o1 = (b & o1) | c;

-x-- ---- ---- ---- ---- ---- ---- (oe term)
--x- x--- ---- ---- ---- ---- ---- (pin 3 and 16)
---- ---- x--- ---- ---- ---- ---- (pin 4)
xxxx xxxx xxxx xxxx xxxx xxxx xxxx
xxxx xxxx xxxx xxxx xxxx xxxx xxxx
xxxx xxxx xxxx xxxx xxxx xxxx xxxx
xxxx xxxx xxxx xxxx xxxx xxxx xxxx
xxxx xxxx xxxx xxxx xxxx xxxx xxxx
x                                     (macrocell fuse)

```

Figure 42: Example programming of a PAL showing only fuses for the top macrocell.

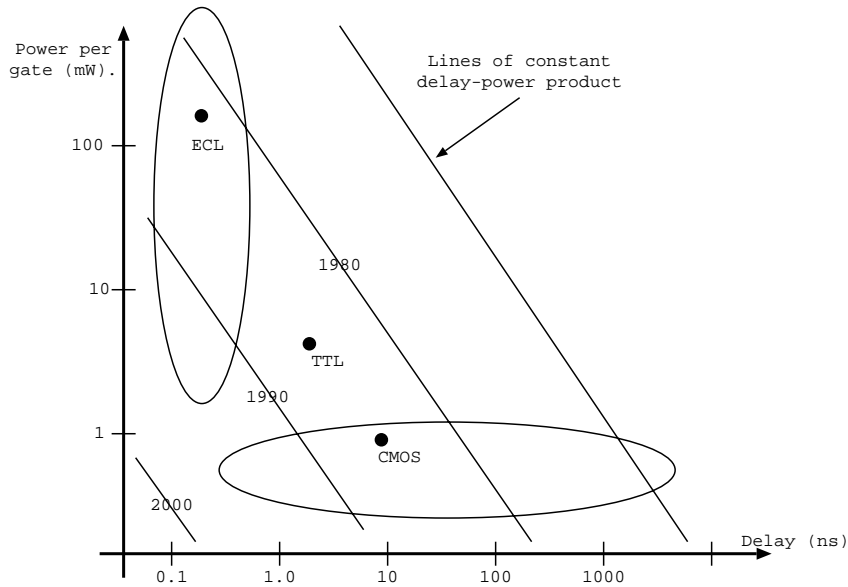


Figure 43: Delay-power style of technology comparison chart.

- saturating or non-saturating.

As the size of transistors on chips is reduced with ever improving fabrication equipment, the number of devices per unit area increases and the delay-power product (power consumption multiplied by propagation delay) decreases.

We here consider the main three circuit technologies of the 70's and 80's using a quad AND gate chip SSI chip for example (shown generically in Figure 51). In the 90's, CMOS is dominating most application areas.

technology	device	propagation delay (ns)	power (mW)	product (pJ)
CMOS	74hc00	7 ns	1 mW	7
TTL	74f00	3.4 ns	5 mW	17
ECL	sp92701	0.8 ns	200 mW	160

The CMOS gate has the property that it consumes virtually no power when outputs are not changing.

The ECL gate is an old technology, but it is still the fastest.

Gates in the core of an IC tend to be one order better in speed or power when compared with the SSI gate used in the example (they have reduced drive requirement).

All of these aspects have to be considered by the system designer and it is therefore a skilled job.

3.7 Fanout and delay estimations.

Figure 44 shows a typical net, driven by a single source. To change the voltage on the net, the source must overcome the stray capacitance and input loads. The fanout of a gate is the number of devices that its output feeds. The term *fanout* is also sometimes used for the maximum number of inputs to other gates a given gate is allowed to feed, and forms part of the design rules for the technology.

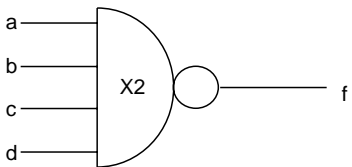
The speed of the output stage of a gate, in terms of its propagation delay, decreases with output load. Normally, the dominant aspect of output load is capacitance, and this is the sum of:

NAND4 Standard Cell

Library: CBG0.5um

4 input NAND gate with x2 drive

Schematic Symbol



Simulator/HDL Call

NAND4X2(f, a, b, c, d);

Logical Function

$F = \text{NOT}(a \& b \& c \& d)$

ELECTRICAL SPECIFICATION

Switching characteristics : Nominal delays (25 deg C, 5 Volt, signal rise and fall 0.5 ns)

Inputs	Outputs	O/P Falling		O/P Rising	
		(ps)	ps/LU	ps	ps/LU
A	F	142	37	198	33
B	F	161	37	249	33
C	F	165	37	293	33
D	F	170	37	326	34

Min and Max delays depend upon temperature range, supply voltage, input edge speed and process spreads. The timing information is for guidance only. Accurate delays are used by the UDC.

CELL PARAMETERS : (One load unit = 49 fF)

Parameters	Pin	Value	Units
Input loading	a	2.1	Load units
	b	2.1	
	c	2.1	
	d	2.0	
Drive capability	f	35	Load units

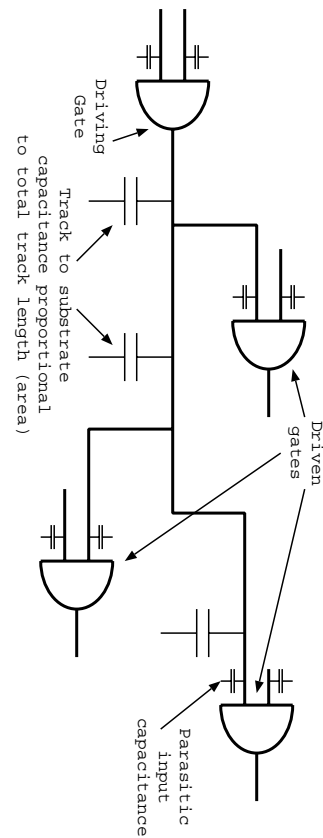


Figure 44: Logic net with tracking and input load capacitances.

Figure 45: An example cell from a manufacturer's cell library.

1. the capacitance proportional to the area of the output conductor
2. the sum of the input capacitances of the devices fed.

To estimate the delay from the input to a gate, through the internal electronics of a gate, through its output structure and down the conductor to the input of the next gate, we must add three things:

1. the internal delay of the gate, termed the intrinsic delay
2. the reduction in speed of the output stage, owing to the fanout/loading, termed the derating delay,
3. the propagation delay down the conductor.

The propagation delay down a conductor obeys standard transmission line formula and depends on the distributed capacitance, inductance and resistance of the conductor material and adjacent insulators. For circuit board traces, resistance can be neglected and the delay is just the *speed of light* in the circuit board material: about 7 inches per nanosecond, or 200 metres per microsecond.

On a chip, the speed of light can be neglected because chips are physically small, but the resistance of the aluminum conductors is sufficiently large to have an effect for critical applications, such as master clock signals.¹

The upshot of this is that on a chip, we can mostly neglect the third aspect of delay, whereas on a circuit board, we need to model certain critical conductors as components in themselves. These have a simple delay model, whose value can be set by post routing back annotation.

The first two aspects of delay may easily be absorbed into the model for a component. We can use the following formula

$$\text{device delay} = (\text{intrinsic delay}) + (\text{output load} \times \text{derating factor}).$$

The intrinsic speed of a device is given in its data sheet, as is the derating factor and the loading factor of its inputs. Typical values for a standard cell array are shown in Figure 45. The output load is the sum of a track dependent part and the fanout dependent part.

The track dependent part is a library constant times the track area.

The load dependent part is the sum of the input loads of all of the devices being fed.

Exercise: How true is the above model of signal delay when we use field programmable gate arrays with high effective track resistance? The conductors on FPGAs consist of many sections of real metallic conductor interconnected by the user-programmed connection points. We then have considerable resistances at points along each conductor's path. How would you estimate the various delays for the staggered arrival of a signal on each part of the net?

3.8 Comparative view of digital logic technologies

Table 2 gives some randomish numbers for what can be achieved with today's chips. These numbers are only accurate to about half an order of magnitude. Technology is always advancing!

For CMOS technology, the number of transistors on a chip and the clock rate both seem to double every three years.

¹Older technologies used polysilicon interconnections which had significant resistance and so the different gates connected to a polysilicon net would experience different arrival times of a signal at their inputs.

Technology	Maximum clock speed	Maximum gate count	Maximum pin count
GaAs bipolar	50 GHz	50	5
GaAs fet	3 GHz	300	100
Si ECL	3 GHz	5000	300
Si CMOS	400 MHz	1 E6	500
Within Si CMOS			
Full-custom	400 MHz	1 E6	500
Standard cell	200 MHz	40000	500
Gate array	100 MHz	20000	500
PAL	75 MHz	500	68
FPGA	35 MHz	4000	200

Table 2: Approximate limits on contemporary technology.

4 Hardware, Software and Design Partition.

A hardware system can be designed using only hardware or hardware with an embedded processor and software in ROM. Perhaps several processors are to be used. A design also must be partitioned into standard components and application specific integrated circuits (ASICs). The choice of which parts of the design are realised in what technology is the design partition problem.

4.1 Embedded Systems

A processor which is permanently connected to a single ROM which contains all of the software that the processor will ever execute is called an *embedded processor*. Here we consider when to use an embedded processor, and then, when using one, the designer has to chose whether to design a new one or use an existing one.

The main difference between a hardware solution and a software solution is the degree of parallelism. Processors typically execute one instruction at a time, reusing the same hardware components again and again. Hardware solutions typically have dedicated circuits for each function. If most of the hardware is likely to be idle for most of the time, a processor is preferred, but if a processor cannot achieve the throughput required, increased parallelism using hardware is preferred. Complex functions normally require a processor, but CAD tools are evolving, allowing complex functions to be expressed alogrithmicly but implemented in logic gates.

High speed processing normally requires dedicated hardware. For instance, consider the error correction performed by a CD player to overcome dirt and scratches. When CD players first came out, this error correction was done with dedicated hardware, but today, microprocessors have increased in speed, and so the function can be done using the processor which is already there to provide other complex functions (e.g. track skip). However, on the latest, 24 speed CD ROM drives for computers, the error correction must be done 24 times faster, and so dedicated hardware is re-introduced.

Processors give flexibility: if the design is likely to be changed in minor ways to form new models, or to provide field upgrades, then using a processor and providing a new software release is a good technique.

Standard processor chips can be very cheap for a given performance. This is because a very great deal of effort is put into their design and they are sold in large quantites, thereby reducing price. Part of the cost of any product is in testing it. A standard processor not only comes with its own test qualification programme, it is able to execute software to help test the rest of the system.

Processors can be placed on an ASIC.

The decision to design and use a custom processor for an application should not be taken lightly. It can be useful, however, when the application is very simple or highly specialised. Examples are signal processing, where a particular algorithm needs to be executed at high speed to track a missile or assist computer vision.

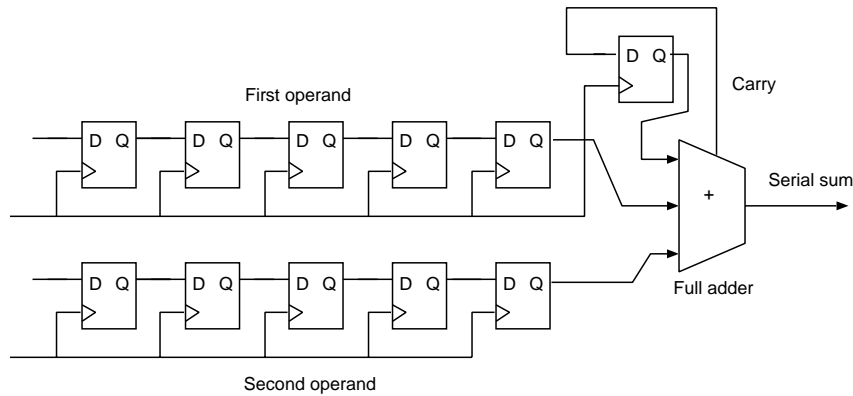


Figure 46: Addition of two integers serially, l.s.b. first

4.2 Logic Synthesis and Collapsing

If one considers an embedded processor connected to a ROM, it may be viewed as one large FSM. Since for any given piece of software, the ROM is unlikely to be full and there are likely to be resources in the processor that are not used by that software, the application of a good quality logic minimiser to the system, while it is in the design database, could trim it greatly. In most real designs, this will not be helpful: for instance, the advantages of full-custom applied to the processor core will be lost. Infact, the minimisation function may be too complex for most algorithms to tackle on today's computers.

However, there are cases that are practical. For instance, many digital signal processing applications do use this approach. A processing algorithm typically consists of multiple stages with names such as pre-emphasis, equalisation, coefficient adaptation, fft, deconvolution and reconstruction. Each of these steps normally has to be done as fast as possible and so parallelism through multiple instances of ALU and register hardware is needed. However, there can be data dependencies in that one step cannot be done until the previous is completed, or in that two steps take the same amount of processing as one other step. In these cases, we desire to reuse the same hardware for each step. The Cathedral DSP compiler is the most famous tool for helping design such circuits. It performs folding of the algorithm to reuse hardware from one step in the next.

Bit-serial arithmetic is also a key approach in custom signal processing. The adder shown in figure 46 requires only about 10 gates to add any number of bits. It operates l.s.b. first so that the carry from one digit is available for the next digit. The clock speed for the calculation does not need to be reduced if longer words must be consumed, but the pipeline delay does go up. Multiplication is also easy in bit serial form, as shown for multiplication by a constant in figure 47. Consider the same functions implemented in naive, broadside implementations.

A battery operated walkman CD ROM player may implement a great deal of DSP. The printing on the unit may advertise '8 times oversampling D-to-A bitstream convertor and MASH bit mapping'. The processing power required for this is as much an Intel Pentium can produce, but through the use of bit serial arithmetic, it takes very little power or silicon area.

The problem with bit-serial processing for the human designer is that it is very tricky to work out where individual bits are at any one time, and the optimisations that are possible cannot readily be seen by the human. Therefore, good quality CAD tools are vital in the design of such systems.

4.3 ASICs

The cost of developing an ASIC has to be compared with the cost of using an existing part. The existing part may not perform the required function exactly, requiring either a design specification change, if possible, and/or some additional *glue logic* to adapt the part to the application.

More than one ASIC may be needed under any of the following conditions:

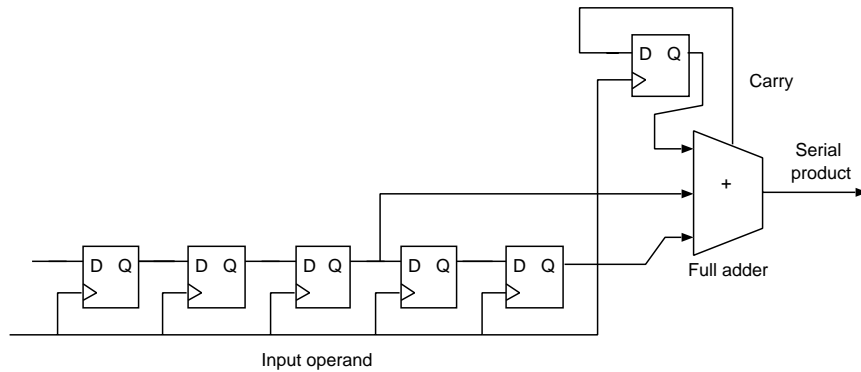


Figure 47: Bit-serial multiplication of an integer by a hardwired constant

- application specific functions are physically distant
- application specific functions require different technologies
- application specific functions are just too big for one ASIC
- it is desired to split the cost and risk or reuse part of the system later on.

Factors to consider on a per chip basis:

- pad count limitation (pad density limit of 15 per mm)
- power consumption limitation (powers above 5 Watts need special attention)
- gate count limitation (above 100 kgates is big for CMOS)
- speed of operation — influences choice of technology
- will it be core or pad bound ?
- special considerations :
 - special static or dynamic RAM needs
 - mixed with analogue parts ?
 - high power handling outputs for load control: e.g. motors.
- availability of a developed module for future reuse

Power dissipation is generally proportional to:

$$\text{Power} \propto (\text{clock speed}) \times (\text{number of gates}) \times (\text{supply voltage squared}).$$

Various technologies can be positioned in the speed-power plane, as shown in Figure 43 and as time goes by, things get better.

4.4 Partitioning example: The Cambridge Fast Ring two chip set.

Two devices were developed for the CFR local-area network, illustrating the almost classical design partition required in high-speed networking. They were never given grander names than the *ECL* chip and the *CMOS* chip. The block diagram for an adaptor card is shown in Figure 48.

The ECL chip clocks at 100 MHz and contains the minimal amount of logic that needs to clock at the full network clock rate. Its functions are:

- implement serial transmission modulator and demodulator

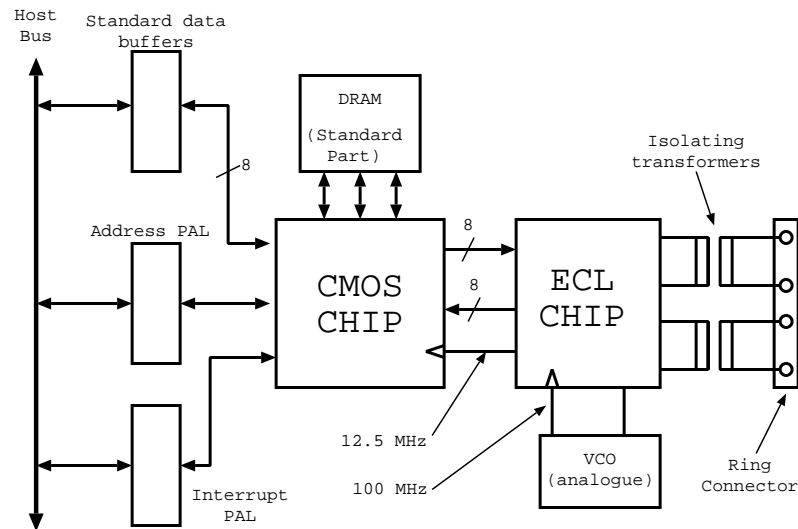


Figure 48: Example of a design partition — the adaptor card for the Cambridge Fast Ring.

- convert from 1 bit wide to 8 bits wide and the other way around
- perform reception byte alignment (when instructed by logic in the CMOS chip).

Other features:

- ECL logic can support analogue line receivers at low additional cost so can receive the incoming signal directly on to the chip.
- ECL logic has high output power if required (1 volt into 25 ohms) and so can drive outgoing twisted pair lines directly.

The CMOS chip clocks at one-eighth the rate and handles the complex logic functions:

- CRC generation
- full/empty bit protocol
- minipacket storage in on-chip RAM
- host processor interface
- ring monitoring and maintenance functions.

The ECL chip has at least 50 times the power consumption of the CMOS chip. The CMOS chip has more than 50 times the gates of the ECL chip.

Two standard parts are used to augment the CFR set: the DRAM chip incorporates a dense memory array which could not have been achieved for anywhere near the same cost onboard the CMOS chip and the VCO (Voltage Controlled Oscillator) device used for clock recovery was left off the ECL chip since it was a difficult-to-design analogue component where the risk of having it on the chip was not desired.

PALs are used to ‘glue’ the network interface itself to a particular host system bus. Only the glue logic needs to be redesigned when a new machine is to be fitted with the chipset. PALs have a short design turn-around time since they are field programmable.

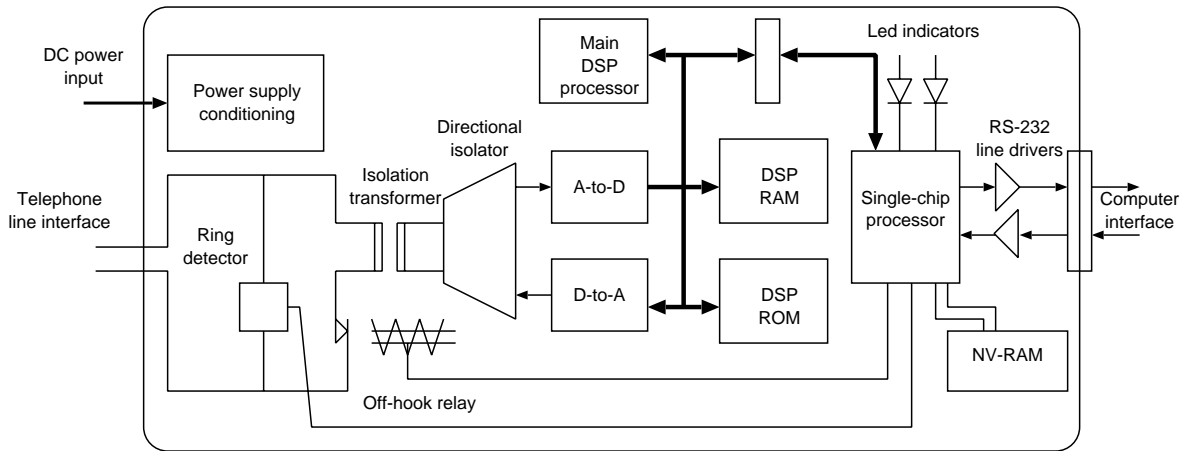


Figure 49: Example of a design partition — an external modem.

4.5 Partitioning example: An external PC Modem.

Figure 49 shows the block diagram of a typical modem. The illustrated device is an external modem, meaning that it sits in a box beside the computer and has an RS-232 serial connection to the computer. It also requires its own power supply.

The device contains a few analog components which behave broadly like a standard telephone, but most of it is digital. A relay is used to connect the device to the line and its contacts mirror the ‘off-hook’ switch which is part of every telephone. It connects a transformer across the line. The relay and transformer provide isolation of the computer ground signal from the line voltages. Similarly the ringing detector often uses a optocoupler to provide isolation. *Clearly, these analog aspects of the design are particular to a modem and are designed by a telephone expert.*

Modems from the 1960’s implemented everything in analog circuitry since microprocessors were not then possible. Today, two microprocessors are often used, but as processing power increases, this can be reduced to one (or sometimes even none, if the main processor of a PC provides the functions needed).

The reason for two processors are interesting and will be discussed in lectures.

Note that the non-volatile RAM requires a special manufacturing processing step and so is not included as a resource on board the single chip processor. Similarly, the RS-232 drivers need to handle voltages of +/- 12 volts and so these cannot be included on chip without increasing the cost of the rest of the chip by using a fabrication process which can handle these voltages. The NV-RAM is used to store the owner’s settings, such as whether to answer an incoming call and what baud rate to attempt a first connection, etc..

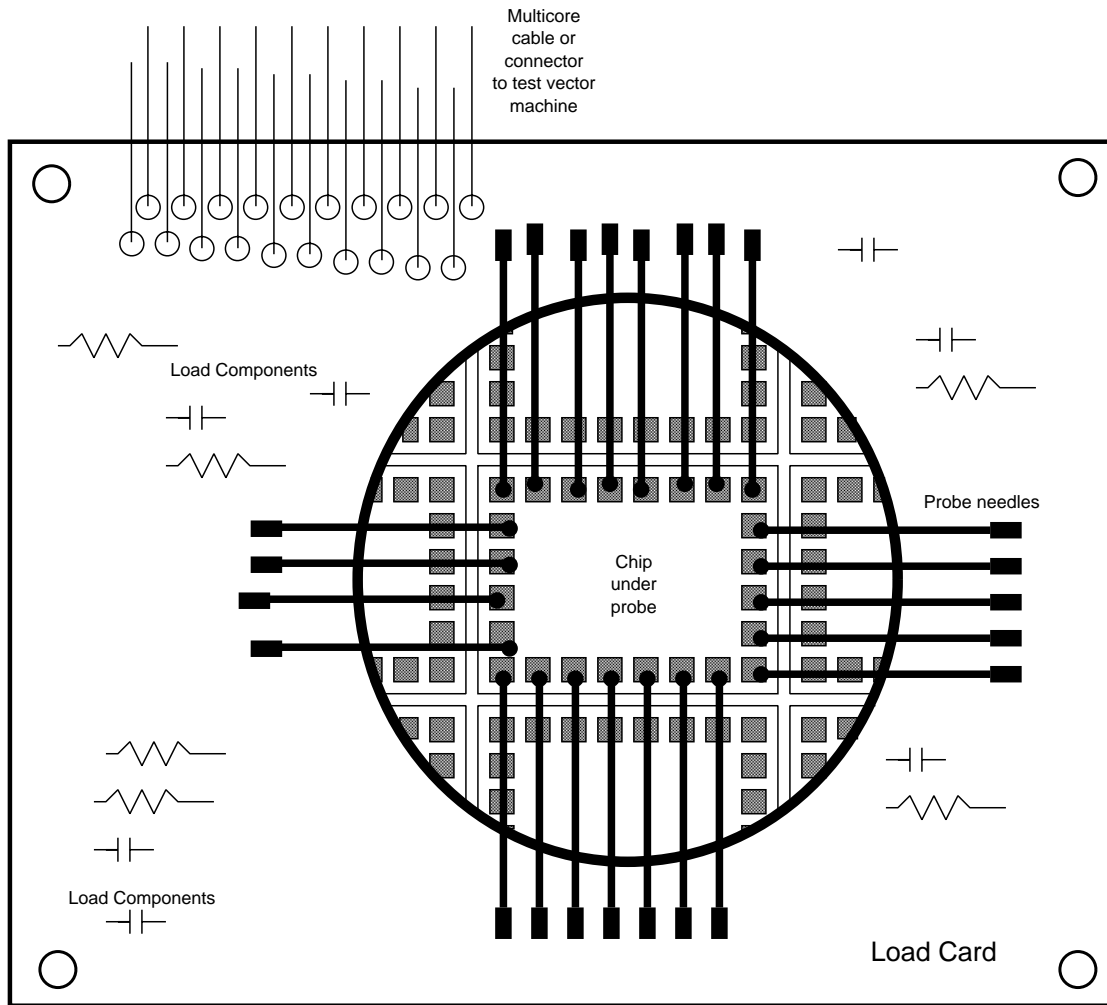


Figure 50: A load card with wafer probe pins for testing a chip before the wafer is diced.

5 Chip, board and system testing.

Note: This section will almost certainly not be lectured or examined.

When a design is complete there will be a production run. For consumer goods, such as a CD player, we might make 5000 units per month for a design life of six months. The aim is that 99.9 percent of units will work.

To ensure the number of failures is acceptable it is necessary for every component and subassembly to be tested before the next stage of assembly, and, in addition, a final test program for the finished product must be used.

If we assume that faults are uniformly distributed over the design and that a reasonable fraction (say half) of the total product functionality is only occasionally used in normal operation, then it is clear that a special test program is required which is able to exercise (nearly) every function of the unit. Since testing time can be a major part of manufacturing time, the test program must run in as short a time as possible.

The fraction of functions covered by a test program is the *fault coverage*, which is normally expressed in percent.

A test program consists of a sequence of stimulus and expected results that must be applied to the device. If the actual results do not match the expected results, then a fault is detected. Of course it is possible for faults to exist which do not show up as a result of a poor test program.

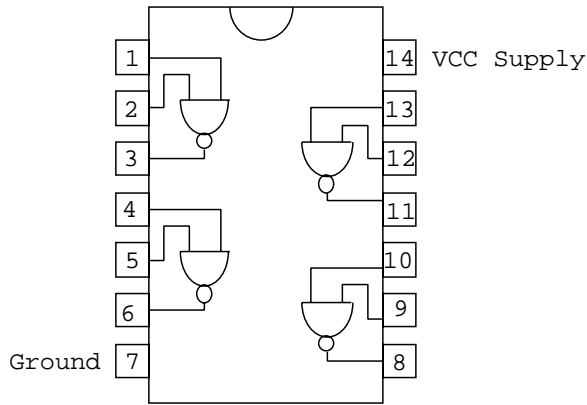


Figure 51: Pin connections of a 7400 quad NAND gate device.

A good test program will achieve 98 to 99.5 fault coverage in a few seconds of testing. Military, aerospace and medical standard equipment may need to have 100 percent coverage.

For a digital chip or a circuit board, a tangible subset of ‘every possible type of fault’ must be used. This is the *stuck at* model of faults. In this model, the only fault considered is for one single net of the whole design to be shorted to the power supply or ground. If there are n nets, there are $2n$ faults that could occur.

5.1 A typical test vector system for digital modules

A test program for a chip or card consists of a set of test vectors generated from simulation or otherwise. The test program is applied to a set of test probe points. These are typically the pins or pads of a chip or a *bed of nails* that a circuit board is pushed on to.

The test program consists of a list of vectors. The length of each vector is equal to the number of probe points. The elements of each vector are an ASCII character taken from one of the following:

- 1 --- apply a logic one to this probe point
- 0 --- apply a logic zero to this probe point
- z --- apply high impedance to this probe point
- H --- expect logic one at this probe point
- L --- expect logic zero at this probe point
- x --- do not care what happens at this point
- p --- a pin not connected to the tester, but connected to power etc.
- c --- clock this pin mid-way through the cycle.

Typically other symbols enable special pulses or varying power supply voltages to be applied. The nature of these special signals is specified in separate tables included with the test program.

The vectors are applied in sequence at some clock rate (e.g. 10 million vectors per second) and there may be 100000 for a large chip or board.

If any of the H or L points do not match, the device under test has failed.

Example — a test of one part of a 7400 quad nand gate shown in Figure 51. We allow spaces in the vectors for clarity. For the first gate, pins 1 and 2 are inputs and 3 is the output. There are four vectors in this rather inadequate program.

```

000 000 0 001 111 1
123 456 7 890 123 4
[ 00H 00x p H00 x00 p ]
[ 01x 00H p H00 x00 p ]
[ 10x 00H p H00 x00 p ]
[ 11L 00x p H00 x00 p ]

```

Exercise: Write a 100 percent coverage test program for a 7400.

5.2 The need for test modes.

For some designs it would take a very long test program to test certain functions. For example, the leap year circuitry in a digital watch might need four years' worth of clock pulses before being exercised. Therefore, the throughput of the testing station would be low.

For other designs, the observability of internal state can be low. This occurs when there is a lot of internal state and few outputs. For example, a credit card PIN number checker need only have one output net, saying good or bad. Many different test sets would have to be presented before the affect of every part of the internal logic could be felt at the output net. This either leads to low fault coverage or long test programs.

Testability can be increased by adding either:

1. Additional outputs to bring out the state of internal nets
2. Additional inputs, which are tied off to one logic level during normal operation, but which during testing can shorten the length of counter cycles or cause more of the internal state to be brought out to the existing pins.

The addition of a single test input is often the preferred solution.

Exercise: Consider the testability of devices with built in redundancy.

5.3 Fault Simulator and Automatic Test Generation

The CAD tool which determines the fault coverage of a given test program against a given design is called the *Fault Simulator*. Its outputs are the fault coverage percentage and the list of stuck at faults not detected by the given test program.

Exercise: Design a fault simulation algorithm whose expected running time is $O(n \times v)$ where there are n nets and v vectors. Can you design a better algorithm ?

A test program generator is a CAD tool which generates a good test program for a design by analysing the structure of the design. It can also inform the designer that her masterpiece is difficult to test, allowing her to make modifications at an early stage.

Manual test programs are often generated from the designer's functional simulations — the ones she has used to determine whether the design meets its specification. Typically the designer concatenates each of these simulations (along with some others which she feels may bump up the fault coverage) and feeds them into the standard simulator. The simulator is run in a mode where it prints a test vector to a file at a certain fraction through every clock event and this file is then the test program. It is clear that a working device will pass this test program, but there is no guarantee that such an approach will yield sufficient fault coverage.

Automatic test program generation can involve sophisticated logic analysis algorithms. The section in Kinniment and McLauchlan on this subject is very good.

5.4 Boundary scan

Test vectors have traditionally been applied to the device under test using a separate wire from the testing machine to each probe point. This is suitable when the device under test is

1. a die on a wafer
2. a packaged chip
3. a printed circuit board.

But today it is desirable to apply these same test programs to the component parts of a complete assembly. Examples are:

1. a multi-board finished product
2. a macrocell placed on a larger ASIC.

With boundary scan, the chip can be put into a boundary scan test mode where the IO pads of a chip are connected into a shift register. The shift register is a chain of flip-flops, one being built into each IO pad for the purpose.

In test mode, the test vector is shifted into the flip-flops. The signal on an input pad is overridden by the value in the scan flip-flop. The test path is then activated using a strobe line which latches the values on the output pads into the scan flip flops. The test vector with output data is then shifted out and examined as the next vector is shifted in.

IO pads are intrinsically big (to get the bond wire on) and slow (owing to the large transistors and diodes used in them). The speed and complexity penalty of adding boundary scan logic to them is therefore low.

Owing to the serial nature of scan path testing, multiple chips can be put in series on the scan path. These chips can be macrocells on one large piece of silicon, or separate pieces of silicon on a circuit card.

Speed of testing can be a problem for long scan chains. Clearly the rate of exchange of test data between the device under test and the testing machine is much lower with a serial interconnection than offered by the multiple parallel paths of the bed-of-nails approach.

5.5 General scan path

General scan path testing is similar to boundary scan, but threads the scan path through all of the flip-flops within the design.

In addition to its normal connections and operations, each flip-flop must support a scan mode and has an extra input for connection to the output of the previous flip-flop in the scan path and a connection to a global control net which enables scan mode.

- Advantage — full testability and observability
- Disadvantage — increase in complexity and reduction in device speed.

5.6 JTAG boundary scan.

The IEEE 1149 (JTAG) international standard defines four wires per chip available for boundary scan. These are:

tms	—	test mode select — put high to enter boundary scan mode
tdi	—	test data input — serial data input
tck	-	test clock to clock each serial bit in
tdo	-	test data output to read back old data as new is shifted in.

For normal operation of the chip, connect the three inputs to logic zero (ground). A fifth wire, test reset, is often found in new implementations.

Figure 52 shows a boundary scan path that is inserted as part of the electronics of each input or output pin. When the test mode select is low, this logic has no effect, but when in test mode, the boundary scan logic takes over from the input pads.

The required logic is based around a shift register with one bit per pin. The JTAG test data input and output are placed at the ends of the shift register and the clock is the JTAG test clock. An internal counter, not shown in the figure, counts the shifted bits until a whole new vector has been

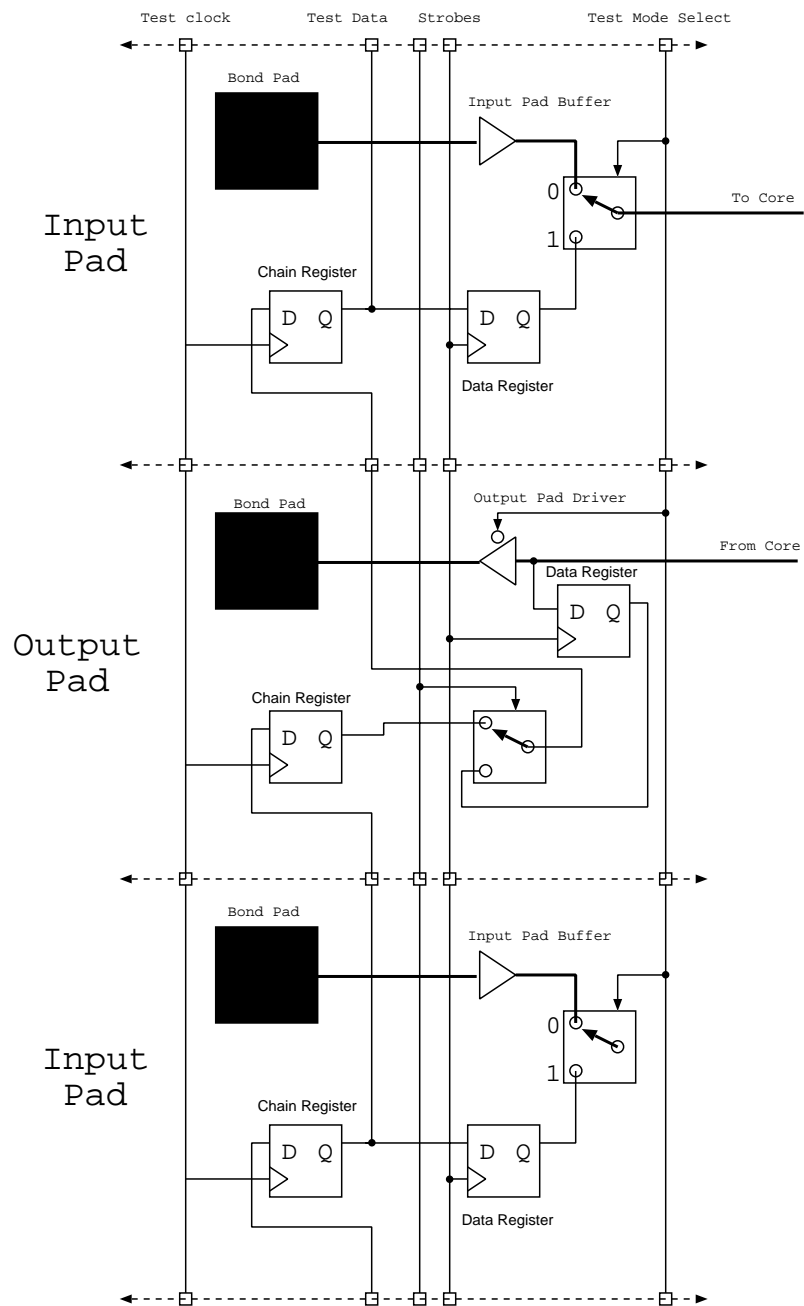


Figure 52: Circuit needed in input and output pads to achieve boundary scan

inserted and then generates strobe signals to transfer data in and out of the holding data registers in each pad.

Disadvantage: The main problem with JTAG is that four pins is quite a large number.

Advantage 1: Chips can be tested in place on their final circuit board (an *in-circuit* test).

Advantage 2: If the ‘chip’ is actually a macrocell built into a larger chip or module, then the ‘chip’ can be tested as an autonomous subunit, reusing its own standard test program.

The JTAG port of a chip is becoming used for additional functions:

- device type and serial number,
- in-circuit programming of reusable FPGAs,
- single step debugging of embedded microprocessors.

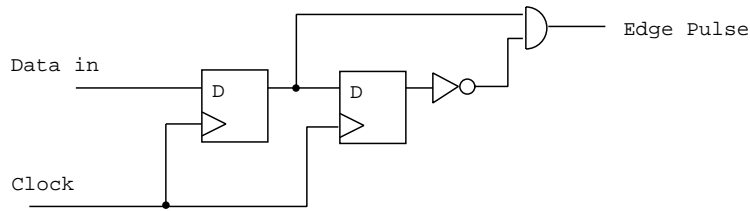


Figure 53: A simple edge detector

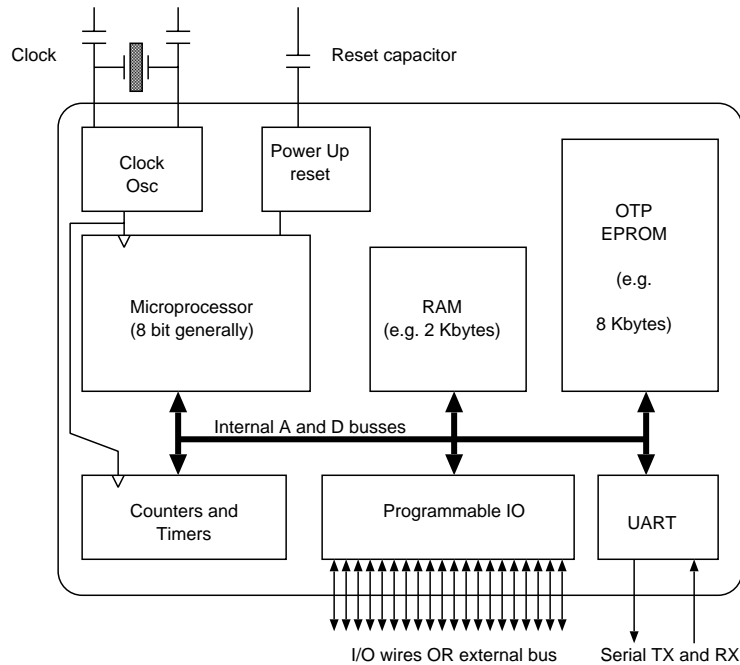


Figure 54: A typical single-chip microcomputer.

6 Further Circuit Structures.

This section introduces some further, fundamental building blocks and includes examples of real products. The important points to note are the reasons for the design partition.

6.1 Single Chip Processor

Figure 54 shows the block diagram of a typical single-chip microcomputer. Such devices are available in 24, 40 or 80 pin packages and can cost only one or two pounds. The device contains the whole of a computer, requiring externally only a power supply and some clock and reset components. A number of programmable I/O pins are provided, but generally another basic mode of operation is provided where 24 or so of them turn into the processor address and data busses for connection of external devices, such as further RAM or EPROM.

Programming is performed in all of the ways available for ROMs. Generally, users will use windowed EPROM devices for development and prototypes and ship product with OTP devices.

A Universal asynchronous receiver and transmitter (UART) is normally always provided to help connect to RS-232 serial interfaces.

These devices are found everywhere today, from keyboards, clock radios, mice, telephones, car window winders, 'computer controlled' cassette recorders, cell phones and modems.

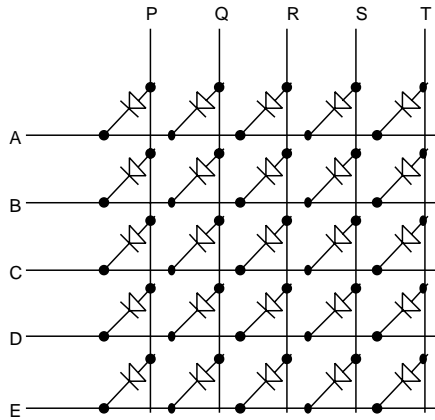


Figure 55: LEDs wired in a matrix to reduce external pin count

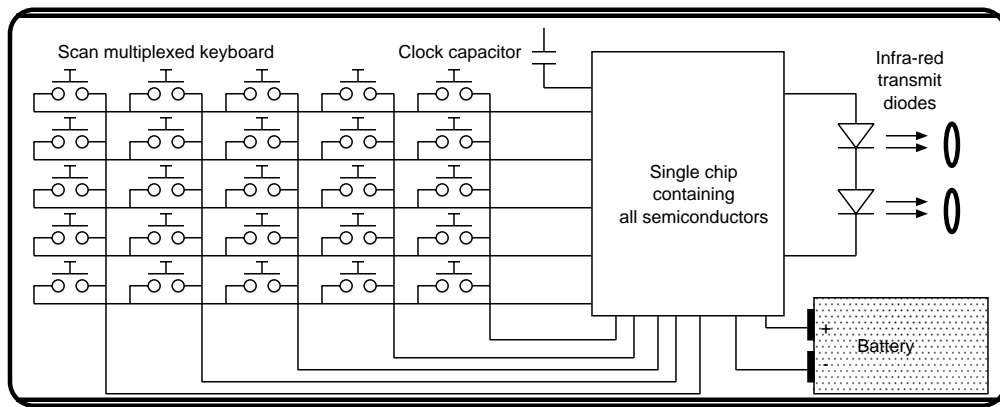


Figure 56: IR Handset Internal Circuit

6.2 Scan multiplexing.

When a large number of leds or switches need to be connected, as in a display or keyboard, the number of connections can normally be reduced by connecting them in a matrix. Figure 55 shows that 25 LEDs can be connected to just ten signals.

The matrix gives a *scan-multiplexed* display or keyboard. In the display, one vertical column line can be driven to logic one at a time and a zero placed on the horizontal lines that should be illuminated in that column. A circuit to repeatedly read out the contents of a RAM and display it is shown in figure 57. Clearly, the desired LEDs are not all on at once, but by scanning the display faster than the human eye can detect flashes (about 50 Hz) and by using sufficiently large currents, so that the elements are brighter than would otherwise be required, this is overcome. The current is set by the value of a series current limiting resistor that is not shown.

For a scan-multiplexed keyboard, the switches take the place of the LEDs. Push-to-make, normally open switches must be used and the user should not press two at once. The scanning circuit must take one row line low in turn. Pull-up resistors keep the column lines at logic one unless a switch to a low row line is pressed.

Figure 56 shows the full PCB circuit of a typical infra-red remote controller. The PCB circuit of a toy electric organ would be identical, but with the IR diode replaced with a loudspeaker: clearly the chip would be different. A pocket calculator also has roughly the same circuit, except there is also a display.

Exercise (long): Sketch the full circuit of the chip in an infra-red remote control handset, as used for TVs and VCRs etc.. The device should have about 50 push keys. Output is through a pair

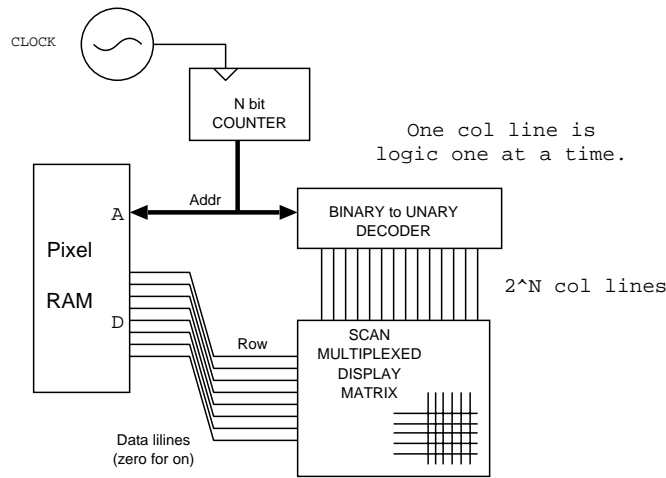


Figure 57: Scan multiplex logic for an LED pixel-mapped display

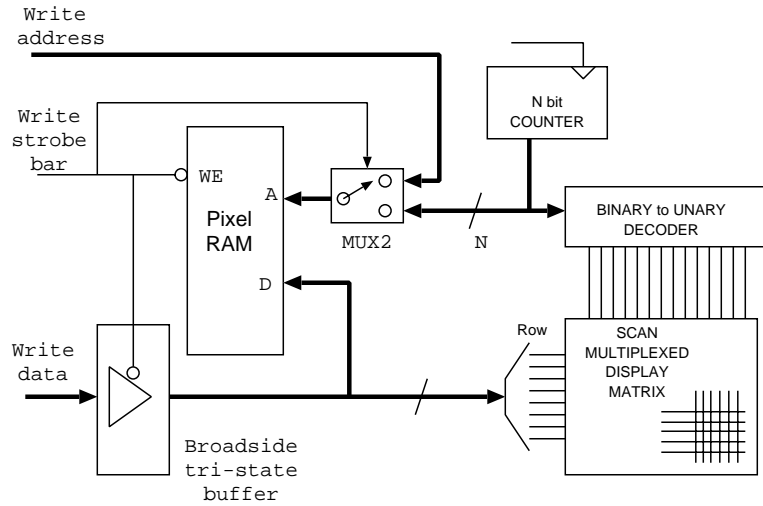


Figure 58: Addition of pseudo dual-porting logic.

of infra-red transmitting diodes which go on or off together. The design must include a ROM which can be programmed at chip manufacture with the desired sequence of pulses needed. The ROM must contain 16 bits for each key and these bits must be transmitted in turn through the IR diode using a 1 millisecond pulse of light for a zero and a 3 millisecond pulse for a one, with a 1 millisecond gap between each pulse.

Exercise: For the IR handset, is it necessary for the clock to be running at all times or only when a key is pressed? This design aspect will affect battery life greatly.

Figure 58 shows how a multiplexor can be added to allow update of the display RAM from a processor or similar. The display will briefly show the wrong data while being updated, but this may be imperceptible to the human. The compromise implied by such *pseudo*-dualporting is accepted instead of having the cost penalty of truly dual-ported storage. (A second independent port to a RAM or register file greatly adds to silicon area).

Figure 59 shows the use of a ROM as a function look up table. All electric guitarists use a distorting amplifier and are picky about the exact nature of the distortion. By using a ROM which contains a function appropriately distorted from the identity function, any desired distortion can be achieved.

The input A-to-D circuit samples the input signal at 44100 samples per second, which is the CD sampling frequency. The D-to-A converter on the output reconstructs an analogue signal.

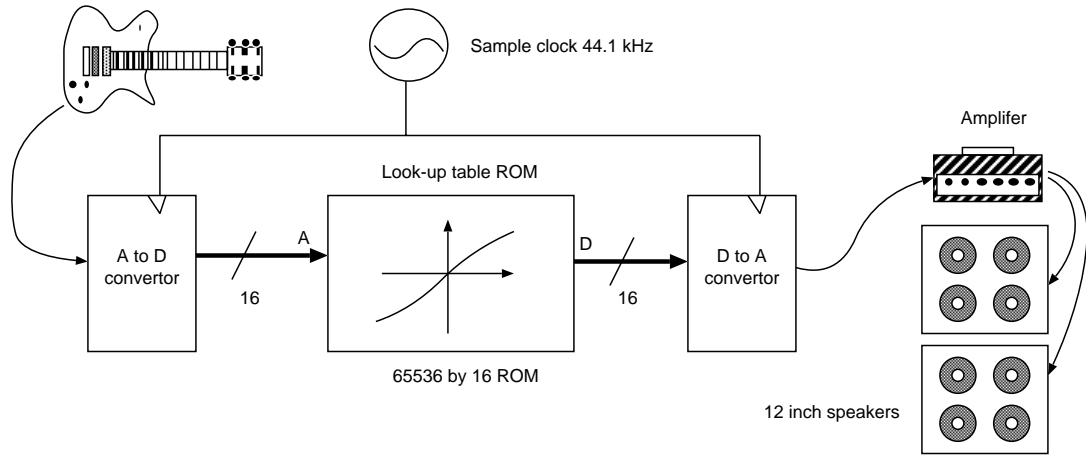


Figure 59: Use of a ROM as a function look-up table.

The analog signals will vary both positively and negatively about zero. Two's complement representations are therefore normally used. Note that to convert from offset binary to two's complement, one just puts an inverter in the most significant bit.

Exercise: Sketch an alternative circuit using a microprocessor instead of the hardwired solution. What pros and cons do the two approaches have? How does a processor help with the user interface?

Figure 60 shows the use of a static RAM to achieve a delay. The counter addresses each location in turn, and at each location, the old contents are read out before the new contents are put back.

Exercise: The circuit shown produces only a delay. For a true echo effect, the original signal must also appear at the output. To achieve multiple echos, the output must be fed back to the input at reduced amplitude. Introduce an adder and any additional logic to achieve these results. Use the fact that division by a constant power of two can be achieved by shifting right while leaving the sign bit unchanged.

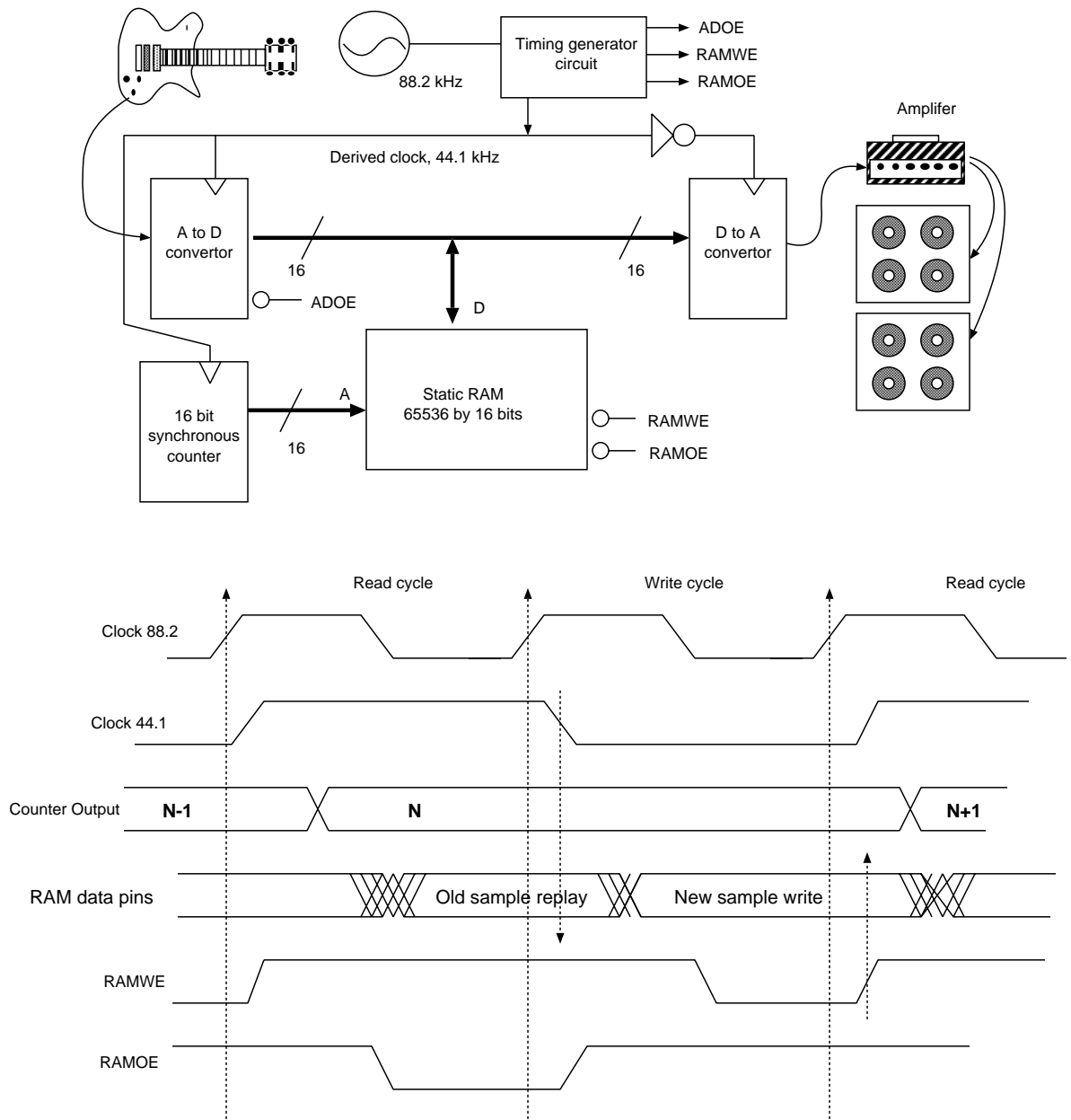


Figure 60: Use of an SRAM to make the delay required for an echo unit.

Example Exam Questions

Q. A circuit is needed to divide a very high frequency clock by an average of 21.6. This is to be done with a counter which accepts a division ratio input to make it divide either by 21 or 22. A further counter, made from slower and cheaper logic, is clocked from the output which generates the control signal to the divider, such that the ratio of 21 and 22 counts overall gives the desired target division ratio. Design it. Prove there are no hazards in your design. Estimate the maximum clock frequency given that the setup time and gate delays of the higher speed logic are 1 nanosecond.

Q. An FPGA is to be used to generate a sequencer for disco lights. The FPGA will receive some pulses that are roughly in time with the music (generated by a crude analog circuit or else by a maniac with some switches) and generate 10 or so output signals to control the lights. There are also some control inputs to select the basic mode (blackout, strobes, dark, light, pattern select etc.). Sketch a block diagram of the whole system and sketch the circuitry of the FPGA. Is using an FPGA a good idea?

Q. What is semi-custom design ? When would you use semi-custom? Why might you use an FPGA to prototype a design that is later to be built in semi-custom? What cost difference would you expect between an FPGA and a semi-custom versions and why? How might your partition be affected by the amount of RAM needed?

Q. A design is required for a wand computer toy. The wand will have a barcode reader which will scan barcodes found on groceries and other goods and accumulate credit inside the wand in an unforgeable, non-volatile way. A group of young friends, each with their own wand, will be able to compete against each other, gathering points from each item scanned, with more points for rare/valuable goods. Design the wand, including its display and controls and a method for wands to connect to each other to exchange/swap barcode sightings.