

- **Embedded System Overview**
- **Classification of Embedded Systems**
- **Hardware and Software in a System**
- **Purpose and Application of Embedded Systems**

1.1 Embedded System Overview

A. Introduction

An embedded system is a combination of computer hardware and software and may have additional mechanical or other parts designed to perform a specific function. Embedded system is a kind of computer system that performs a dedicated function and/or is intended for use with a specific embedded software application. Embedded systems are devices which are used to control, monitor or assist the operation of an equipment, machinery or plant. An embedded system may be defined as a computer system designed for specific control functions within a larger system, often with real time computing constraints. It may be embedded as a part of a complete device often including hardware and mechanical parts.

B. Characteristics

Single functioned: As embedded systems are designed for specific control functions, it usually executes a specific program to carry out the specific function repeatedly. In some cases, exceptions may occur but in general all embedded systems are supposed to carry out single specified function. One case is where an embedded system program is updated with a newer version. A second case is when several programs are swapped in and out of a system due to some constraint.

Tightly constrained: Feasibility and utility of the embedded system are measured in terms of cost, size, performance, power and other parameters. These all parameters are referred as design metric. All computing systems are constraint with design metric but it is more tightly constraint in embedded systems. Embedded systems in general must be economic, small possible size, fast enough to process data in real time and must consume minimum power to extend battery life or prevent the necessity of a cooling fan.

Reactive and real time: In general embedded systems must continually react to changes in the system's environment and must compute certain results in real time without delay. A delay in computation and slow response may result a failure in the operation of the system.

C. Components

An embedded system has mainly three components hardware, application software and real time operating system (small scale embedded system may not require RTOS).

Hardware: It represents the physical component of the system which interacts with each other to perform the specific task. Processor, RAM, ROM, ADC, DAC, Timers, Ports etc are some of the hardware components of the embedded system.

Application Software: The application software may perform concurrently the series of tasks or multiple tasks. Generally they are written in Assembly, C, C++, Java etc.

RTOS: It supervises the application software and provides a mechanism to let the processor run a process as per scheduling and do the context-switch between various tasks. RTOS defines the way the system works and sets the rules during the execution of the application software. Win CE, VxWorks, Embedded Linux etc.

D. Design Metrics

A design metric is a measurable feature of a system's implementation. Some of the commonly used metrics include:

- **NRE cost** (nonrecurring engineering cost): It represents the monetary cost for designing the system. Large number of units can be produced without any additional design cost. Since the cost doesn't occur more than once for a particular system, it is termed as nonrecurring.
- **Unit cost:** It is the monetary cost of manufacturing each unit of the system excluding NRE cost.
- **Size:** It is the physical space required by the system. For software it is measured in terms of bytes and for hardware it is measured in terms of no of gates or transistors.
- **Performance:** It represents the execution time of the system.
- **Power:** The amount of power consumed by the system, which may determine the lifetime of a battery, or the cooling requirements of the IC.
- **Flexibility:** The ability to change the functionality of the system without incurring heavy NRE cost.
- **Time to prototype:** The time needed to build a working version of the system, which may be bigger or more expensive than the final system implementation. It can be used to verify the system's usefulness and correctness and to refine the system's functionality.
- **Time to market:** The time required to develop a system to the point that it can be released and sold to customers.
- **Maintainability:** The ability to modify the system after its initial release.

- **Correctness:** We can check the functionality throughout the process of designing the system and we can insert test circuitry to check that manufacturing was correct.
- **Safety:** The system is supposed to cause no harm.

The Time to Market Design Metric

Introduction of an embedded system to the marketplace significantly affects the overall system profitability. The market window, period during which the product have highest sales, for products is getting shorter, so a short delay on introduction of product to the marketplace can render huge loss. Using a simplified model of revenue as shown in the figure below, we will deduce the loss of revenue that can occur due to delayed entry of a product in the market.

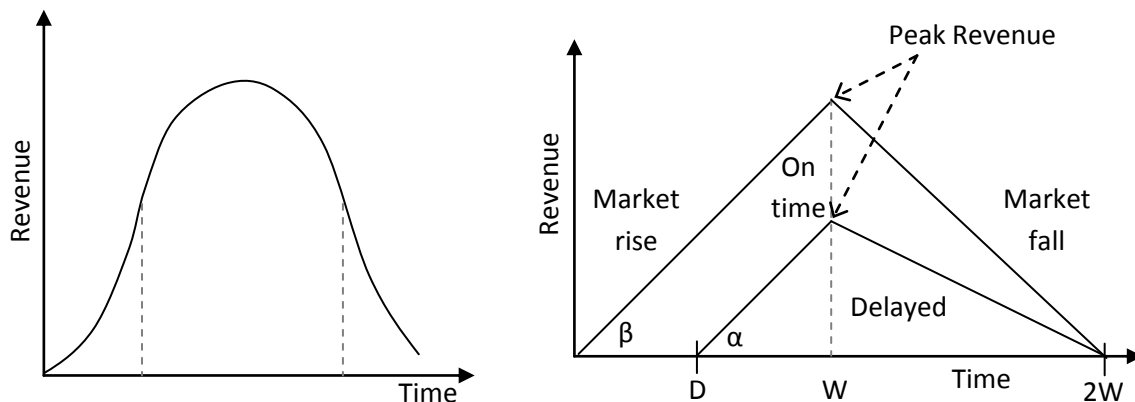


Figure 1.1: Market window and simplified revenue model for loss calculation for delayed entry

This model assumes the peak of the market occurs at the halfway point, denoted as W , of the product life. The peak is same for delayed entry. The revenue for an on-time market entry is the area of the triangle labeled *On-time*, and for delayed entry is the area of triangle labeled *Delayed*. The difference between the areas of two triangles gives the revenue loss for a delayed entry.

$$\text{Revenue Loss} = ((\text{On time} - \text{Delayed}) / \text{On time}) * 100$$

$$\begin{aligned} \text{Area of On time triangle} &= \frac{1}{2} * \text{base} * \text{height} \\ &= \frac{1}{2} * 2 * W * W * \tan \beta \quad (\text{Assuming, market rise angle is } \beta) \\ &= W^2 \tan \beta \end{aligned}$$

$$\text{Area of Delayed entry triangle} = \frac{1}{2} * (2W - D) * (W - D) * \tan \alpha$$

Assuming $\beta = \alpha$, and on solving we get,

$$\text{Revenue Loss} = (D(3W - D) / 2W^2) * 100\%$$

The NRE and Unit Cost Design Metrics

The NRE cost is the one time monetary cost of designing the system, whereas the unit cost represents the monetary cost of manufacturing each copy of the system, excluding NRE cost.

$$\text{Total Cost} = \text{NRE cost} + \text{unit cost} * \# \text{ of units}$$

$$\text{Per-product Cost} = \text{total cost} / \# \text{ of units}$$

$$= \text{NRE cost} / \# \text{ of units} + \text{unit cost}$$

The larger the volume, the lower the per-product cost, since the NRE cost can be distributed over more products. The per-product cost of the product approaches the unit cost for very large volume.

For Example, let us consider products using three different technologies; Technology A with NRE Cost of \$2000 and unit cost of \$100, B with NRE cost of \$30000 and unit cost \$30, and C with NRE cost of \$100000 and unit cost \$2

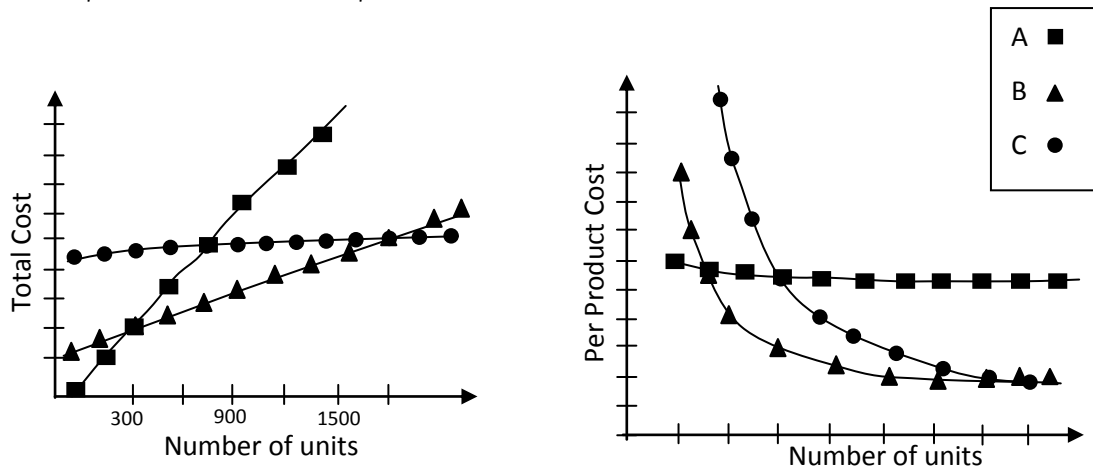


Figure 1.2: Plot of total and per product cost as a function of volume for products A, B and C

The plot on the left, total cost versus volume, shows that line of technology A and B intersect at 300 which implies technology A yields the lowest total cost for low volumes, less than 300 units. Technology B Yields lowest total cost for volume between 300 and 1800, since the line of technology B and C meets at 1800. Furthermore, technology C yields the lowest cost for volumes above 1800. The plot on the right, per-product cost versus volume, illustrates how larger volume will amortize NRE costs resulting in lower per-product costs. For example, for technology C with a volume of 50000, the per-product cost is \$4 but if we consider a volume of 200000, the per-product cost will reduce to \$2.5.

The Performance Design Metric

Performance of a system is a measure of duration the system takes to execute our desired tasks. Though the performance of a system is governed by clock frequency or instructions per second, the main measures of performance are:

Latency or response time: the time between the start of the task's execution and the end.

Throughput: the number of tasks that can be processed per unit time.

Speedup is a common method of comparing the performance of two systems. The speedup of system A over system B is determined by:

Speedup of A over B = performance of A/performance of B

E. Example of an Embedded System – A Digital Camera

A digital camera can be taken as embedded system as it performs only a single function of capturing image. It is tightly constrained as it is affordable, portable, and consumes less power. And as it is fast enough to process numeral images in milliseconds, it exhibits real time feature. But however, a simple digital camera may not possess high degree of reactive attribute. On the contrary, few contemporary digital cameras are capable of detecting human expressions.

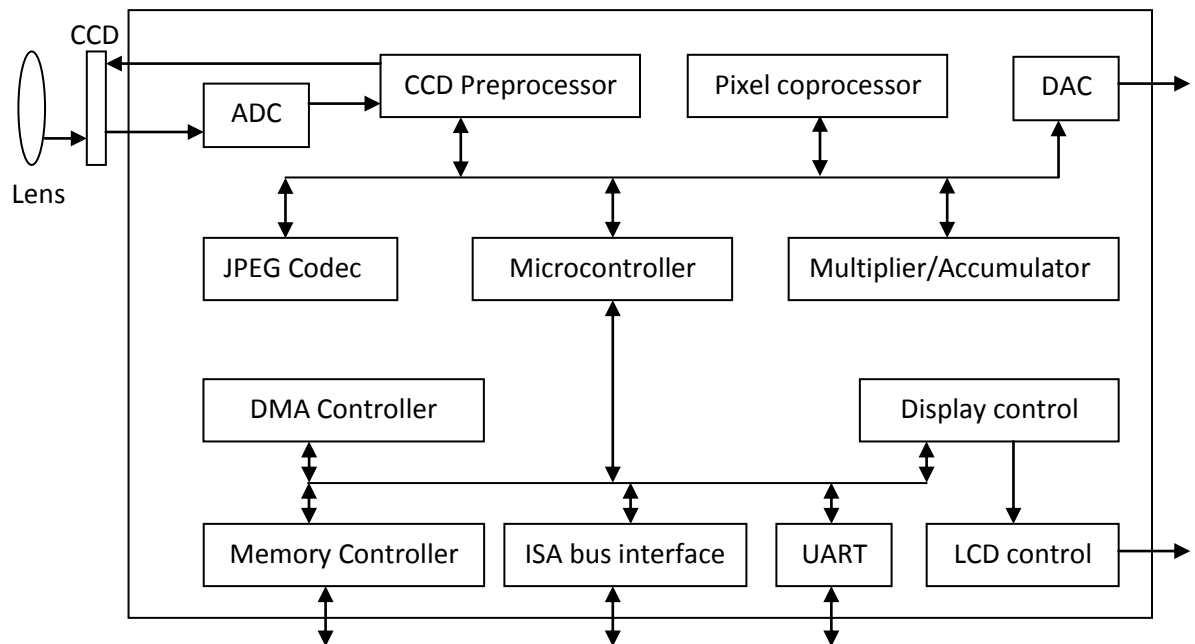


Figure 1.3: Block diagram of a typical digital camera

The different components of the digital camera chip are listed below:

- **Charge - Coupled Device (CCD):** It contains an array of light sensitive photocells that capture an image. A CCD is an integrated circuit etched onto a silicon surface forming light sensitive elements called pixels. Photons incident on this surface generate charge that can be read by electronics and turned into a digital copy of the light patterns falling on the device. Photons striking a silicon surface create free electrons through photoelectric effect. And the electrons are gathered in the place where they are generated and count them in some manner to create an image. It is accomplished by positively biasing discrete areas to attract electrons generated while the photons come onto the surface.
- **Analog to digital converter:** It converts analog images to digital images. The main job is to classify the voltages of the pixels into levels of brightness and to assign each level to a binary number, consisting of zeroes and ones. Most digital cameras use at least an 8-bit ADC, which allows for up to 256 values for the brightness of a single pixel.
- **Digital to analog converter:** It converts digital images to analog images.
- **CCD preprocessor:** it provides commands to the CCD to read the image. Some other functions of preprocessor are over scan correction and trimming, bias removal, dark current removal, and flat fielding.
- **JPEG (Joint Photographic Expert Group) Codec:** It compresses and decompresses an image using the JPEG compression standard which enables compact storage of images in the limited memory of the camera.
- **Pixel coprocessor:** It helps in rapidly displaying images. A pixel coprocessor is required in digital cameras for displaying images directly or after operations such as rotate right, rotates left, rotate up, rotate down, shift to next, shift to previous.
- **Memory Controller:** It controls the access to a memory chip. It is a digital circuit which manages the flow of data going to and from the main memory.
- **DMA (Direct Memory Access) Controller:** It enables direct memory access by other devices while the microcontroller is performing other functions.
- **UART (Universal Asynchronous Receiver and Transmitter):** It enables communication with a PC's serial port for uploading video frames. It is the microchip with programming that controls a computer's interface to its attached serial devices.
- **ISA (Industry Standard Architecture) bus interface:** It enables faster connection with a PC's ISA bus. Most PC's have an ISA slot on the main board that accepts either an 8 bit or a 16 bit ISA printed circuit card.

- **LCD Control:** It controls the display of images on the camera's liquid crystal device.
- **Display Control:** It also controls the display of images on the camera's liquid crystal device.
- **Multiplier/Accumulator:** It assists with certain digital signal processing.
- **Microcontroller:** It is the processor that controls the activities of other circuits within the system.

1.2 Classification of Embedded Systems

Classification based on Generation

A. First Generation

Embedded systems were designed using 8 bit microprocessors or 4 bit microcontrollers. The hardware circuits were simple and the firmware was developed using assembly code. Motor controller using 8085 can be an example of first generation embedded system.

B. Second Generation

In second generation, the systems were built using 16 bit microprocessors and 8/16 bit microcontrollers. More complex and powerful instructions were available for the designer to work with. Some systems involved embedded operating systems for their operation. Data Acquisition Systems can be an example of second generation embedded systems.

C. Third Generation

The systems were designed with more advanced processor technology in the form of 32 bit processors and 16 bit microcontrollers. Along with complex and powerful instruction sets, instruction pipelining was introduced for better performance. Dedicated embedded real time operating system implementation was another important feature in this generation. Also, the concept of application specific processors like Digital Signal Processors (DSP) and Application Specific Integrated Circuits (ASIC) came into existence.

D. Fourth Generation

Fourth generation was marked with the advent of System on Chips (SoC), reconfigurable processors and multicore processors. These embedded systems used high performance real time embedded operating systems for its operation.

Classification Based on Complexity and Performance

A. Small Scale Embedded Systems

These systems are designed with a single 8 or 16 bit microcontroller (8051 family, PIC16F8X, Hitachi H8). They have little hardware and software complexities and involve board level design. They may be battery operated. While developing embedded software for these system, an editor, assembler and cross assembler specific to the microcontroller or processor are used as the main programming tool. Usually C language is used for developing these systems. To develop such systems, the size requirement of the software must not exceed the available memory while the hardware design must be done in such a way that the power dissipation must be limited when the system is running continuously. Automatic vending machine, stepper motor controller for a robotics system etc can be the examples of small scale embedded systems.

B. Medium Scale Embedded Systems

These systems are designed with a single or few 16 or 32 bit microcontrollers (8051MX, PIC16F876) or DSPs or Reduced Instruction Set Computers (RISCs). It may also employ the readily available single purpose processors and IPs for various functions, for example: bus interfacing, encryption, deciphering and so on. These systems have both hardware and software complexities. For software design, the programming tools used is RTOS, Source code engineering tool, Simulator, Debugger and Integrated Development Environment (IDE). Software tools also provide the solutions to the hardware complexities. Some of the examples of medium scale embedded systems are Computer networking systems, signal tracking system etc.

C. Sophisticated Embedded Systems or Large Scale Embedded Systems

These systems have enormous hardware and software complexities and may need scalable processors or configurable processors and programmable logic arrays. They are used for cutting edge applications that need hardware and software co-design and integration in the final system. They are constrained by the processing speeds available in their hardware units. Certain software functions are implemented in the hardware to obtain additional speed by saving time. Some of the functions of the hardware resources in the system are also implemented by the software. These systems generally implement high performance real time operating system. Development tools for these systems may not be readily available at a reasonable cost or may not be available at all. In some cases, a compiler or retargetable compiler might have to be developed for these. (A retargetable compiler is one that configures according to the given target configuration in a system). Embedded System for wireless LAN & for convergent technology devices is one of the sophisticated embedded systems.

1.3 Hardware and Software in a system

In an embedded system, hardware of the system represents the single purpose processor whereas the software represents the general purpose Processor.

A. Single Purpose Processor

Single Purpose processor is a digital circuit designed to execute exactly one program. It does not require a program memory since the program to run on the processor is fixed (only one) and it can be built into the digital circuit. The datapath contains only the essential components for the specified task. JPEG codec, Display controller, DMA controller etc can be taken as the examples of single purpose processor.

Design metric benefits of single purpose processor

- Performance may be fast, size and power may be small
- Unit cost may be low for large quantities

Design metric drawbacks of single purpose processor

- Design time and NRE costs may be high, flexibility low, unit cost high for small quantities
- Performance may not match general purpose processors for some applications.

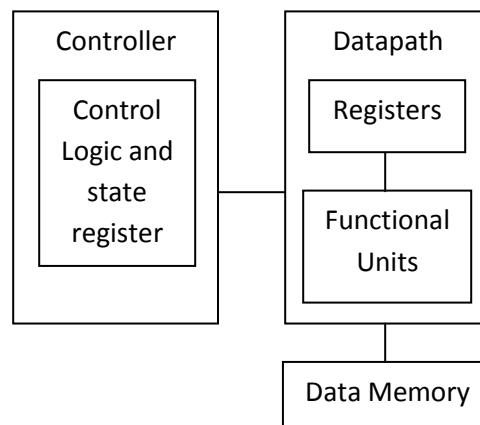


Figure 1.4: Block Diagram of Single Purpose Processor

B. General Purpose Processor

The required functionality is carried out by programming the processor's memory. In this context, a programmable device is built that is suitable for a variety of applications. Microprocessors are the examples of general purpose processor. Such processors include:

- Program memory: The program cannot be built into the digital circuit in general purpose processors since the program likely to run on the processor will be unknown.
- General Datapath: The datapath must be general enough to handle a variety of computations, so such datapath typically has a large register file and one or more general purpose arithmetic logic units (ALUs).

Design metric benefits of general purpose processor

- Time to market and NRE costs are low because designer must only write a program but does not have to deal with any digital design.
- Flexibility is high because changing functionality requires changing only the program.
- Unit cost may be low in small quantities.
- Performance may be fast for computation intensive applications.

Design metric demerits of general purpose processors

- Unit cost may be relatively high for large quantities, since in large quantities we could design our own processor and amortize NRE costs.
- Size and power may be large due to unnecessary processor hardware.
- Performance may be slow for certain applications.

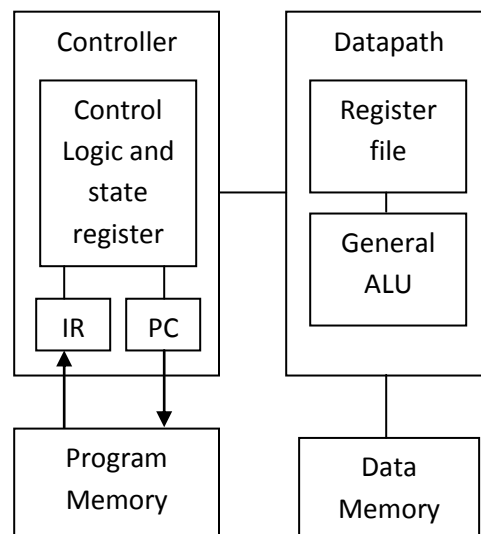


Figure 1.5: Block Diagram of General Purpose Processor

C. Application Specific Processors

Application specific processors are programmable processors optimized for a particular class of applications. It generally includes program memory, optimized datapath and special functional units. These processors provide optimum level of performance maintaining appropriate size and power consumption. Microcontrollers for controlling application and digital signal processors (DSPs) for huge data processing application are examples of application specific processors.

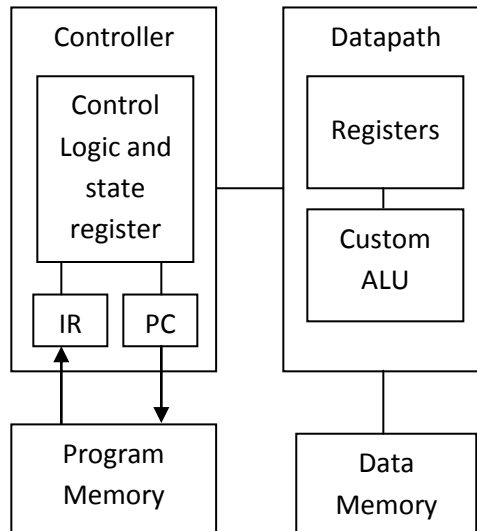


Figure 1.6: Block diagram of Application Specific Processors

1.4 Purpose and Application of Embedded Systems

Purpose of Embedded Systems

A. Data collection

In embedded systems, the data is collected from other external devices for storage, analysis, manipulation or transmission. Data may be in analog or digital form. Systems working with digital data require analog to digital converters if the collected data is in analog form. The collected data can be used for meaningful purpose based on the functionality of the embedded system. For instance, a digital camera collects data, stores it and finally provides graphical representation of data in the form of captured image.

B. Data communication

An embedded system is required to connect two or more devices which may be at close vicinity or at remote location. The communication between devices can be done via wired line medium or wireless medium. Embedded systems are incorporated with different wireless modules or wire-line

modules for communication purpose. Bluetooth, ZigBee, Wi-Fi, and GSM are few wireless modules. And for wire-line purpose, an embedded system may have RS-232, SPI, I2C, USB and other serial and parallel protocols. Some embedded systems like network hubs, routers, etc act as mediators in data communication and provide various features including data security.

C. Data Processing

The collected data in embedded system is subjected to some sort of processing for which embedded systems are attributed with data processing modules. Speech coder, audio video codec, etc can be few examples of data processing unit. Data processing includes the manipulation of data for appropriate purpose.

D. Monitoring

Many embedded systems are incorporated with sensors to check the state of the different parameters. The parameters can be current, voltage, temperature, humidity, etc which are continuously monitored and appropriate processing or controlling of devices is done. However, the value of the parameters cannot be controlled by the system itself. The values of parameters are used for some controlling purpose or for some graphical representation purpose or simply stored for further analysis and processing.

E. Control

For control purpose, actuators along with sensors are present in the embedded systems. The sensor connected in input port detects the change in the desired parameter and the actuators at output ports are controlled accordingly to implement the desired functionality. Electric Motors are examples of actuators. In an object avoiding robot, ultrasonic sensor senses the presence of certain kind of object and the motor is rotated accordingly to avoid the collision.

F. Application specific user interface

To provide the better user interface based on application has been one of the concerns of contemporary embedded systems. Keypads, simple LCD modules, speakers, etc are basic and common interface for users. However, sensitive touch pad along with high definition display has been the sophisticated interface implemented in current scenario.

Applications of Embedded Systems

1. Household appliances: Microwave ovens, Television, DVD players and recorders.
2. Consumer electronics: cell phones, cameras, video games

3. Office automation: fax machines, printers, scanners
4. Business equipment: alarm systems, card readers,
5. Automobiles: engine controller, fuel injection, antilock brakes.
6. Networking: Modem, Network cards, Network switches and routers
7. Medical equipments
8. Aerospace research
9. Integrated systems in aircrafts and missiles
10. Industrial and Military applications

- **Combinational Logic**
- **Sequential Logic**
- **Custom Single-Purpose Processor Design**
- **Optimizing Custom Single-Purpose Processors**

2.1 Combination Logic

Combinational circuit is a digital circuit whose output is purely a function of its present inputs. Combination logic circuits are made up from basic gates or universal gates that are combined or connected together to produce more complex switching circuits. In general, logic gates are the building blocks of combinational logic circuits. It has no memory block. Some of the examples of the combinational circuits are decoder, multiplexer, adder, ROM etc.

CMOS Transistors

A transistor, which acts as a simple on/off switch, is the basic electrical component in digital system. More abstract components, logic gates, are formed with the combinations of transistors. In Complementary Metal Oxide Semiconductor (CMOS), the gate voltage controls the flow of current from source to drain. The nMOS conducts when gate is at high voltage (5v) whereas pMOS conducts when gate is at low voltage (0v). The symbol for nMOS and pMOS is shown in the figure 2.1.

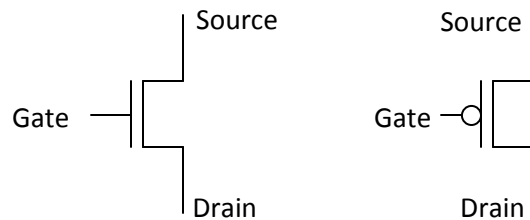


Figure 2.1: nMOS and pMOS transistors

Different gates and boolean functions can be realized using nMOS and pMOS.

A. Inverter: When $x = 0$, transistor T1 conducts but T2 does not. So, output is logic 1. And when $x = 1$, T2 conducts but T1 does not.

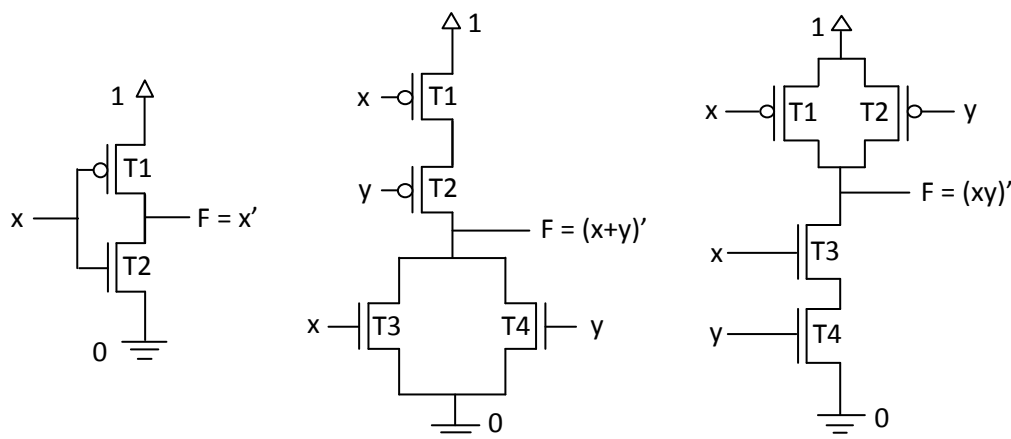


Figure 2.2: Inverter, NOR and NAND gate using nMOS and pMOS

B. NOR Gate

When $x = 0$ and $y = 0$, then T1 and T2 conduct but T3 and T4 don't. So, F is connected to Vcc.

When $x = 1$ and $y = 0$, then T2 and T3 conduct but T1 and T4 don't. So, F is connected to ground.

When $x = 0$ and $y = 1$, then T1 and T4 conduct but T2 and T3 don't. So, F is connected to ground.

When $x = 1$ and $y = 1$, then T3 and T4 conduct but T1 and T2 don't. So, F is connected to ground.

When atleast one of the two inputs is high then the output is connected to ground. And when both inputs are low then the output is connected to Vcc.

C. NAND Gate

When $x = 0$ and $y = 0$, then T1 and T2 conduct but T3 and T4 don't. So, F is connected to Vcc.

When $x = 1$ and $y = 0$, then T2 and T3 conduct but T1 and T4 don't. So, F is connected to Vcc.

When $x = 0$ and $y = 1$, then T1 and T4 conduct but T2 and T3 don't. So, F is connected to Vcc.

When $x = 1$ and $y = 1$, then T3 and T4 conduct but T1 and T2 don't. So, F is connected to ground.

When atleast one out of two inputs is low then the output is connected to Vcc. And when both inputs are high then the output is connected to ground.

Basic Logic Gates

- The **NOT** (Inverter) gate simply complements the input.
- The **AND** gate outputs 1 if and only if all of its inputs are 1.
- The **OR** gate outputs 1 if at least one of its inputs is 1.
- The **XOR** (exclusive-OR) gate outputs 1 when only one of its inputs is 1.
- The **NAND**, **NOR** and **XNOR** gates outputs the complement of AND, OR and XOR respectively.

Basic Combinational Logic Design

In Combinational Design, output is purely a function of its present inputs and has no memory of past inputs. We can use basic logic gates to design combinational circuits. In such design, outputs are described in terms of inputs.

General steps for combinational Logic Design

- The description is translated into a truth table with all possible combinations of input values.
- The input values lies on the left of the truth table and the corresponding output values of the inputs lies on the right of the truth table.
- For each output, we have to derive the equations. The equation may contain number of combinations of the inputs. The number of combinations depends on the number of high (1)

value on each column of the output. Rows of the inputs are used to derive the equation corresponding to the high output of the column. And the equation must be further minimized.

- Another way to derive minimized equation directly is by using k-map. It is always better to use k-map unless the design is too simple (when the output column consists of only one high value).
- The final equation is translated to an equivalent circuit diagram using logic gates.

Combinational Logic Design Example

Example 1: In an alarm system of a bank, three sensors are implemented and the alarm is triggered when at least two sensors detect the change. Assuming sensors to output digital values, design a combinational logic circuit for alarm system.

Solution: Let a, b, c represent the three sensors and y represents the buzzer for alarm. The output y should be high when two or more than two inputs are high. The truth table and its corresponding combinational design are shown below.

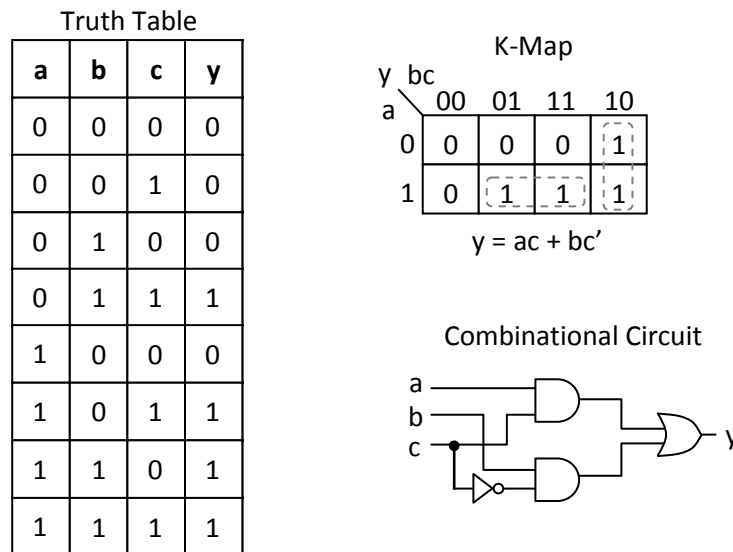


Figure 2.3: Truth table, K-map, and combinational circuit for bank alarm system

RT-Level Combinational Components

Register-transfer or RT level components are generally used when the design of the circuit becomes complex. As the number of input increase, the complexity of the design increase. One of the ways to reduce design complexities is by using RT-level components. Multiplexers, decoder, adder are the examples of RT-Level Components.

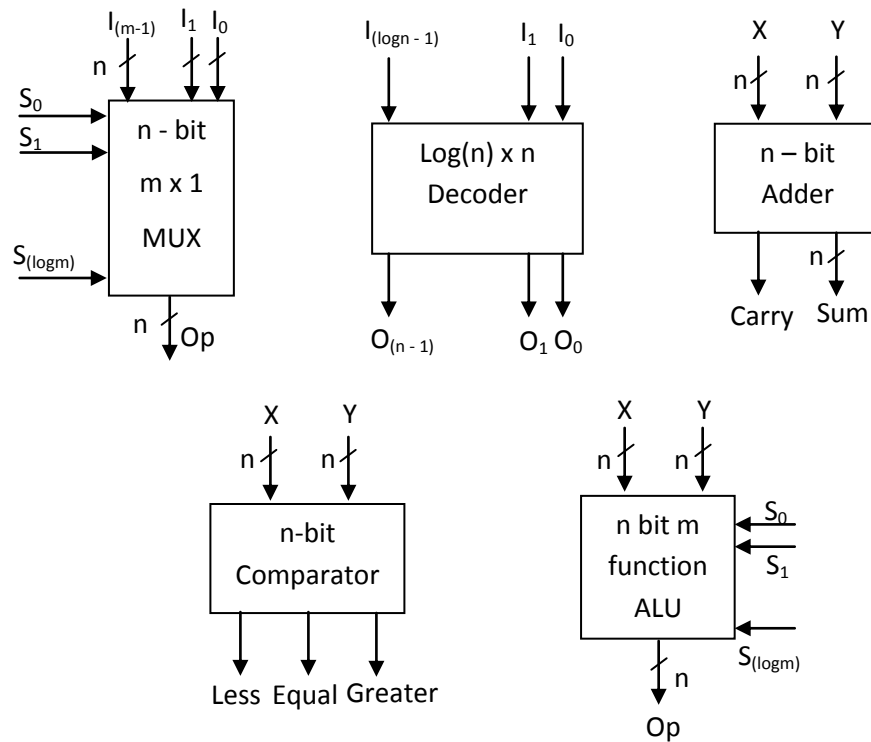


Figure 2.4: Few commonly used RT-Level Combinational Components

- A **multiplexer** allows only one of its data inputs to pass through to the output. For $m \times 1$ multiplexer there are m data inputs and one data output with $\log_2 m$ select lines. The value of select line determines which input data to pass through to the output. It can be used for parallel to serial conversion.
- A **decoder** allows exactly one of the output lines to be high at a given time for a particular input. For n input lines there will be 2^n output lines. A decoder can be used for coding the addressing lines in the memory. It can be used to convert binary to a suitable form.
- An **adder** is used to add two n – bit inputs producing an n -bit sum along with a carry of 1 bit.
- A **comparator** allows to compare two n -bit binary inputs, generating the corresponding output based on whether one input is less than, equal to, or greater than another input.
- An **arithmetic-logic unit (ALU)** performs variety of arithmetic and logic functions on its n – bit inputs. The select line is used to select which function is to be carried out. If there are 2^m functions that can be done by ALU then there must be at least m select lines.
- A **shifter** is another example which is used to shift the bits of the input right or left. It can be used as a divider or multiplier. For example shifting 0110 (6) to the right would give 0011 (3).

2.2 Sequential Logic

A sequential circuit is a digital circuit whose outputs are a function of not only the present inputs but also the past inputs. The output of a sequential logic depends on its present internal state and the present inputs. Hence a sequential logic circuit has some kind of memory. Logic gates and flip flops are the basic building blocks of sequential logic circuits. Flip flop is an example of sequential logic circuit.

A flip flop stores a single bit. The different types of flip flops are listed below.

- **D-flip flop:** It has two inputs D and clock, when clock is high, value of D is stored in flip flop and same will be the value of the output Q. When clock is low, previously stored bit is maintained ignoring the value of input D.
- **SR flip flop:** It has three inputs S (set), R (reset) and clock. When clock is low, the previously stored bit is maintained ignoring the values of input at S and R. When clock is high, the output varies with inputs S and R. If S is high, the output Q will be high and high bit (1) will be stored by the flip flop. If R is high, then low bit (0) will be stored. The output will not change if both the inputs are low but the undefined condition will occur if both the inputs are high.
- **JK flip flop:** Its operation is similar to that of SR flip flop but when both the inputs J and K is high, the stored bit toggles either from high to low or low to high.

Flip flops are generally designed to be edge triggered to prevent the unexpected behavior from signal glitches, the inputs are checked either at the rising edge or falling edge of the clock. Glitches represent an undesirable transition that occurs before the signal settles to its intended value.

RT Level Sequential Components

Generally RT level sequential components are required for designing complex sequential systems. Counters and registers are examples of RT level sequential components.

- A **register** stores n bits from its n bit data input which also appears at its output. A register usually has at least two control inputs, clock and load. For a rising edge triggered register, the inputs are only stored when load is high and clock is rising from 0 to 1. Another control input clear may be used to resets all bits to 0 regardless of the value of input. Since all n bits of the registers can be stored in parallel, we refer this type of register as a parallel load register.
- A **shift register** stores n bits from its one bit data input with at least two control inputs clock and shift. When clock is rising and shift is 1, the nth bit of input is stored in the (n-1)th bit, and (n-

1)th bit of input is stored in the (n-2)th bit and so on down to the second bit being stored in the first bit. The first bit is shifted out appearing as an output bit. It has one bit output and the input must be shifted into the register serially.

- A **counter** is a register that adds binary 1 to its stored binary value. In general, a counter has a clear, count and load as a control inputs. Clear resets all stored bits to 0 and a count input enables incrementing on each clock edge. It often has parallel load data input and associated load control signal. A common counter feature is both up and down counting which required an additional control input to indicate the count direction.

A small triangle in the block represents the clock input for any sequential logic. Control inputs in sequential logic can be either synchronous or asynchronous. A synchronous input value only has an effect during a clock edge while an asynchronous input value affects the circuit independent of the clock. Clear control lines are asynchronous inputs while load, shift count control lines are synchronous inputs.

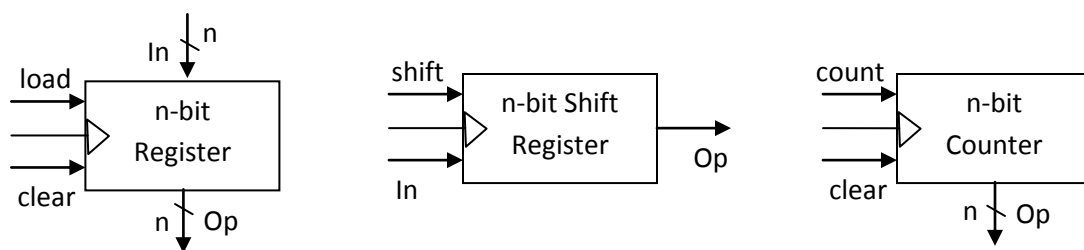


Figure 2.5: RT-Level Sequential Components

Sequential Logic Design

1. Translate the problem description to a state diagram, also called a finite state machine (FSM).
2. In FSM, each circle represents a state where desired output values are listed next to each. Whereas the input conditions which cause a transition from one state to another are listed next to each arc.
3. Draw an implementation model which implements the FSM using a state register to store the current state and combinational logic to generate the required output values and next state.
4. Assign each state a unique binary value, and create a truth table for the combinational logic. The external inputs and the bits coming from the state registers are fed to the combinational logic as inputs. Whereas, the external output values along with the state bits to be loaded into the state register acts as the output of the combinational logic.

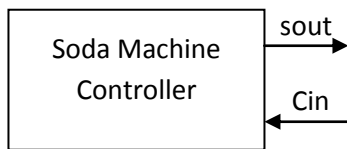
5. The output values change only with the current state, so we list the external output values only for each possible state, regardless of the change in external input values.
6. Now, we can have a truth table, with the help of which we can proceed with combinational design by generating minimized output equations using k-map. And finally, drawing the combinational logic circuit.

Sequential Logic Design Example

Example 1: Design a soda machine controller, given that a soda costs 75 cents and your machine accepts quarters only. Draw a black-box view, come up with a state diagram and state table, minimize the logic, and then draw the final circuit.

Solution: The coin must be entered three times to get a soda out of the machine. Throughout the design, Cin represents the coin input and sout indicates the soda output whereas Q1, Q0 represent current state and I1, I0 represent next state.

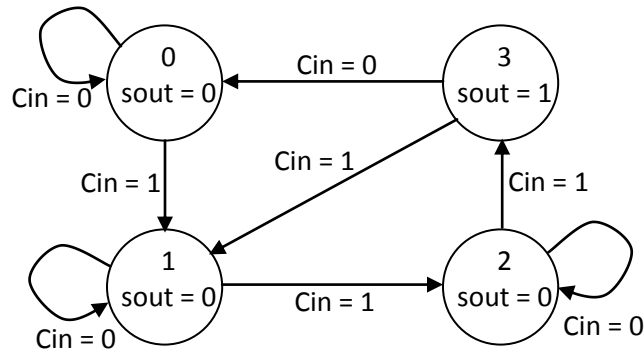
A. Black Box View



C. State Table

Inputs			Outputs		
Q1	Q0	Cin	I1	I0	Sout
0	0	0	0	0	0
0	0	1	0	1	
0	1	0	0	1	0
0	1	1	1	0	
1	0	0	1	0	0
1	0	1	1	1	
1	1	0	0	0	1
1	1	1	0	1	

B. State Diagram



D. K-map

I1	Q1Q0			
	00	01	11	10
Cin 0	0	0	0	1
1	0	1	0	1

$$I1 = Q1Q0' + Q1'Q0Cin$$

I0	Q1Q0			
	00	01	11	10
Cin 0	0	1	0	0
1	1	0	1	1

$$I0 = Q1'Q0Cin' + Q1Cin + Q0'Cin$$

sout	Q1Q0			
	00	01	11	10
Cin 0	0	0	1	0
1	0	0	1	0

$$sout = Q1Q0$$

E. Combinational Circuit

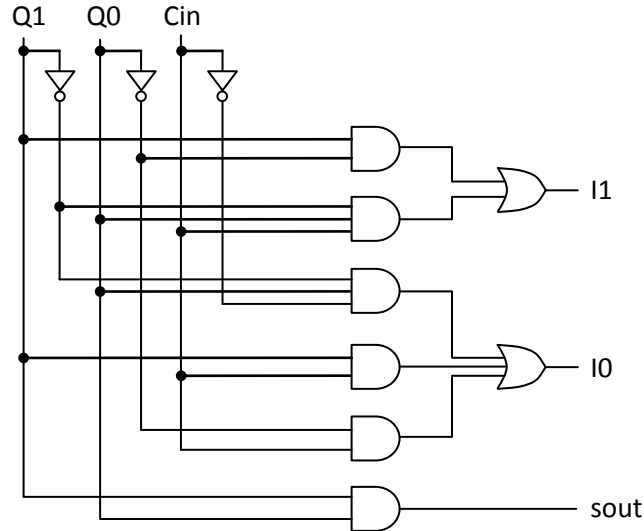


Figure 2.6: Soda machine controller design

2.3 Custom Single-Purpose Processor Design

A basic processor consists of a controller and a datapath.

Datapath

- It stores and manipulates a system's data.
- It contains register units, functional units and connection units like wires & multiplexors.
- The datapath can be configured to read data from particular registers, feed that data through functional units configured to carry out particular operations like add or shift, and store the operation results back into particular registers.
- Examples of data include binary numbers representing external conditions like temperature or speed, characters to be displayed on a screen.

Controller

- It sets the datapath control inputs, like register load and multiplexor select signals, of the register units, functional units, and connection units to obtain the desired configuration at a particular time.
- It monitors external control inputs as well as datapath control outputs, known as status signals, coming from functional units, and it sets external control outputs as well.

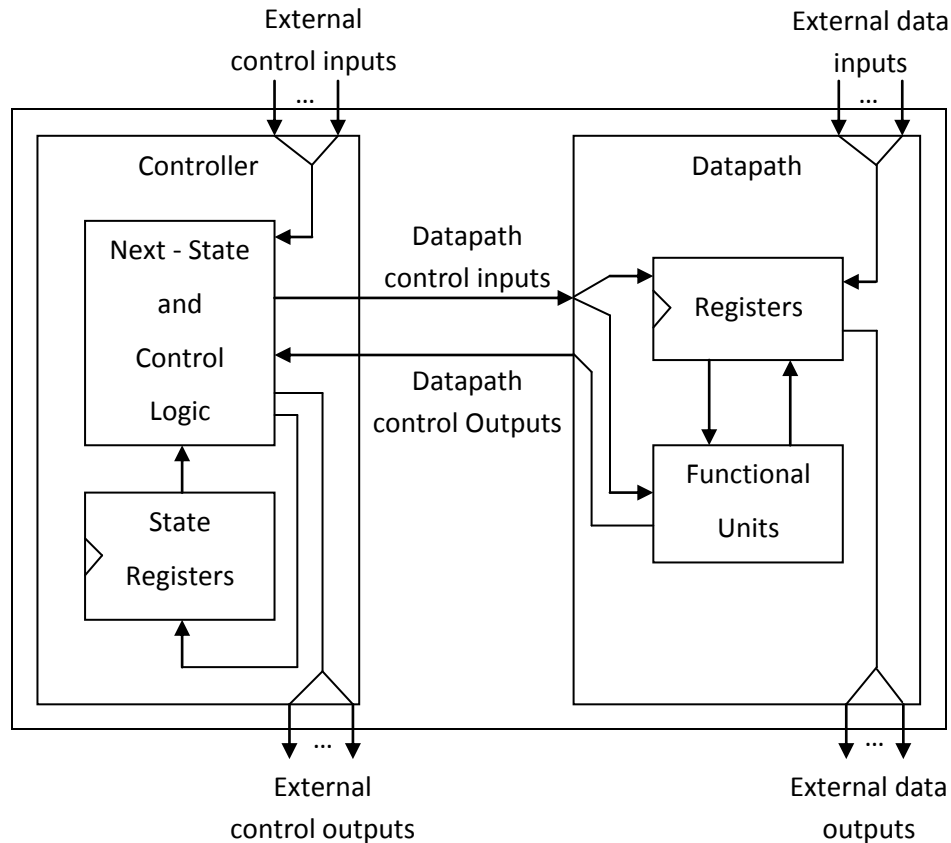


Figure 2.7: Internal View of controller and datapath of Single Purpose Processor

Steps for designing Single-Purpose Processor

- 1. Draw a Black Box Diagram:** Black box diagram is a simple box with external interfaces of a system. It generally includes input and output signals along with few control signals.
- 2. Write the functionality or program:** The functionality or program is a code which provides the solution to the defined problem.
 - The input signals are assigned to a variable.
 - Number of temporary variables may be used based on requirement.
 - The final result is assigned to the output port.
- 3. Design a Finite State Machine with Data (FSMD):** The code is converted into equivalent complex state diagram which is known as Finite State Machine with Data. In FSMD, Templates are used to represent various constructs of program. The templates for assignment, branch statement and loop statement are discussed below.

- **Assignment Statement:** For this statement, a single state is used with statement representing its action. Generally, a single arrow is used to connect to next state. The template used for statement $C = A + B$ is shown as an example.

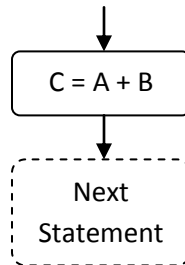


Figure 2.8: Template for assignment statement

- **Branch Statement:** It can be represented by using condition state C, join state J, and few other states in between C and J state. State C and State J are with no actions, left empty. But states between C state and J state contain actions. Its template can vary depending on number of conditions defined in the problem. However, for each true condition, there can be several states representing actions. Conditions are written along side with the arrow that connects the C state and states of each branch. Last states of each branch are connected to the J state.

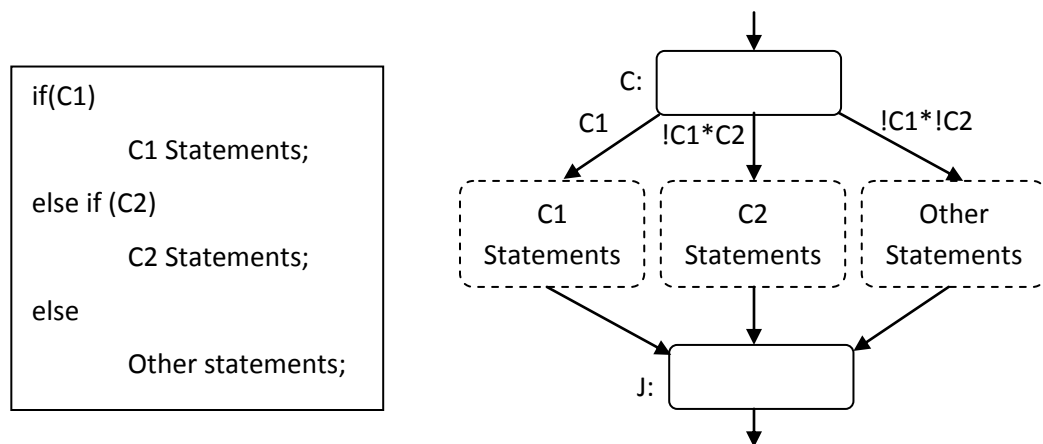


Figure 2.9: Template for branch statement

- **Loop Statement:** Its template consists of Condition State C, Join State J, and other states representing statements of loop. Condition is written alongside arrow connecting condition state and state of first statement of loop. The last state of loop is connected to the J state which is connected back to condition state. Complement condition is used

alongside arrow connecting C state and next statement outside of loop. The template for the loop statement is shown in the figure below.

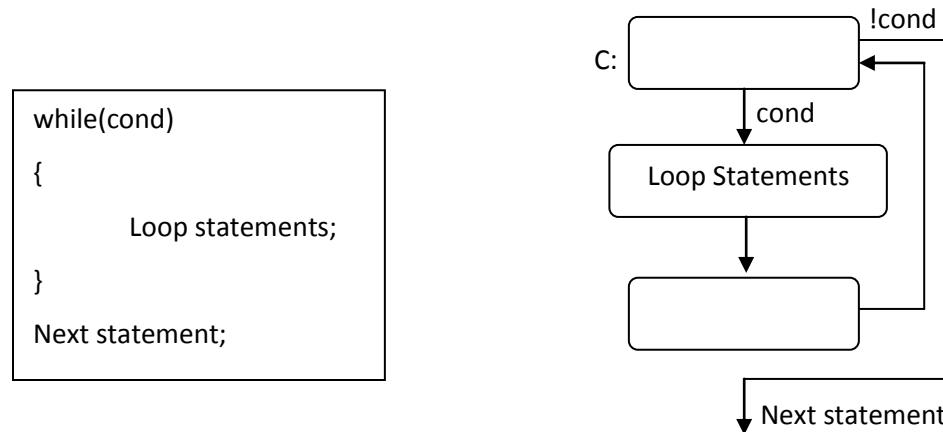


Figure 2.10: Template for loop statement

4. Build a Datapath: The datapath is build based on functionality of the system. Following steps are needed to be taken into considerations while developing a datapath.

- **Registers:** The number of registers to be used is defined by the number of variables used in the functionality. Registers are assigned to inputs, temporary variables and output.
- **Functional Units:** Blocks representing arithmetic and logical operations are defined within the datapath.
- **Connections:** The connections among ports, registers and functional units are done based on operands used in various assignments and comparison of functionality code. Appropriate multiplexor is required when the value in register can be assigned from more than one source. The sources may be an input port, a functional unit, or another register.
- **Control inputs and outputs:** Input control signals are generally required by registers and multiplexor. Register load signal is used in case of register while selection line signals for multiplexor. Control output is produced by logical units of the datapath. Each control singles are given a unique identifier.

5. Develop a Finite State Machine (FSM): The states and transitions for FSM are same as that of FSMD. However, the complex actions and conditions of FSMD are replaced by Boolean expressions using the control signals defined within datapath. For every register write operations (assignment statement, arithmetic statements), register load signal is asserted and

corresponding multiplexor selection line is activated if there are two or more sources for a given register. Also the logical operations are replaced by the control signals of its corresponding functional block.

Example 1: Design a single purpose processor that calculates the Greatest Common Divisor (GCD) of two numbers. Include FSMD, Datapath and FSM in the design.

➔ Initially, the black box view diagram is drawn and then followed by the functionality which is converted into FSMD using appropriate templates.

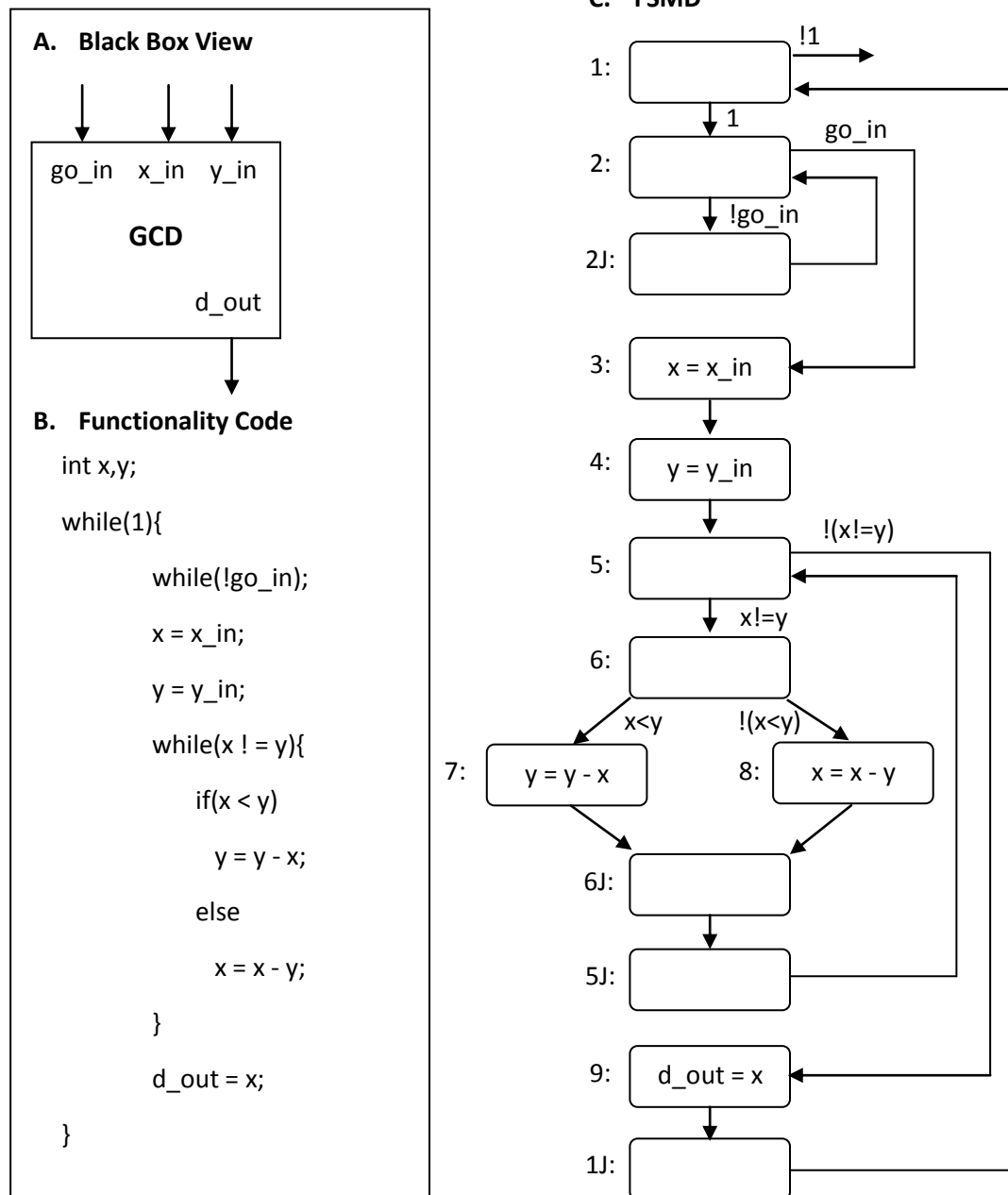


Figure 2.11: Black box view, functionality and FSMD diagram of GCD processor

D. Datapath for GCD processor:

- **Number of Registers:** Two inputs x_{in} and y_{in} assigned to variables x and y , final result assigned to d_{out} , and no other temporary variables are used. Hence, three registers x , y and d are required.
- **Functional Blocks:** The arithmetic and logical operation involved in the functionality are $x-y$, $y-x$, $x!=y$ and $x<y$. Hence, two subtractors and two comparing blocks are required.
- **Connections and MUX requirement:** The value in register x has two sources, x_{in} and $x-y$, so it requires a multiplexor of 2×1 . Similar is the case for register y . For connections, the output of registers x and y are connected to inputs of subtracting blocks and comparing blocks. Also, the line representing x_{in} and $x-y$ are connected to the inputs of mux whose output is fed to register x . Similarly, y_{in} and $y-x$ are connected to the register y through mux. And, the output of x register is connected to input of register d . All connections must be done so as to represent the corresponding operation in the functionality.
- **Control Signals:** Unique identifier for various control signals is assigned.
 - Load signal of registers: x_{ld} for register x , y_{ld} for register y and d_{ld} for register d .
 - Selection lines of multiplexor: x_{sel} for multiplexor associated with register x and y_{sel} for multiplexor associated with register y .
 - Signals from logical block: x_{neq_y} and x_{lt_y} are used for x not equal to y and x less than y respectively.

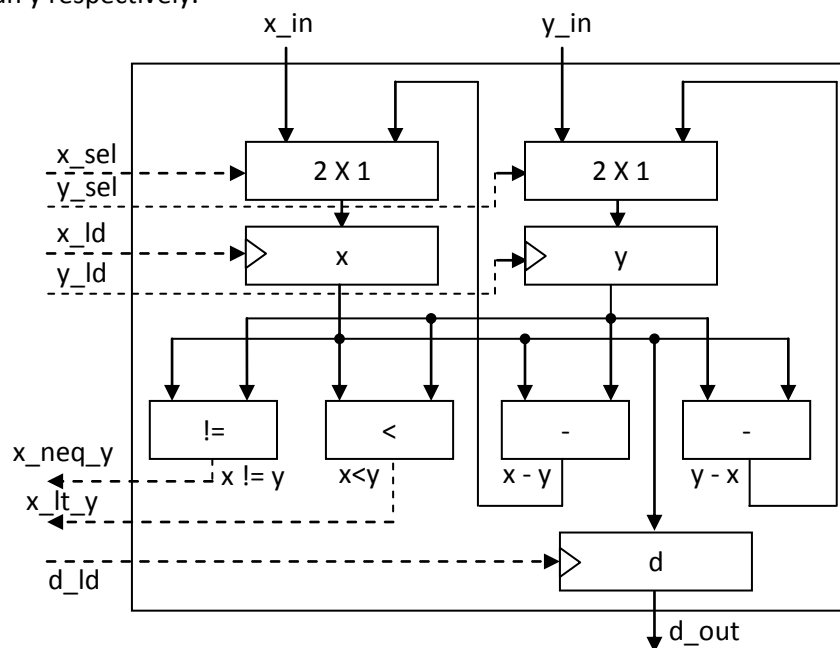


Figure 2.12: Datapath of GCD Processor

E. Finite State Machine for GCD processor

All actions and conditions are replaced by equivalent Boolean expressions as used in datapath. For example, action $x = x_{in}$ is replaced by expressions $x_{sel} = 0$ and $x_{ld} = 1$. $x_{sel} = 0$ will connect the input line x_{in} to register x and $x_{ld} = 1$ will load the value of x_{in} into x . In case of $d_{out} = x$, only $d_{ld} = 1$ is used as it has only one source and no multiplexor is used. And condition $x < y$ is replaced by x_{lt}_y . The identifiers for control signals, however, used in FSMD must match with the one that is defined in datapath.

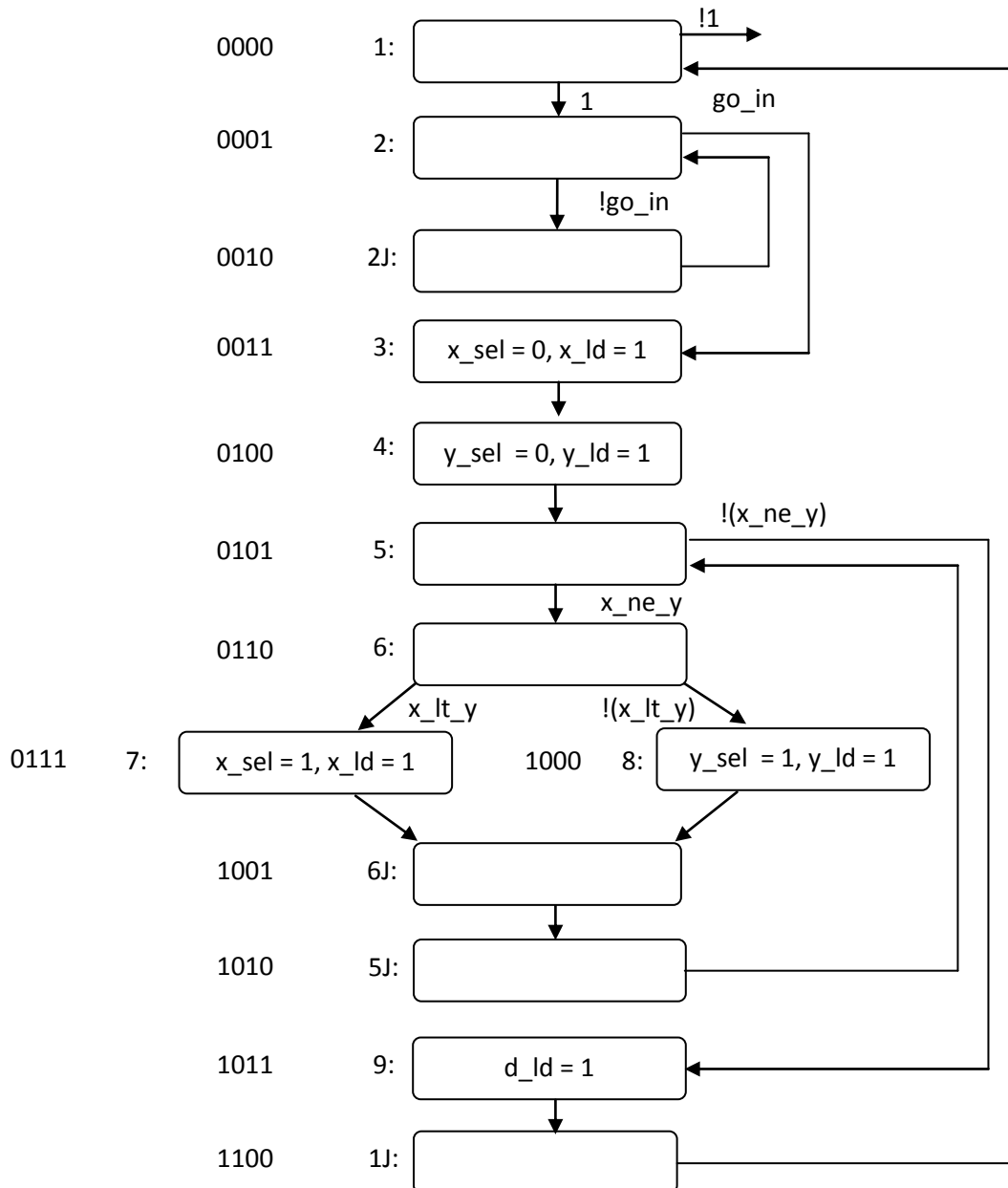


Figure 2.13: FSM of GCD Processor

2.4 Optimizing Custom Single-Purpose Processors

Optimization is the task of making design metric values the best possible. Optimization can be done by simplifying the resulting design of any system utilizing various techniques. Different states in the FSM can be removed which does nothing and are redundant. Also, we can share a component for same operations in different states and hence minizing the size of the system as well as its cost. Other various factors can be considered for optimum design but some simple optimization that can be applied are discussed further.

Optimizing the Original Program

We should analyze different program attributes and try to develop alternative algorithm that are more efficient. We can analyze the algorithm in terms of time complexity and space complexity. Number of computations can be a form of time complexity whereas the size of variables required corresponds space complexity.

Lets take the example of GCD:

<pre> int x,y; while(1){ while(!go_in); x = x_in; y = y_in; while(x != y){ if(x < y) y = y-x; else x = x-y; } d_out = x; } </pre>	<pre> int x,y,r; while(1){ while(!go_in); x = x_in; y = y_in; while(y != 0){ r = x % y; x = y; y = r; } d_out = x; } </pre>
--	---

To compute GCD of 42 and 8, it takes 9 iterations to complete the operation, x and y will take different values as (42, 8), (34, 8), (26, 8), (18, 8), (10, 8), (2, 8), (2, 6), (2, 4), (2, 2).

To compute GCD of 42 and 8, it takes 3 iterations to complete the operation, x and y will take values as (42, 8), (8, 2), (2, 0). If y is greater than x, it will take 4 iterations, one more than previous.

Optimizing the FSMD

Each state in an FSMD is assigned with operations from the desired program; this process is also termed as scheduling. The scheduling process can be improved by following methods.

- **Merge States:** States with independent operations can be merged.
- **Eliminate State:** States with constants on transitions can be eliminated since transition to be taken will be fixed as defined by constants. And some states without any operation can also be eliminated.
- **Separate States:** States which require complex operations can be broken into smaller states to reduce hardware size.

Considering the example of GCD:

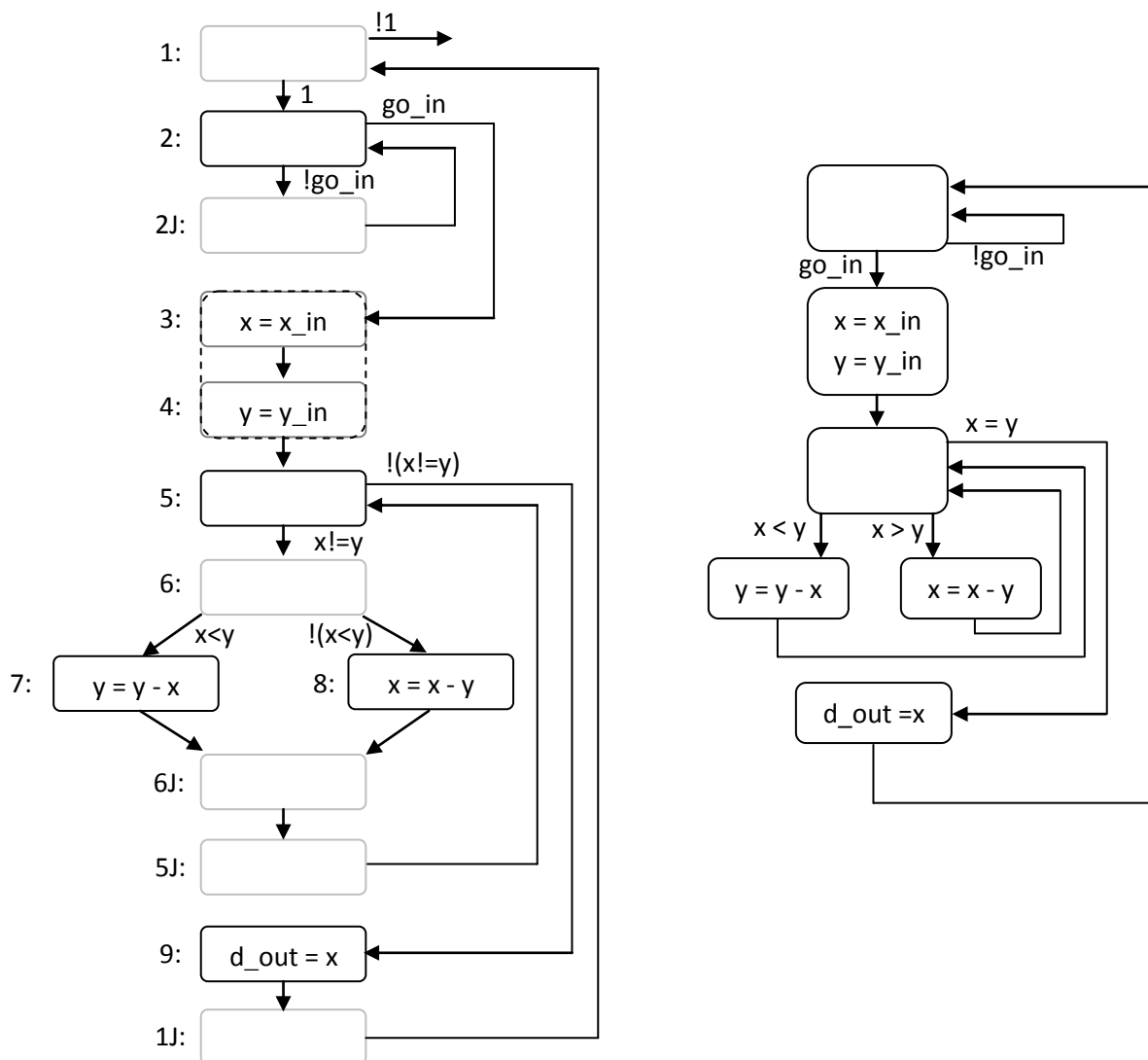


Figure 2.14: Optimized FSMD from original FSMD

The following actions are taken to optimize the original FSMD.

- Eliminate state 1 – transitions have constant values
- Merge state 2 and 2J – loop has no body
- Merge state 3 and 4 – operations are independent of each other
- Merge state 5 and 6 – transitions from state 6 can be done in state 5
- Eliminate state 5J and 6J – transitions from each state can be done from state 7 and 8 respectively
- Eliminate state 1J – transitions from state 1J can be done directly from state 9

Consider the operation $p = a * b * c * d$, if we use single state for this particular operation then three multipliers are required which renders system expensive and bulky. So the operation can be broken down as $x = a * b$, $y = c * d$ and $p = x * y$ with each operations having its own state. Thus, only one multiplier would be required in the system.

Optimizing Datapath

During the datapath design, the task of selecting a RT components for particular operation is termed as allocation. Whereas the task of mapping operations from the FSMD to allocated components is termed as binding. The optimization in datapath design can be done by following ways.

- **Sharing of Functional Units:** Single functional unit can be shared if same operations occur in different states. For example, in computation of GCD there were two subtractor used for two subtraction, rather a single subtractor can be used with the help of the multiplexor. Hence one to one mapping is not necessary.
- **Use of Multi-functional Units:** A variety of operations can be performed by ALU hence it can be shared for different operations occurring in different states.

Optimizing the FSM

Optimization in FSM can be done by:

- **State Encoding:** It is the task of assigning a unique bit pattern to each state in an FSM. The size of the register as well as the size of the combinational logic varies for different encodings. For example, if we have four states then it can be encoded as 00, 01, 10, 11 but it can also be encoded as 11, 10, 01, 00. If the number of state is large the number of ways of state encoding will be very large, hence CAD tools are used to determine the most efficient encodings.

- State Minimization:** It is the task of merging equivalent states into a single state. Two states are equivalent if those two states generate the same outputs and transition to the same next state, for any given input combinations. Merging equivalent states yield exactly the same output behaviour.

Few Solved Examples

Problem 1: Design a combinational logic circuit for the given problem whose description is given
as: y is 1 if a is 1, or b and c are 1. z is 1 if b or c is 1, but not both (or, a, b, and c are 1).

Solution: Initially, truth table is formed by writing down all possibilities of inputs followed by writing the outputs as defined by the given problem. Then, the K – map is used to minimize the equations and finally the combinational circuit is drawn.

A. Truth table

Inputs			Outputs	
a	b	c	y	z
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	1	0
1	0	1	1	1
1	1	0	1	1
1	1	1	1	1

B. K - map

y bc
 a

	00	01	11	10
0	0	0	1	0
1	1	1	1	1

 $y = a + bc$

$$y = a + bc$$

		bc			
		00	01	11	10
a	0	0	1	0	1
	1	0	1	1	1

$z = ab + b'c + bc'$

$$z = ab + b'c + bc'$$

C. Combinational Circuit

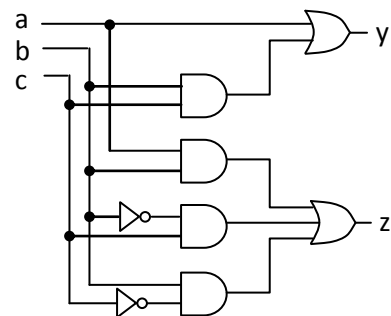


Figure 2.15: Truth table, K-map, and combinational circuit

Problem 2: Design a 2-bit comparator with a single output “less than”, using the combinational design technique described in the chapter. Start from a truth table, use K-maps to minimize logic and draw the final circuit.

Solution: As the comparator is 2 – bit, there must be total of four inputs; two inputs each of two bits. And only less than condition is to be checked, so only single output must be defined. Then the general steps for designing a combinational logic circuit is followed.

A. Truth table

a_1	a_0	b_1	b_0	It
0	0	0	0	0
0	0	0	1	1
0	0	1	0	1
0	0	1	1	1
0	1	0	0	0
0	1	0	1	0
0	1	1	0	1
0	1	1	1	1
1	0	0	0	0
1	0	0	1	0
1	0	1	0	0
1	0	1	1	1
1	1	0	0	0
1	1	0	1	0
1	1	1	0	0
1	1	1	1	0

B. K - map

		a_1a_0			
		00	01	11	10
b_1b_0	00	0	0	0	0
	01	1	0	0	0
	11	1	1	0	1
	10	1	1	0	0

$$It = b_1a_1' + b_0a_1'a_0' + b_1b_0a_0'$$

C. Combinational Circuit

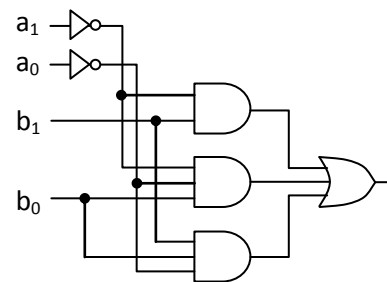
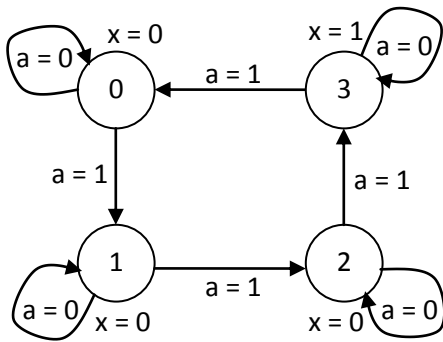


Figure 2.16: Truth table, K-map and combinational circuit for two bit comparator

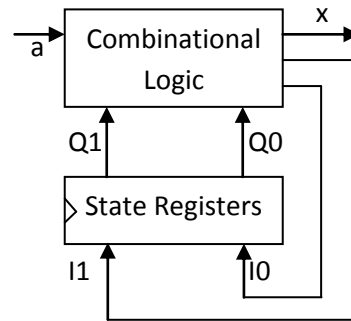
Problem 3: Construct a pulse divider. Slow down your pre-existing pulse so that you output a 1 every four pulses detected.

Solution:

A. State Diagram



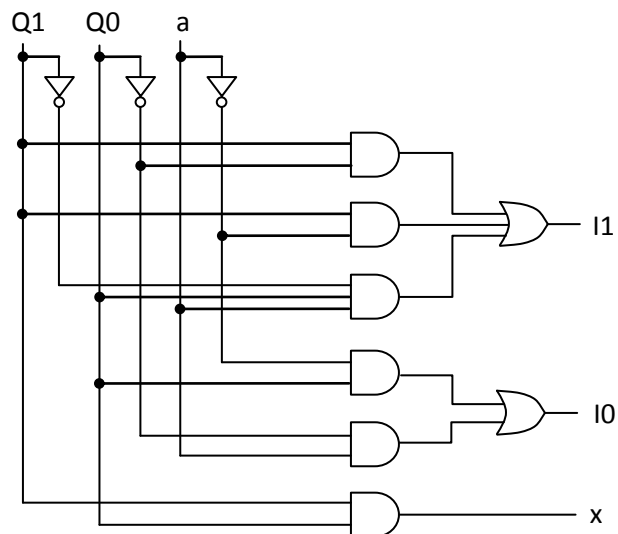
B. Implementation Model



C. State Table

Inputs			Outputs		
Q1	Q0	a	I1	I0	x
0	0	0	0	0	0
0	0	1	0	1	
0	1	0	0	1	0
0	1	1	1	0	
1	0	0	1	0	0
1	0	1	1	1	
1	1	0	1	1	1
1	1	1	0	1	

E. Combinational Circuit



D. K-map

I1

	Q1Q0			
	00	01	11	10
a 0	0	0	1	1
1	0	1	0	1

$$I1 = Q1'Q0a + Q1a' + Q1Q0'$$

I0

	Q1Q0			
	00	01	11	10
a 0	0	1	1	0
1	1	0	0	1

$$I0 = Q0a' + Q0'a$$

x

	Q1Q0			
	00	01	11	10
a 0	0	0	1	0
1	0	0	1	0

$$x = Q1Q0$$

Figure 2.17: Pulse Divider – State diagram, state table, K-map, combinational circuit

Problem 4: Design a single purpose processor that calculates x to the power n (x^n). Include FSMD, Datapath and FSM in the design.

Solution:

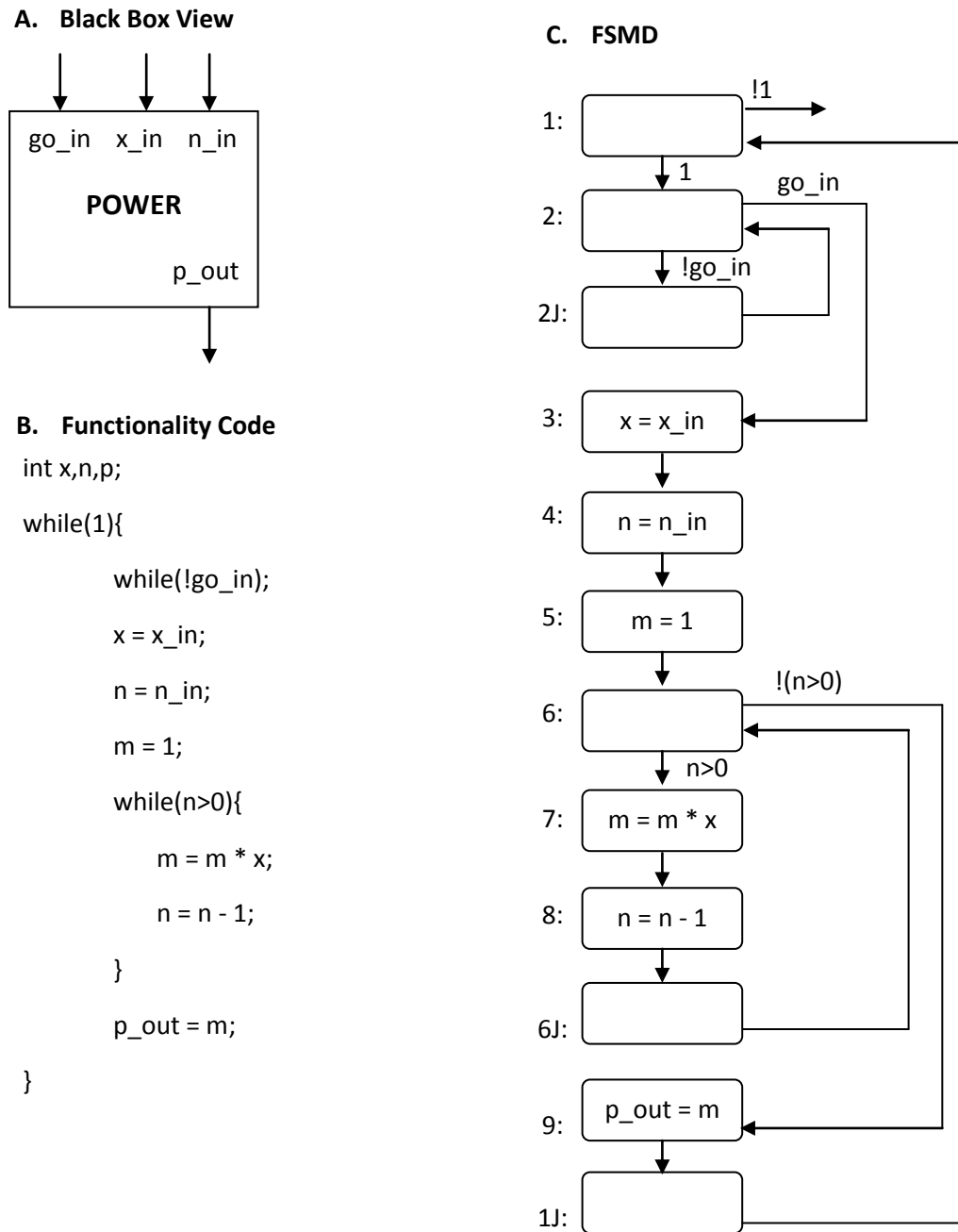


Figure 2.18: The black box view, functionality and FSMD for processor that calculates x^n .

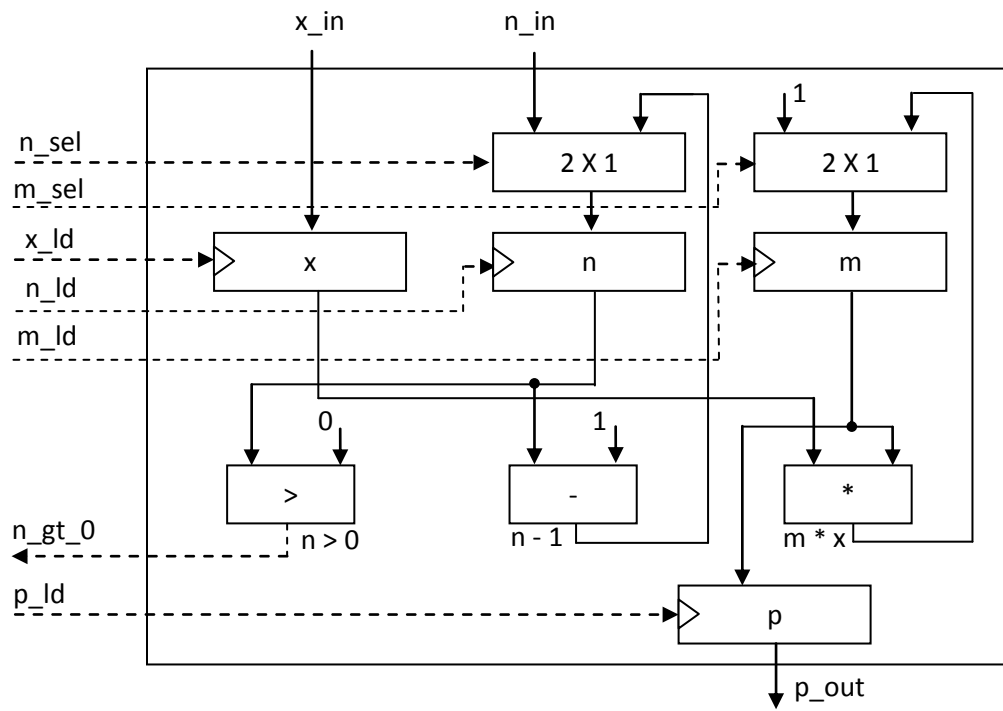
D. Datapath for the processor that calculates x^n 

Figure 2.19: Datapath

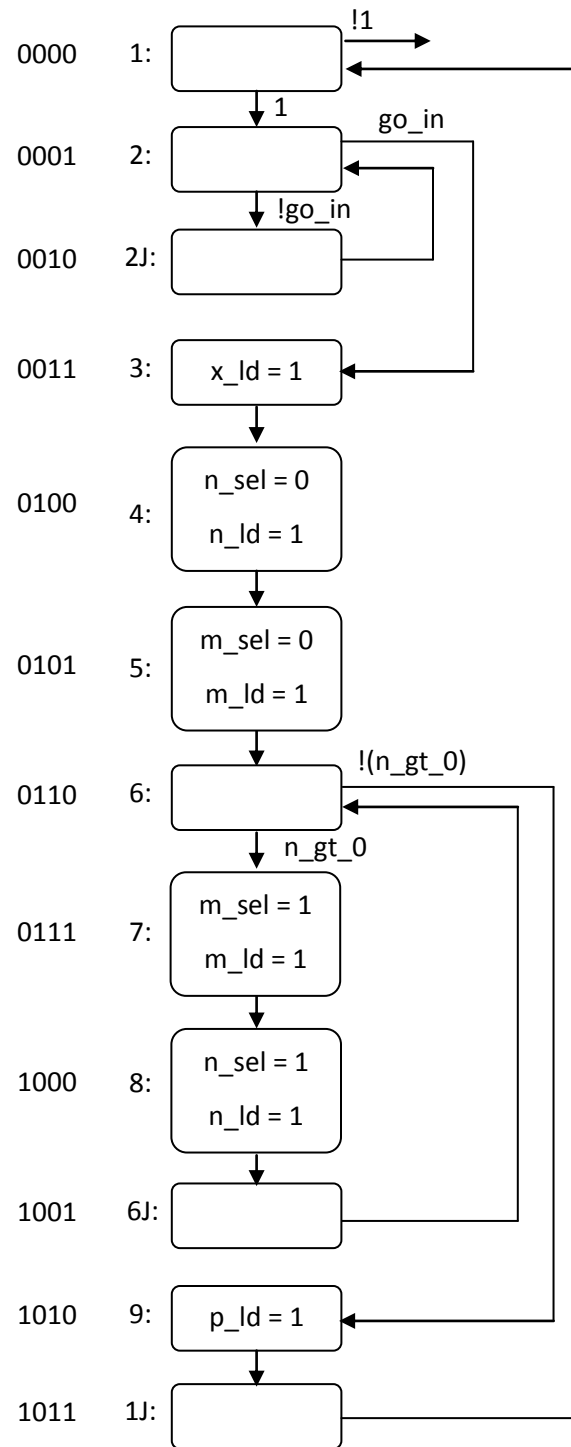
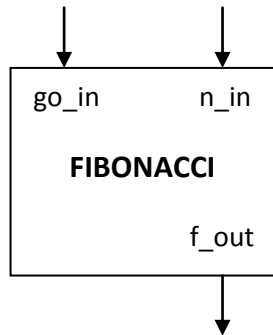
E. FSM of the processor that calculates x to the power n 

Figure 2.20: FSM Controller

Problem 5: Design a single purpose processor that generates Fibonacci series up to n places. Start with a function that computes desired result, translate the function into a state diagram, sketch a probable datapath, and draw FSM diagram.

Solution:

A. Black Box View



B. Functionality Code

```

int ft, st, nt, count, n;
while(1){
    while(!go_in);
    n = n_in;
    ft = 0;
    st = 1;
    count = 1;
    while(count <= n){
        f_out = ft;
        nt = ft + st;
        ft = st;
        st = nt;
        count++;
    }
}
  
```

C. FSMD

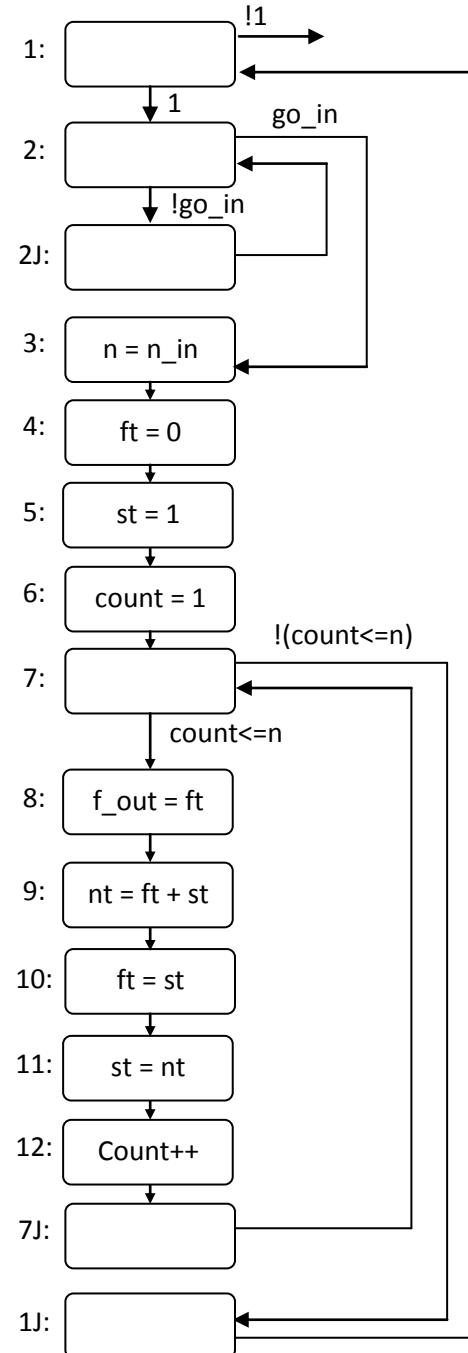


Figure 2.21: Fibonacci series generator – the black box view, functionality and FSMD

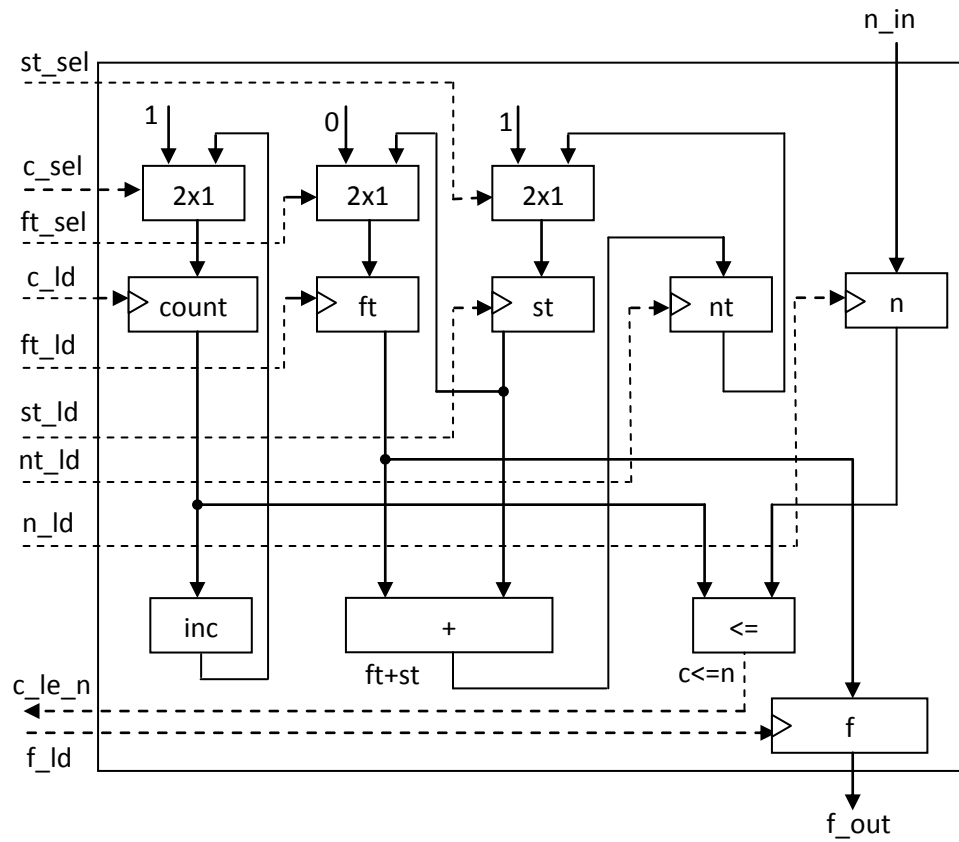
D. Datapath of the processor that generates Fibonacci series:

Figure 2.22: Datapath for Fibonacci series generator

E. FSM controller for Fibonacci series generator

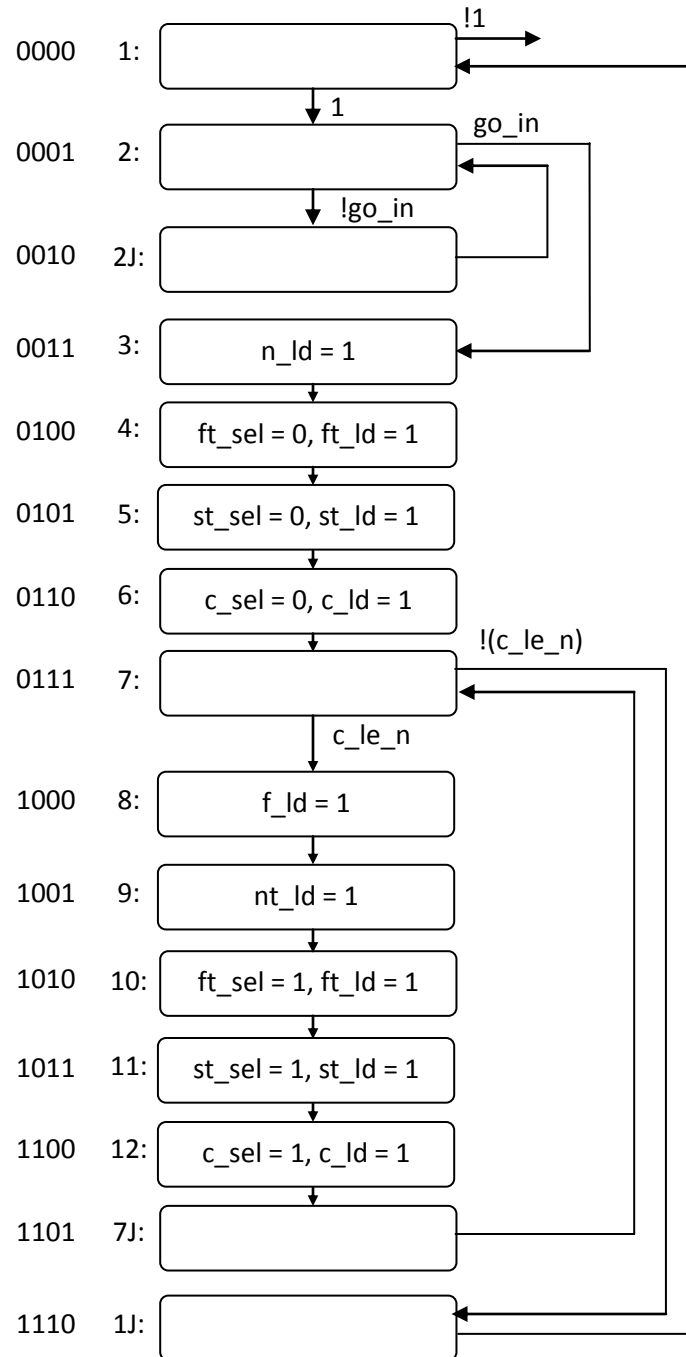


Figure 2.23: FSM of Fibonacci series generator

- **Basic Architecture**
- **Operation**
- **Programmer's View**
- **Development Environment**
- **Application-Specific Instruction Set Processors**
- **Selecting a Microprocessor**
- **General-Purpose Processor Design**

3.1 Basic Architecture

A general-purpose processor is a programmable digital system which consists of a datapath and a controller which are tightly linked with a memory. Figure 3.1 shows the various components in the architecture of the general-purpose processor.

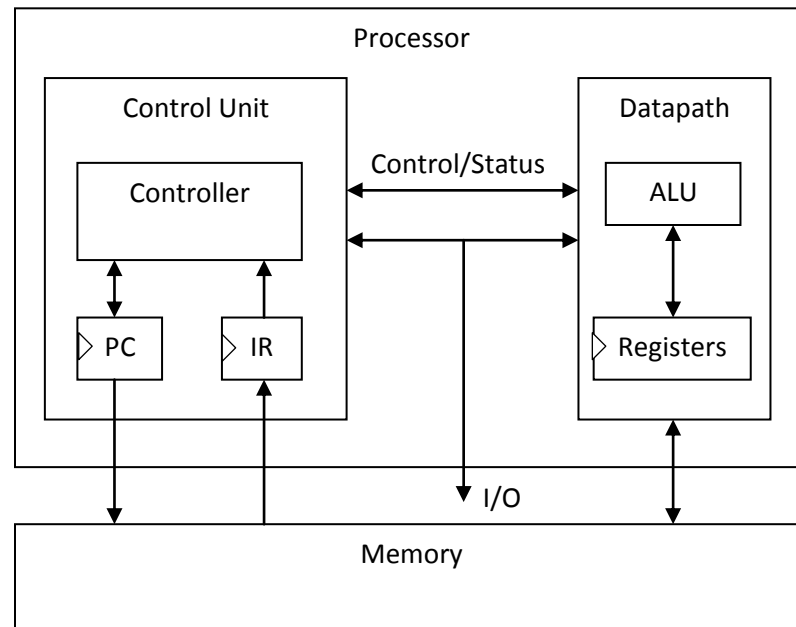


Figure 3.1: Basic Architecture of general-purpose processor

DATAPATH

Datapath consists of the circuitry for transforming data and for temporary data storage. It contains an arithmetic-logic unit which manipulates data through various operations such as addition, subtraction, logical AND, logical OR, rotating, shifting etc. ALU also generates status signals to represent various conditions such as carry, zero, sign, parity and so on. Such information is stored in status register.

Data path contains registers to store temporary data and different status generated by operations. The temporary data may be the data from memory for ALU to process, or the data that needs to be moved from one memory to another memory, or the data from ALU that needs further processing by ALU or needed storage. For data transfer within datapath, internal bus is used. But movement of data from and to memory is done by external bus.

CONTROL UNIT

The control unit consists of circuitry to generate control signals to carry out various operations. It consists of controller, Program Counter (PC) and Instruction Register (IR).

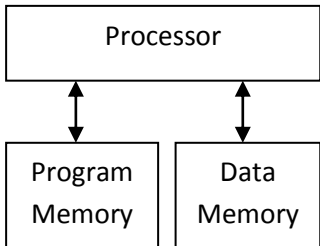
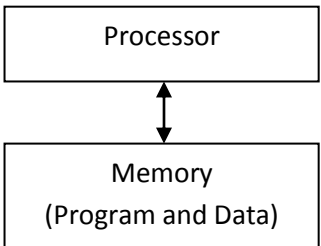
Controller consists of a state register and control logic. It sequences through the states and generates the control signals to read instructions into the Instruction Register, and control the flow of data between ALU, registers of datapath and memory. Controller also determines the next value of Program Counter. For non-branch instruction, the value of Program Counter is incremented. But for branch instruction, status signals from datapath and content of Instruction Register are evaluated for next address of program counter.

Program Counter is used to hold the address of the next program instruction to be fetched, while an Instruction Register is used to hold the fetched instruction. The bit width of Program Counter indicates the address size of memory which in turn can be used to determine the number of directly accessible memory locations. For example, A 16 bit PC represents address size of 16 bit and $2^{16} = 65536$ addressable memory locations.

MEMORY

Memory is used to store information for medium or long term. Information can be data or program. Program information is the set of instructions that is used to carry out desired function. Data are the information used by the program for various purposes.

There are two memory architectures based on program and data storage.

SN	Harvard Architecture	Princeton Architecture
1.	Distinct data and program memory space	Data and program share memory space
2.	Improved performance: Data and instructions can be fetched simultaneously	Data and instructions cannot be fetched simultaneously
3.	More connecting wires	Less connecting wires
4.	<p>Block Diagram</p> 	<p>Block Diagram</p> 

3.2 Operation

Instruction Execution

Instructions are the sets of code that carry out particular function. For each instruction, the controller sequences through several stages. Each stage may consist of one or more clock cycles.

The various stages or sub-operations can be explained as:

- **Fetch Instruction:** The next instruction to be executed is loaded into Instruction Register from memory. The address of the memory where instruction resides is given by program counter.
- **Decode Instruction:** Instruction in the instruction register may represent various operations based on op-code and may include register or memory as operands. In this stage, the operation to be done by the instruction is determined.
- **Fetch Operand:** For a given operation, operand can be a register or memory. In operations including registers, the required data are loaded into registers as specified by the instruction.
- **Execute Operation:** The ALU handles the arithmetic and logical operations defined by the instructions. The loaded registers are fed to the inputs of ALU to carry out the operation.
- **Store results:** The destination to store results may be either register or memory. After the execution of operation, the final data is loaded into register or memory as defined by the instruction.

Pipelining

Pipelining is implemented to increase the throughput of the system. In pipeline, the given task is divided into various stages and multiple stages which are independent of each other are executed simultaneously. For efficient instruction pipeline, different stages must be of almost same length and each instruction must require same number of cycles to complete its execution.

Branching instructions can be an obstacle for efficient pipeline as next instruction to be executed will only be known after execution stage of branch instruction. This problem, however, can be addressed using various techniques. One simple method is to stall the pipeline when there is an occurrence of branching instruction. The pre-fetch of next instructions is not done in this method rather waited for execute stage to complete first. Another popular method is to use

branch prediction. In this method, the branch is guessed and the next instruction is fetched correspondingly. If the guess is correct, then it results in efficient pipeline. But, however, if the prediction is not correct then all pre-fetched instructions in the pipeline must be ignored. The following diagram shows an example of an instruction pipeline having five stages.

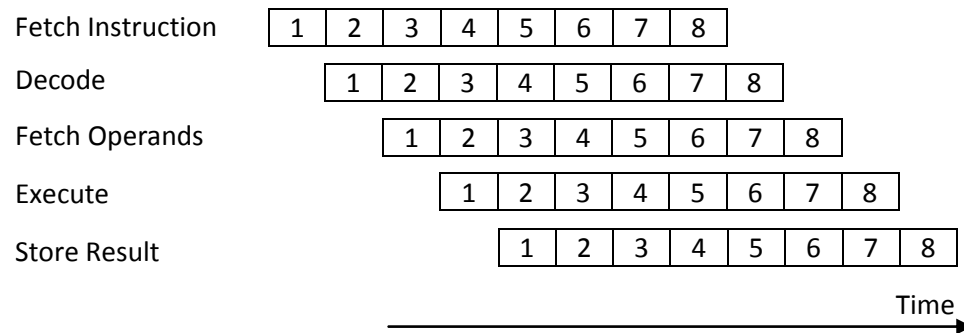


Figure 3.2: Eight instructions in execution using instruction pipeline

In figure 3.2, there are 8 instructions in the pipeline with each instruction divided into five stages and each requiring equal time to complete. In absence of pipeline, the total time required to complete eight instructions would be $8 \times 5 = 40$ clock cycles, assuming each stage to complete in one cycle. However, with pipeline implementation, the total completion time required is 12 clock cycles. In this way, pipeline helps to improve the performance of the system.

Superscalar and Very Long Instruction Word (VLIW) Architectures

Multiple ALU architecture is implemented in superscalar architectures to improve the performance of the system. Such systems can execute two or more scalar operations in parallel, which increase the requirement of ALU in the processor. It may require extensive hardware to detect multiple independent instructions that can be executed simultaneously. Instructions in such architectural systems are ordered statically (at compile time) or dynamically (during runtime).

Very Long Instruction Word (VLIW) architecture is a type of static superscalar architecture. It contains multiple independent instructions in a single word. Several operations are encoded in a single machine instruction. The compiler detects and schedules the instructions.

3.3 Programmer's View

Software programmer does not require detail understanding of architecture of the system; instead they need to know what instructions are available and how they are used. Embedded system programmer, however, needs to know certain information, if not all, about the system, as the programs may include assembly level language.

Instructions in the program may be of different level. Firstly, in machine level the codes are in binary form. Another level is assembly level in which mnemonics are used to represent instructions which are processor specific in nature. Next level of programs follow structured languages which are processor independent. Programmer of assembly language and machine language must have information about architecture of the processor.

In embedded system design, the programmer must be aware of the following:

A. Instruction Set

The instruction set is a list of instructions which represent the bit configurations for operations that can be carried out by the processor. Assembly language programmer must be aware of the available instruction set. Since embedded system design may require some portion of assembly code to be written, programmer of embedded system must know the instruction set available for the processor they are working on.

Every instruction, in general, consists of op-code and operand field. Op-code field specifies the operation to be done. Different types of instructions are explained briefly.

- Data-transfer instructions move data from memory to register, register to memory or register or input/output ports.
- Arithmetic and logical instructions cause ALU to carry out certain operations involving registers and store the final result back to register. ADD, SUB, AND, OR etc are few examples of arithmetic and logical operations.
- Branch instructions change the flow of program and it also determines the address of the next instruction to be executed. Branch instructions may be unconditional jumps, conditional jumps or procedure call and return instructions.

An operand field specifies the location of actual data that takes part in an operation. The number of operands per instructions varies among processors and its instruction type. Addressing modes are used to represent data location and its accessing mechanism. The simple instruction format is shown in the figure below.

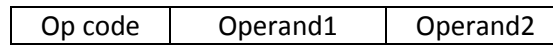


Figure 3.3: Simple two address instruction format

Commonly used addressing modes are explained in the following paragraph.

Immediate Addressing Mode: The operand field contains the actual data.

Register Direct Addressing Mode: The operand field contains the address of the datapath register in which the data is stored.

Register Indirect Addressing Mode: The operand field contains the address of a register, which in turn contains the effective address of the data in memory.

Direct Addressing Mode: The operand field contains the effective address of actual data that is used in operation

Indirect Addressing Mode: The operand field contains the address of a memory location, which in turn contains the address of a memory location in which actual data is available.

Implicit or Implied Addressing Mode: The operand field is not used in this mode; the register to be used in operation is defined implicitly. In general, accumulator is used as an implicit register.

Displacement Addressing Mode: The operand is added to a particular register to obtain the effective address of the data. In index addressing, index registers are used. While in relative addressing, value of operand is added to the current address to determine the actual address.

The operations of few addressing modes can be visualized using following figure.

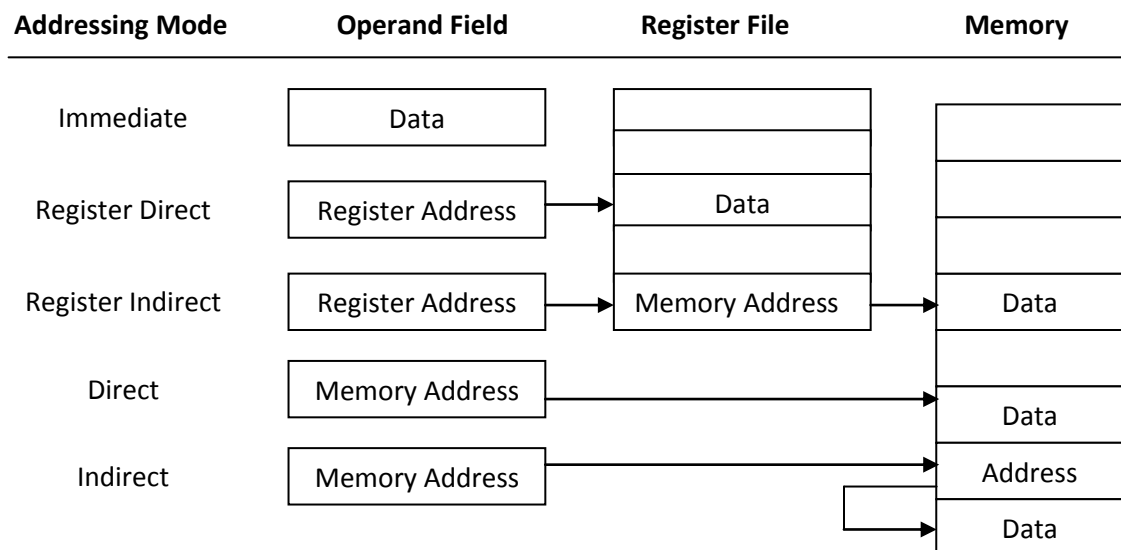


Figure 3.4: Addressing Modes

B. Program and Data Memory Space

The programmer in embedded system design must be aware of the size available for program and for data. Programs must be written within the defined memory space limits. For example, in microcontrollers, the on-chip memory for program and data are fixed. So, one should be able to write the code efficiently so as not to exceed the memory limit.

C. Available Registers

Programmer of embedded system design must be informed about the number of registers available for general purpose and specific purpose. For example, multiplication in 8051 microcontroller can be done using accumulator and register B. Information about accumulator and Register B is not required for structured language programmer. However, various special function registers used for configuring timers, serial communication, and interrupts must be known to every programmer.

D. Input Output Facility

Every processor facilitates programmer with input output pins to communicate with external devices. Programmer working with processors must be alert about the number of input output pins available and their functions. In parallel I/O, port can be read or written to using specific function register. Also, communications can be done through system bus in which address and data ports can be activated by certain instructions.

E. Interrupts

Interrupt is a facility provided to the user in which the processor serves the device which requires urgent attention. It causes processor to suspend execution of the current program and starts executing interrupt service routine that does the function required by the device which interrupts the processor. The programmer should be aware of the types of interrupts supported by the processor and must write interrupt service routine when required.

F. Operating System

An operating system is a layer of software that provides low-level services to the application layer. Few services involve loading and executing of programs, sharing and allocating system resources, and synchronization mechanism. Another important service is process scheduling in which the high priority process is executed first. Other services include handling hardware interrupts, and provide device drivers.

High level applications invoke operating system using system call. When a program requires service from operating system, it generates a predefined software interrupt that is served by the operating system. Values required to the services are typically passed as the parameters in the program. CPU registers are involved for information exchange among application programs and operating system.

3.4 Development Environment

Processors along with different development tools are used for the development of software or an embedded system. Processor that is used to write and debug the program is commonly referred as **development processor**. Desktop Computer can be taken as an example of development processor. Such processors may not be a part of embedded system's implementations. But the processor in which our program is loaded is referred as **target processor**. AVR, 8051, PIC microcontrollers or 8085, 8086 microprocessor can be few examples of target processor. Such processors are always a part of system implementations. Various tools for the software development as well as embedded systems development are described in the following paragraphs.

Tools for Implementation Phase

Assemblers convert assembly instructions to binary machine instructions. It replaces op-code and operand mnemonics by binary equivalent. It also translates symbolic labels into actual addresses. It generates a equivalent binary code for a single machine instruction, so it follows one to one mapping principle.

Compilers convert high level programs to machine programs. Each high-level constructs may be translated to several machine instructions. Hence, it may not follow one to one mapping principle. Cross compilers are those compilers which run on one processor but generate the code for a different processor.

Linker combines object files into a single executable file, or another object file. It allows creation of a program in separately assembled or compiled files. It combines machine instructions of user code and instructions from standard library.

Tools for Verification Phase

Debuggers are programs that are used to test and debug the targeted program. These are programs that run on development processor but execute code designed for target processors. It simulates the function of the target processors and allows evaluation and correction of

programs in development processor. Stepping, breakpoints, watch values are few debugging techniques supported by various debuggers. These debuggers are also known as instruction set simulators (ISS) or virtual machines (VM). Design cycle for debuggers is fast as compared to other tools, since the program is tested in development processor. But, these tools can, however, lead to inaccuracy as it does not interact with the actual system.

Emulator can be a hardware or software that enables one system to behave like another system. It consists of debugger coupled with a board connected to development processor. The board consists of target processor or device similar to target processor and support circuitry. It supports debugging of program while it executes on target processor. It also enables one to control and monitor the program's execution in actual embedded system circuit. Since the code must be downloaded into emulator hardware in each test, the design cycle is little longer compared to debugger. But it leads to accurate testing as it interacts with the rest of the system components as well.

Device Programmers are the devices with the help of which binary machine programs are loaded into target processor's memory. Using this tool, the program can be tested in its realistic form which results in high accuracy as program runs on actual system. The design cycle, however, is longest since the target processor is removed from the system, programmed using programmer and returned to the system. If the device programmer can be made in-build within the system, the design cycle will be reduced.

DESIGN FLOW

Every software or system development process includes implementation and verification phase. During implementation phase, various implementation tools such as assembler, compiler are used while verification tools such as debugger, programmer are used in verification phase.

Software Development Process

For a software development, the development processor as well as the target processor may be common. And the development tools are available in a single package which is referred as Integrated Development Environment (IDE).

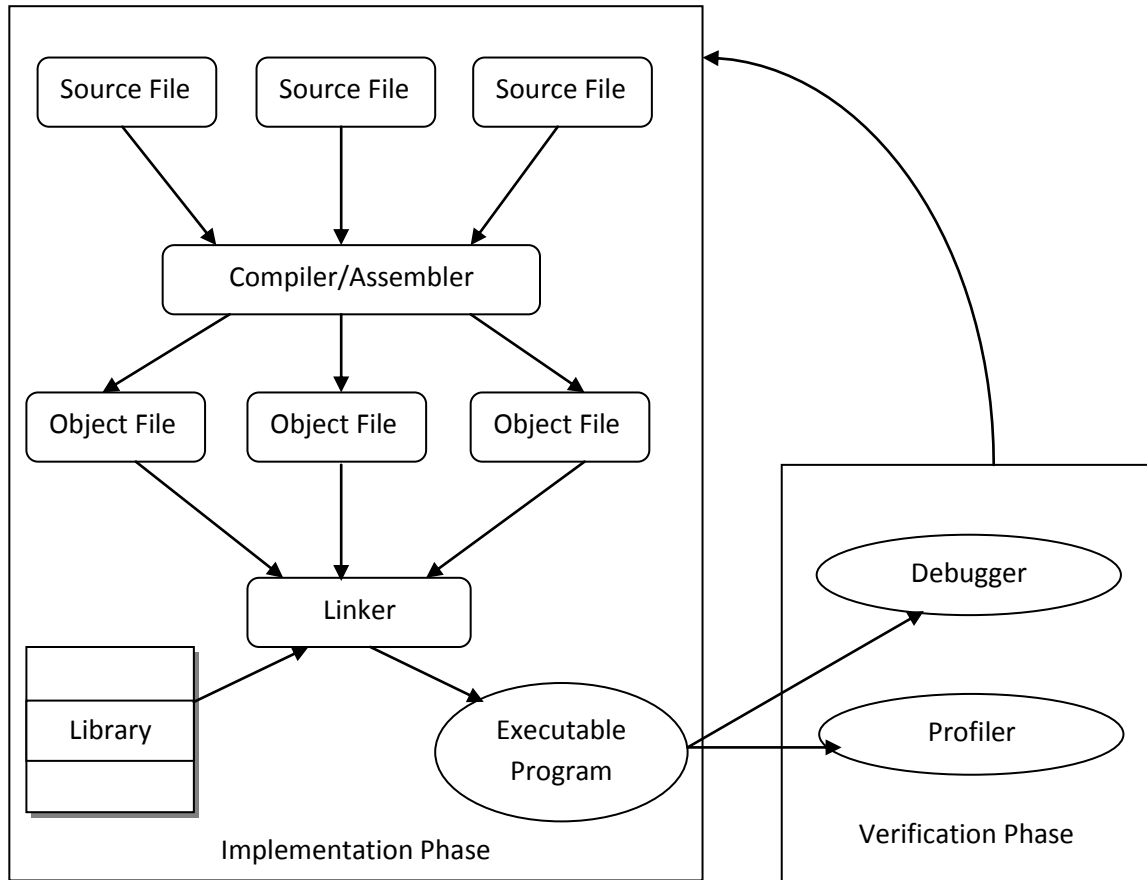


Figure 3.5: Software Development Process

Implementation Phase: Source code is written using an editor, and then the code is compiled/assembled using compiler/assembler. Finally, with the help of linker all required files are combined into a final executable file.

Verification Phase: The executable file is run under the command of a debugger. All possible inputs, especially boundary cases, are used to check the behavior of program. Profilers can be used for performance analysis of the program. Time and space complexity can be analyzed. Time complexity includes duration of execution of program whereas space complexity includes memory usage.

Embedded System Development Process

In case of embedded system design, the target and development processors are different in almost all systems. The Integrated Development Environment (IDE) tools for various processors are available for implementation phase. Though the implementation phase for embedded

system is similar to that of software implementation phase, the verification phase differs drastically.

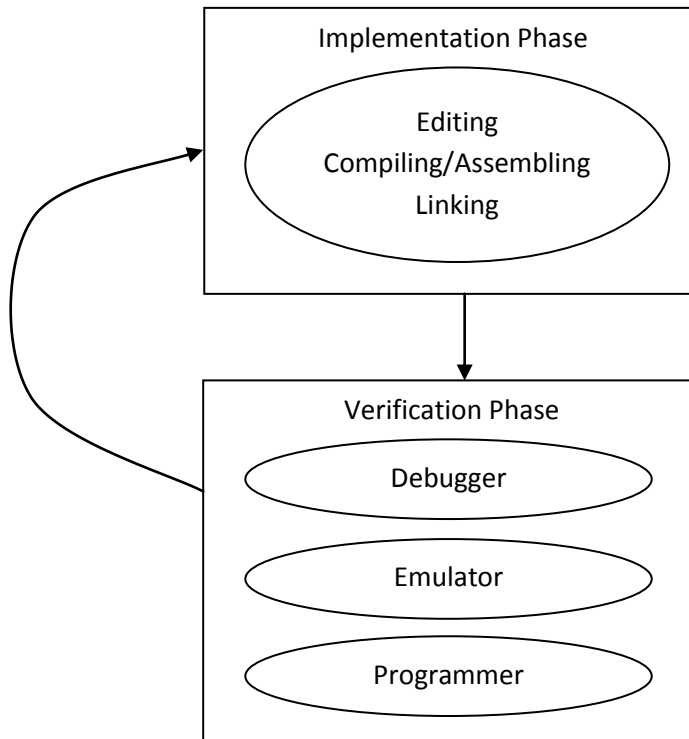


Figure 3.6: Embedded System Development Process

Implementation Phase: The process of editing, compiling/assembling and linking the program is same as that of software development process. However, development processors use cross compilers or cross assembler. As those compilers run on development processor, for example PC, and generate the file for target processor, for example hex file for microcontrollers.

Verification Phase: Embedded system works in conjunction with other components as well as with real time environment, so debugging a program requires control over time and environment. Based on requirement and availability, debuggers, emulators or device programmers can be used for verification. Code may be simulated on development processor using debuggers or code may be checked by loading into emulator hardware. Also, programmer can be used to load the code directly into the target processor.

3.5 Application-Specific Instruction Set Processors

Application-Specific Instruction Set processors are specific to the particular application domain. They can be programmed based on requirement of particular arena, which makes it more flexible. Also other constraints such as performance, power, cost, and size are efficient enough to develop a system. However, instruction set processor and its associated software tools are expensive to develop. It can be categorized as microcontrollers, digital signal processors and less general application specific instruction set processors.

Microcontrollers

Microcontrollers are specific to applications that perform a large amount of control oriented tasks. The following are few general features of microcontrollers.

- It includes several peripheral devices such as timers, analog to digital converters, serial communication devices, and so on.
- It generally contains program and data memory on the same IC. Various peripherals along with memory incorporated within the same IC result in compact and low-power implementation.
- It provides the programmer direct access to number of pins of the IC. Access to pins enable programmer to interface with other devices such as sensor, actuators, LCDs, and other different devices that may be used in the system.
- Some specialized instructions may be available. Such facility improves the performance of the system.

Digital Signal Processors (DSP)

These are processors which are specific to applications that process large amounts of data. The source of large amount of data includes image captured by a camera, voice packet through a network router, audio clip played by an instrument. Few features, out of many, are listed below.

- It may contain numerous register files, memory blocks, multipliers and other arithmetic units.
- It facilitates with instructions that are applicable uniquely to digital signal processing. Filtering and transforming vectors can be two examples.

- Frequently used arithmetic functions are implemented using hardware. It results in faster execution of arithmetic functions compared to software implementation.
- Some special digital signal processors allow concurrent execution of functions which boost the performance of the system.
- It incorporates many peripherals specific to signal processing. It may include ADC, DAC, PWN, DMA controllers, Timers and Counters.

Less-General ASIP

These are developed to perform some very domain specific processing while allowing some degree of programmability. Processors designed for networking hardware can be taken as an example of less-general ASIP.

3.6 Selecting a Microprocessor

In any embedded system, a designer must select the microprocessor based on technical and nontechnical aspects.

- Technical aspects: Selection of processor must be done based on required speed within limited power, size, and cost.
- Non technical aspects: Before selecting microprocessor, one must be aware of development environment, prior expertise of processor, licensing arrangements and so on.

Comparing Speed

Speed of processors can be measured and compared using various methods.

A. Clock Speed of Processor

Speed can be compared based on clock speeds of processors, but the number of instructions per clock cycle may differ. So, it may not be a efficient method unless processors to be compared have same number of instructions per cycle.

B. Instruction per second

The speed can be evaluated using number of instructions executed per second. But the complexity of available instruction sets may differ creating some hindrance in speed

comparison. For example, to perform same operation, one processor may require 200 instructions while another may require 300 instructions.

C. Dhrystone benchmark

It is a program that runs on different processors and evaluates their performance based on execution of certain operations. Dhrystone benchmark performs no useful work rather checks the integer arithmetic and string-handling capabilities of the processor on which the benchmark runs on. Since processors can execute such operations thousands of times in a second, speed of processor may be expressed in terms of Dhrystones per seconds.

D. Millions of Instruction Per Second (MIPS)

It is a general measure of computing performance and the amount of work a processor can do. MIPS can be useful when comparing performance of processors having similar architecture. The origin of MIPS is based on VAX 11/780 which could execute one million instructions per second or could execute 1757 Dhrystones per second. Hence, 1 MIPS = 1757 Dhrystones/sec. Also, performance of other computers were measured based on VAX 11/780.

3.7 General-Purpose Processor Design

General-purpose processor can be designed using the design technique of single-purpose processor as general-purpose processor is a type of single-purpose processor which process instructions stored in program memory. The design starts with the design of instruction set, followed by creating a FSMD and datapath, and finally the controller is developed. All steps are explained in details with example.

A. Instruction Set Design

Instruction set defines various operations that can be done by the processor. It also determines the size of memory required along with number of registers to be used.

- Initially, how many and what kind of operations are to be included must be considered.
- Then, how many, what types, and location of operands to be used must be selected.
- And finally the size and format of instruction must be set.

B. Creating a FSMD

FSMD represents the state diagram of the given functionality which is based on instruction set. The following steps are required to generate a FSMD.

- First, RESET state is defined, which can be used to clear various registers.

- Next, FETCH state is used where instruction from memory is loaded in instruction register. In this state, program counter is increment after each instruction fetch.
- Then, DECODE state is used as a transition state before execution of instructions. In this state, no operation is done but it adds extra cycle necessary for instruction register to get updated.
- Finally, EXECUTE state is defined based on the operation represented by opcode. The number of EXECUTE state is given by the number of instructions available. The operation to be performed is detected just before the start of this state and hence the actual instruction operations are carried out in this state.

C. Building a Datapath

To carry out various operations of FSMD, datapath must be build. The following steps are required.

- For each declared variable, we need to instantiate a storage device
- Then, instantiate functional units to carry out the FSMD operations. General purpose ALU is, usually, implemented in the design.
- Connecting different components within the datapath. When more than one input appear at any ports then appropriate multiplexor is needed.
- Finally, unique identifiers are created for every control signal.

D. Development of Controller or FSM

- Rewrite the FSMD states without any instructions or operations
- Equivalent binary operations on control signals must be written in each state rather than the operation. Each FSMD operation must be replaced by binary operations

EXAMPLE: Design a general purpose processor with four data transfer instruction, two arithmetic operations and one jump instruction.

→ The following are the considerations made in the design.

- 16 bit instruction size, which has direct impact on memory and register selections.
- Instruction Register (IR) and Program Counter (PC) of 16 bit,
- Memory of 64K x 16 bit,
- Register file of 16 x 16 bit

A. Instruction set design

Instruction	First Byte		Second Byte		Operation
MOV Rn, direct	0000	Rn	Direct		$Rn = M(\text{direct})$
MOV direct, Rn	0001	Rn	Direct		$M(\text{direct}) = Rn$
MOV @Rn, Rm	0010	Rn	Rm		$M(Rn) = Rm$
MOV Rn, #imm	0011	Rn	Immediate		$Rn = \text{immediate}$
ADD Rn, Rm	0100	Rn	Rm		$Rn = Rn + Rm$
SUB Rn, Rm	0101	Rn	Rm		$Rn = Rn - Rm$
JZ Rn, relative	0110	Rn	Relative		$PC = PC + \text{relative (if Rn is 0)}$

Figure 3.7: A simple instruction set

From the above instruction set, the various means of data transfer and operations can be analyzed which may be useful in developing FSM and datapath.

- The address of memory location are available from
 - Instruction Register: In instruction MOV Rn, direct and MOV direct, Rn, the direct address is used which is available in IR as lower bytes.
 - Register: In instruction MOV @Rn, Rm, the address of memory is given by value of register.
(Address is also given by PC to load the instruction into IR)
- The value in register can be loaded from:
 - Memory: In instruction MOV Rn, direct, register is loaded from memory whose address is given by lower eight bits of IR.
 - Instruction Register: In instruction MOV Rn, #imm, the immediate value of IR is loaded into register.
 - ALU: After execution of ADD Rn, Rm and SUB Rn, Rm, the final result is stored in register.
- Three operations are performed by ALU
 - Addition, subtraction and comparison

B. FSMD for given instruction set

In FSMD, the basic stages of instruction cycle are implemented as states. It includes RESET, FETCH, DECODE and EXECUTE state. The RESET, FETCH and DECODE states are common to almost every design. The EXECUTE state, however, differs when the number and type of instructions are different.

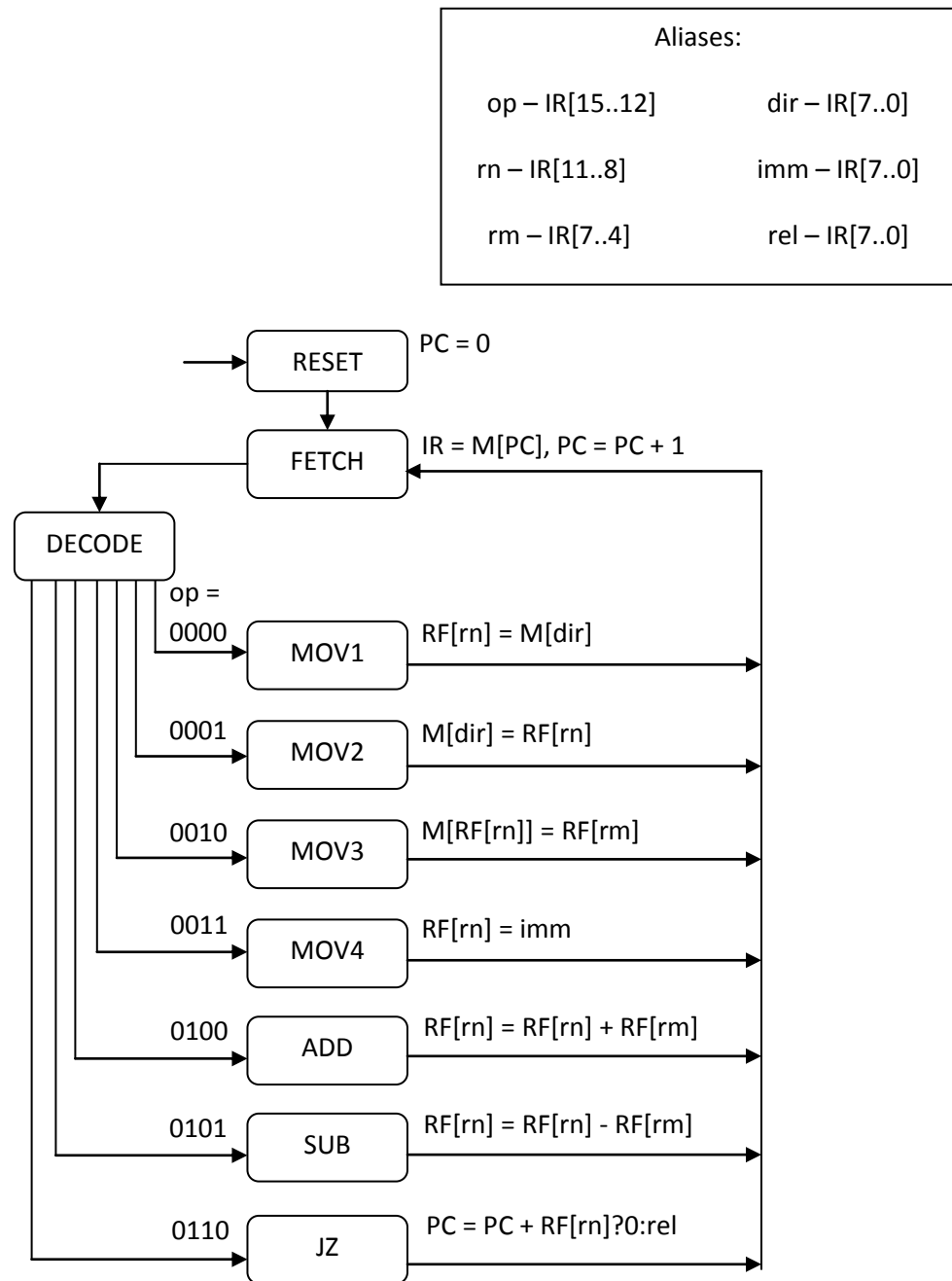


Figure 3.8: Finite State Machine with Data (FSMD)

C. Datapath for FSMD

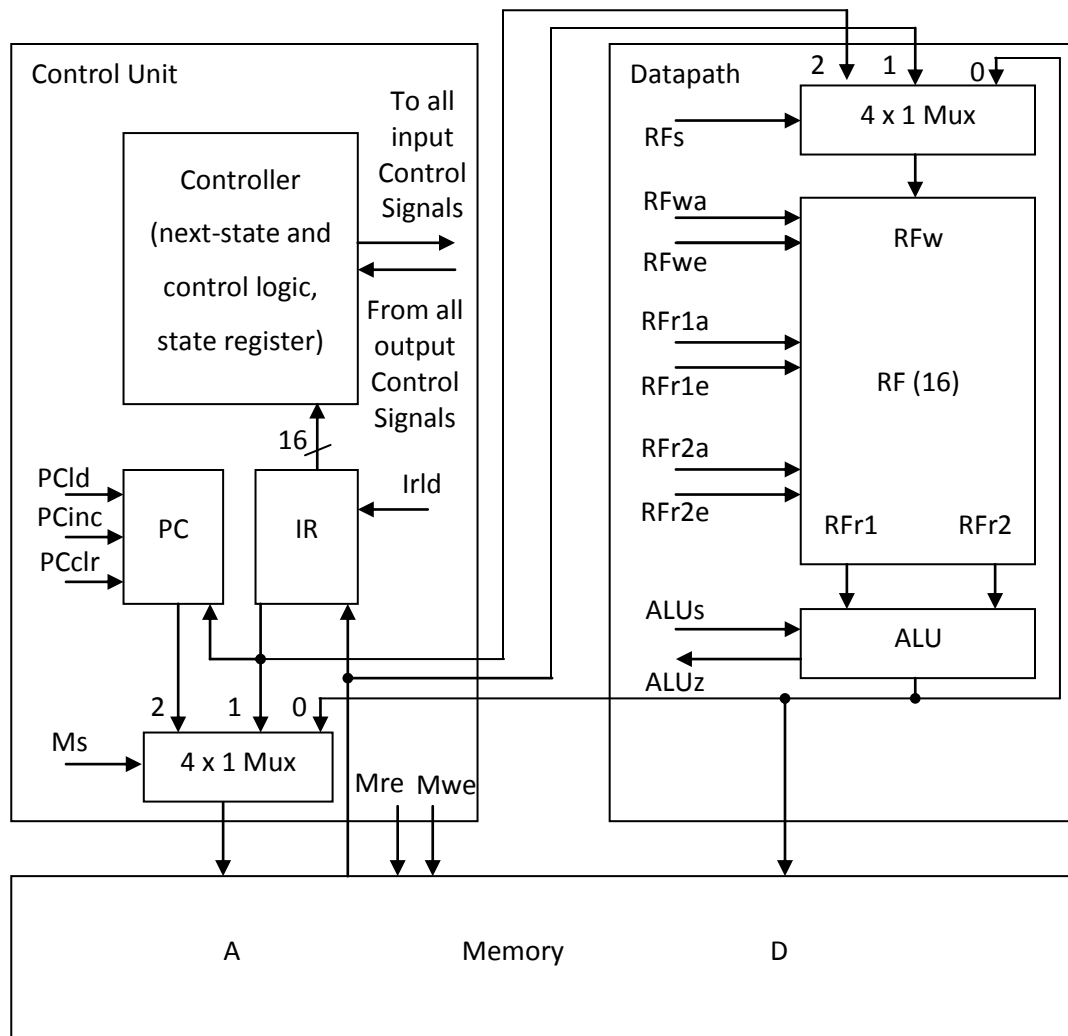


Figure 3.9: Datapath of our simple general purpose processor

Components in datapath

- Register file of 16x16 and a general purpose ALU.
- Multiplexer of 4x1, since the register in register file can have three sources; Immediate data from IR, data from Memory, and data from ALU

Components in Control Unit

- Controller for next-state and control logic, state register, Program Counter, Instruction Register
- Multiplexer of 4x1, since memory address can be selected from three sources; address from PC, direct address from IR, and address from register.

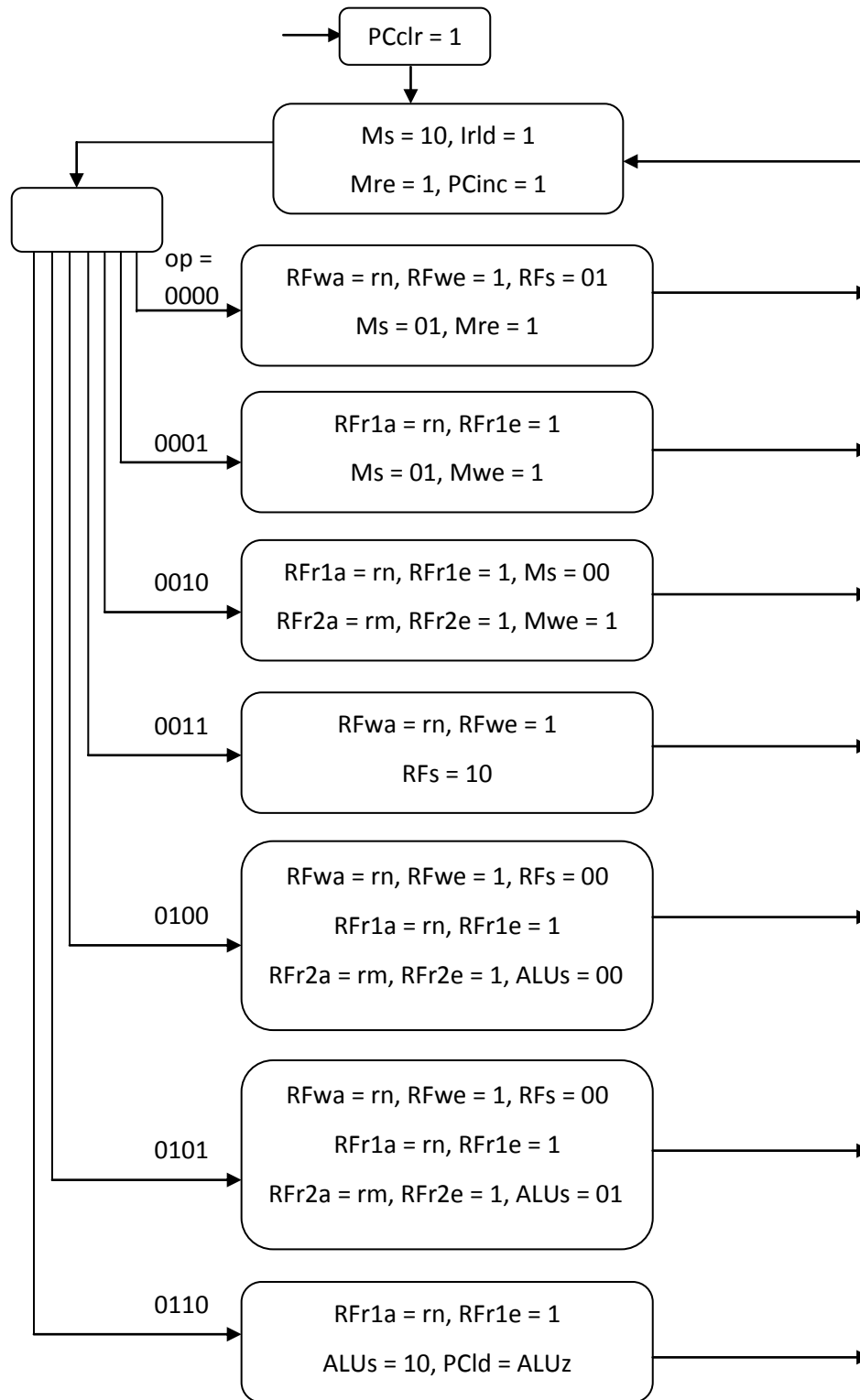
D. Finite State Machine (FSM) Design

Figure 3.10: Finite State Machine (FSM)

Converting FSMD operations to FSM operations

Example 1: MOV Rn, direct \rightarrow $RF[rn] = M[dir]$

It means to read the content of memory of address dir (8 lower bits of IR) and write it into one of registers of register file. Value of rn gives the address of register in register file.

- Address of memory is directly available in IR, using multiplexer selection $M_s = 01$ will select address from IR. For a memory read operation, M_{re} must be set ($M_{re} = 1$).
- The value is to be written into register file, so $RF_{wa} = rn$ selects a register from register file and RF_{we} enables the write operation. Set $RF_s = 01$, as data is coming from memory.

Example 2: ADD Rn, Rm \rightarrow $RF[rn] = RF[rn] + RF[rm]$

Here, values from two registers are read and then added using ALU. The final result is stored in register. Address of registers to be selected is given by rn and rm for read operation while value of rn gives the address of register for write operation.

- Selection of registers for read operation: $RF_{r1a} = rn$ and $RF_{r2a} = rm$ select two registers while $RF_{r1e} = 1$ and $RF_{r2e} = 1$ enable both registers for read operation.
- Adding the value of registers using ALU: $ALU_s = 00$ represent the addition of two registers.
- Selection of register for write operation: $RF_{wa} = rn$ selects the register and $RF_{we} = 1$ enables the write operation.

- **Memory Write Ability and Storage Permanence**
- **Common Memory Types**
- **Composing Memory**
- **Memory Hierarchy and Cache**

4.1 Introduction

A memory stores large numbers of bits. For m words of each n bits, memory can store total of $m \cdot n$ bits. To access each word, address input signals are defined. $\log_2(m)$ address inputs are required to select m words. Also, if there are k address inputs then the memory can have 2^k words.

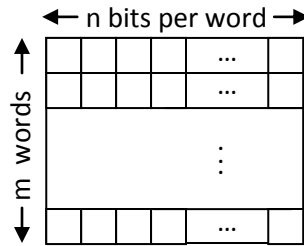


Figure 4.1: $m \times n$ memory

For example: A 4096 x 8 memory

- Stores 32768 bits
- Requires 12 ($2^{12} = 4096$) address signals
- Eight input/output data signals.

A memory access may refer to memory read – retrieve the word of a particular address, or memory write – store a word in a particular address. Control input signal r/w is used to indicate the type of access. Another control input signal, enable, which when asserted, is used to access the memory. Multiport memory supports multiple accesses to different locations simultaneously. Multiport memory systems have multiple sets of control lines, address lines, and data lines.

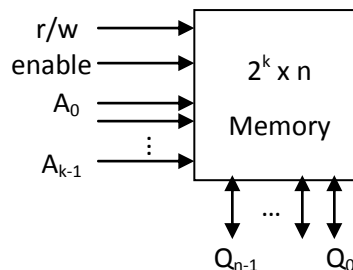


Figure 4.2: External View of Memory

Conventionally, ROM is referred as a memory that a processor can only read, and it holds stored bits even without a power source. Whereas, RAM is referred as a memory that a processor can both read and write but loses its stored bits if power is removed. But contemporarily, advanced ROMs, EEPROM and Flash, can be read as well as programmed and advanced RAMs, NVRAMs, can hold their bits even when power is removed. Advancement of memory have blurred the distinction

between the RAM and ROM. Different memories are differentiated based on two characteristics, write ability and storage permanence.

4.2 Memory Write Ability and Storage Permanence

Write Ability refers to the manner and speed that a particular memory can be written. Every memory must have a way to write bits onto it but the manner and speed of such writing varies among different memory. In-system programmable is used to categorize memories into two along the write ability axis. In-system programmable memory can be programmed by a processor whereas non in-system programmable memory must be programmed by some external means.

Range of Write Ability

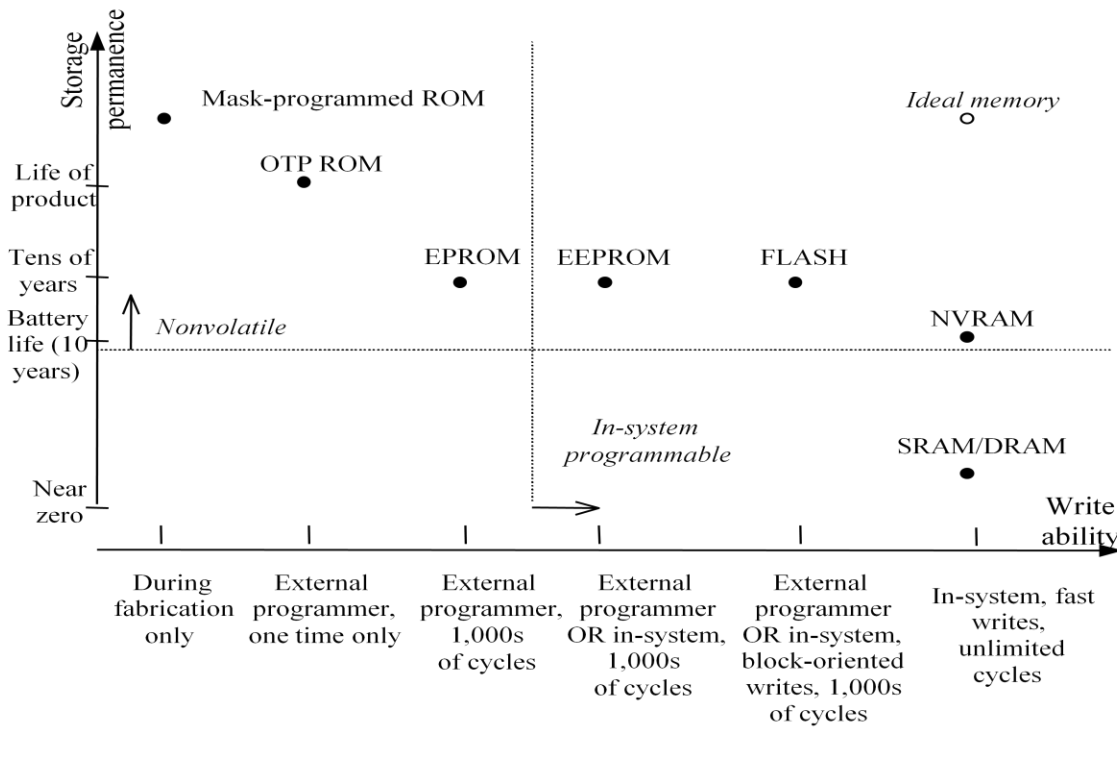
- High End – processor can write to memory simply and quickly by setting its address lines, data input bits and control lines appropriately. Example: RAM
- Middle Range – processor can write to memory a bit slower compared to high end. Example: EEPROM, FLASH
- Lower Range – special device called programmer is used to write into the memory. The device must apply suitable voltage levels to write to the memory. E.g.: EPROM, OTP ROM
- Low End – bits are stored during fabrication. Example: Mask-programmed ROM

Storage Permanence refers to the ability of memory to hold its stored bits after those bits have been written. Volatile and Nonvolatile are commonly used to divide memory types into two categories along the storage permanence axis. Non volatile memory can hold its bits even after power is no longer supplied. On the contrary, volatile memory requires continual power to retain its data.

Range of Storage Permanence

- Low End – memory in this range begins to lose its bits almost immediately after those bits are written and therefore it must be refreshed periodically. Example: DRAM
- Lower Range – memory holds bit as long as power is applied to the memory. Example: SRAM
- Middle Range – memory in this range holds bits for days, months, or even years after the memory power source has been turned off. Example: NVRAM
- High End – memory in this end will never lose its bits, as long as the memory chip is not damaged. Example: Mask Programmed ROM

Two characteristic, write ability and storage permanence, are important and desired in any system but it creates a trade-off. Write ability and storage permanence tend to be inversely proportional to one another. Moreover, highly writable memory, usually, requires more area and/or power than less-writable memory.



Write ability and storage permanence of memories, showing relative degrees along each axis (not to scale).

Figure 4.3: Various memories based on write ability and storage permanence

4.3 Common Memory Types

Read Only Memory (ROM): It is a nonvolatile memory, that can be read from but cannot be written to, by a processor, but it can be programmed by setting the bits within the memory. Traditionally, ROM is programmed off-line, when it is not actively involved within an embedded system.

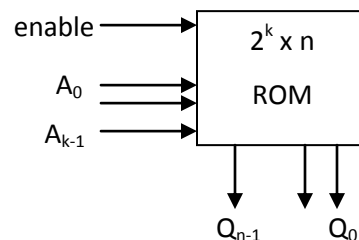


Figure 4.4: External Block Diagram

Uses of ROM

- store a software program for a general purpose processor
- To store constant data, like large lookup tables of strings or numbers
- to implement a combinational circuit

Example 1: Symbolic View of the internal design of an 8x4 ROM

- Horizontal lines = words (8), Vertical lines = data (4)
- Word line connected to data line via the programmable connections
- Circles on data and word lines are connected to represent high logic(1)
- Wired-OR represents all word lines are ORed together.
- If word 3 needs to be read then the input of decoder is set to 011 which makes the word 3 line high and other word lines low, since the data lines 0 and 3 are not connected to the high word 3 line, the output of the ROM will be 0110.

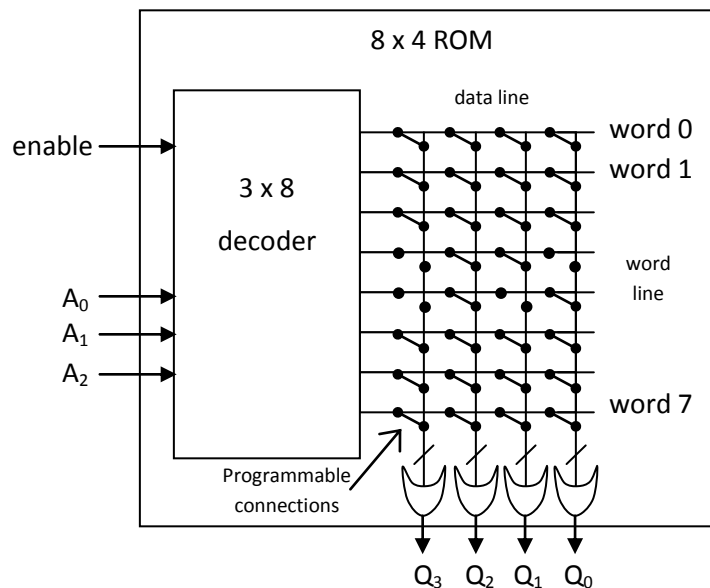


Figure 4.5: Internal View of an 8x4 ROM

Example 2: Implement the following combinational functions using a ROM

$$y = a'b'c' + a'bc' + ab'c + abc, z = a'b'c + a'bc' + a'bc + ab'c + abc$$

Solution: Three inputs a, b and c is taken as address lines. So, for three inputs the decoder of 3 x 8 must be used resulting in eight word lines. And there are two outputs, so there must be two data

lines. Hence, a ROM of 8 x 2 is required. Initially, the truth table, if not given, is formed from the given functions. The programming connections are done based on the output of the functions.

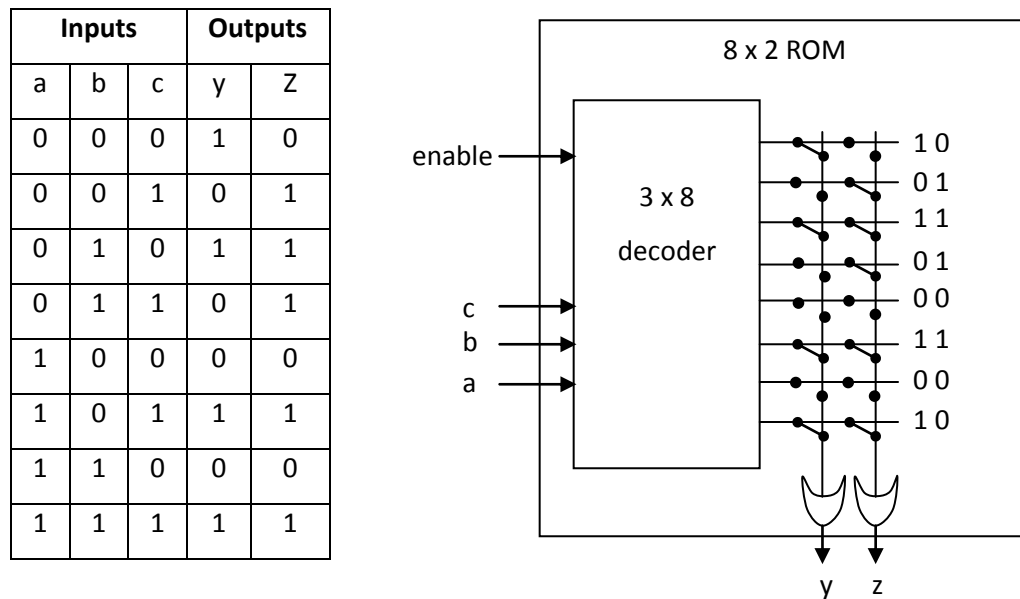


Figure 4.6: Truth table for given function and its implementation using ROM

TYPES OF ROM

A. Mask-Programmed ROM

- Connection is programmed during fabrication, by creating an appropriate set of masks.
- It has extremely low write ability. Once fabricated, its content cannot be reprogrammed or changed.
- It has highest storage permanence. Stored bits will never change unless the chip is damaged.
- It is used in such embedded systems whose design has been finalized and large numbers of unit are needed to be manufactured.

B. One-Time Programmable ROM – OTP ROM

- Connection is programmed using a device called programmer that configures each programmable connection according to the file provided by user. Programmer blows fuses by passing a large current wherever a connection is not required. The blown out fuses cannot be reestablished, hence it is referred as one time programmable ROM.
- It has lowest write ability of all PROMs, since it can be programmed only once.

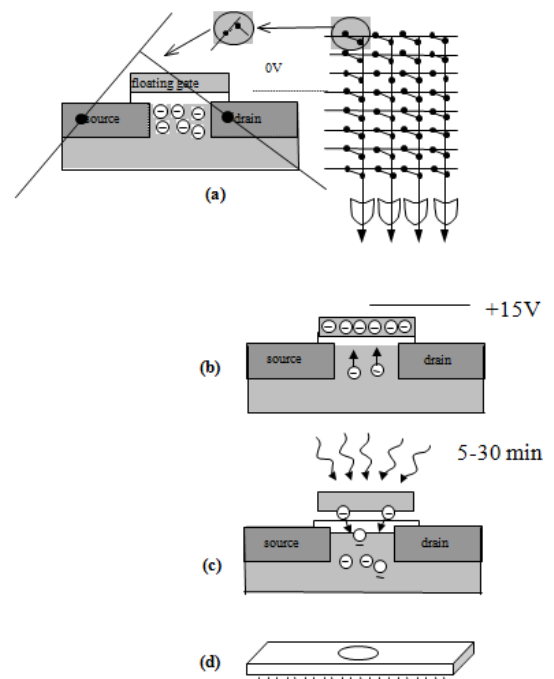
- It has very high storage permanence, since its stored bits won't change unless, some more fuses are blown out using programmer.
- It is cheap which makes it more suitable in final products compared to other types of PROM. Also compared to mask-programmed ROM, time-to-market constraints and unit costs make OTP-ROM a better choice.

C. Erasable Programmable ROM – EPROM

- EPROM uses a MOS transistor as its programmable component. The transistor has a floating gate surrounded by insulator. When high voltage (12v – 25v) is applied, it causes electrons to tunnel through the insulator into the gate. When the high voltage is removed, the electrons cannot escape and hence the gate has been charged and programming has occurred. To erase the program, the electrons must be excited enough to escape the gate which is done by exposing UV light for 5 – 30 minute. For the UV light to reach the chip, EPROMs are provided with a small quartz window in the package.
- Reading an EPROM is much faster than writing, since reading doesn't require programming.
- EPROMs have improved write ability and can be reprogrammed thousands of times.
- EPROMs have reduced storage permanence. They hold their stored bits for about 10 years.
- Electrical noise or radiations causes stored bits of the chip subject to undesirable changes and hence EPROMs are scarcely used in production. It offers a better choice in the testing phase of the system rather than in production.

• Internal Operation of EPROM

- a) Negative charges form a channel between source and drain storing a logic 1
- b) Large positive voltage at gate causes negative charges to move out of channel and get trapped in floating gate storing a logic 0
- c) Shining UV rays on surface of floating gate causes negative charges to return to channel from floating gate restoring the logic 1
- d) An EPROM package showing quartz window through which UV light can pass



D. Electrically Erasable Programmable ROM – EEPROM

- EEPROM is programmed and erased electronically, using higher than normal voltage. Electronic erasing requires seconds, rather than many minutes required for EPROMs. Moreover, individual words can be erased and reprogrammed in case of EEPROM, whereas EPROM can only be erased in their entirety.
- It is in-system programmable since circuit providing higher than normal voltage levels for erasing and programming is built into the embedded system. EEPROM is built with a built in memory controller which hides internal memory access details for the memory user and provides a simple memory interface to the user. The memory controller contains the circuitry and single purpose processor required to erase and program the word at the user specified address.
- EEPROM provides better write ability compared to EPROM, it can be reprogrammed tens of thousands of times.
- EEPROM has storage permanence on a par with EPROM, about 10 years.
- Writing is slower, since it involves the process of erasing and programming. Busy pin is available to indicate that the EEPROM is busy in writing.
- EEPROM can be used to serve as the program memory for a microprocessor. It can also be used to store data than an embedded system should save after the system is off.

E. Flash Memory

- It is an extension of EEPROM which uses the same floating-gate principle along with same write ability and storage permanence.
- It improves the performance of a system with its fast erase ability, in which large blocks of memory can be erased all at once.
- Writing to a single word in flash may be slower than writing to a single word in EEPROM, since an entire block will need to be read, updated and written back.

Random Access Memory – RAM

It is a memory that can be both read and written easily. Typically RAM is volatile, since it loses its content after the power is removed. The internal structure of RAM is comparatively complex than of ROMs.

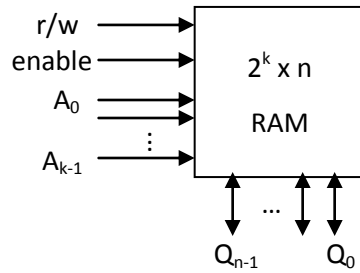


Figure 4.7: External View of RAM

Example 1: Sketch the internal structure of a 6 x 6 RAM

- Each word consists of a number of memory cells, each storing one bit.
- Each input data line and output data line is connected to every cell in its column.
- Output of a memory cell being ORed with the output data line of each column.
- The read/write input is connected to every cell

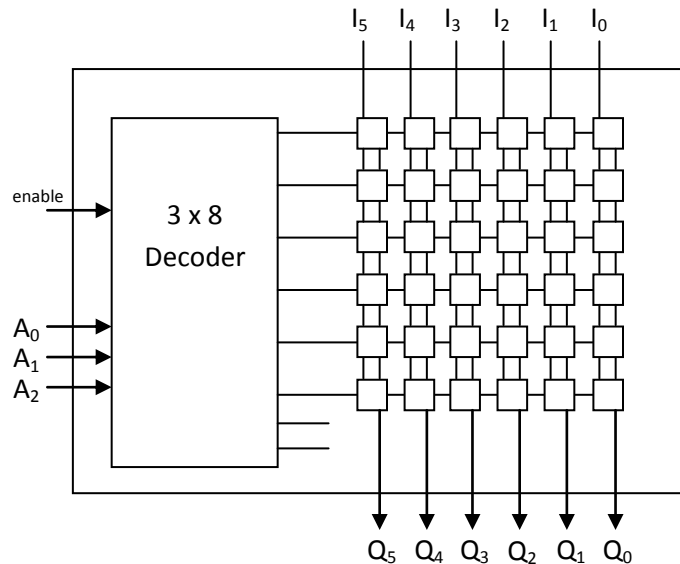


Figure 4.8: Internal Structure of 6 x 6 RAM

Types of RAM

A. Static RAM – SRAM

- Uses a memory cell consisting of a flip flop to store a bit
- Requires about six transistors to represent a single bit
- It holds data as long as power is supplied hence called static RAM.
- Generally used for high-performance parts of a system. E.g. Cache memory

B. Dynamic RAM

- Uses a memory cell consisting of a MOS transistor and capacitor to store a bit
- Requires only one transistor, resulting in more compact memory than SRAM
- Each cell must be charged (refreshed) regularly, since the charge stored in capacitor leaks gradually causing the loss of data.
- DRAM access tends to be slower than SRAM, since accessing a DRAM word results in the word's being stored in a buffer and then being written back to the word's cell.

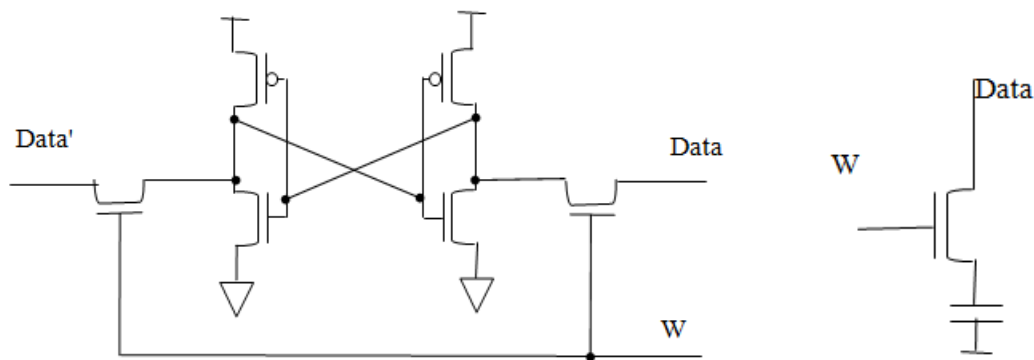


Figure 4.9: SRAM and DRAM

C. Pseudo-Static RAM – PSRAM

- These are DRAMs with a memory refresh controller built in.
- PSRAM may be busy refreshing itself when accessed, which could slow access time and add some system complexity.
- It is a popular low-cost high-density memory alternative to SRAM.

D. Nonvolatile RAM – NVRAM

- It holds data even after external power is removed.
 - **Battery-Backed RAM:** Contains a static RAM with permanent battery connected. When power is removed or drops below a certain threshold, the internal battery maintains power and the memory continues to store its bits. There is no limit on the number of times the Battery-Backed RAM can be written to.
 - **Static RAM with EEPROM or FLASH:** This type of NVRAM stores its complete RAM contents into the EEPROM just before the power is turned off. The data is reloaded into RAM after the power is turned back in.

Example: HM6264 and 27C256 RAM/ROM Devices

- Low-cost low-capacity memory devices used in 8-bit microcontroller-based embedded systems
- The first two numeric digits indicate whether the device is RAM (62) or ROM (27), whereas the subsequent digits give the memory capacity in kilobits.
- Placing a memory address on the address-bus and asserting the read signal output enable (OE) performs a read operation.
- Placing some data and a memory address on the data and address busses and asserting the write signal enable (WE) performs a write operation.

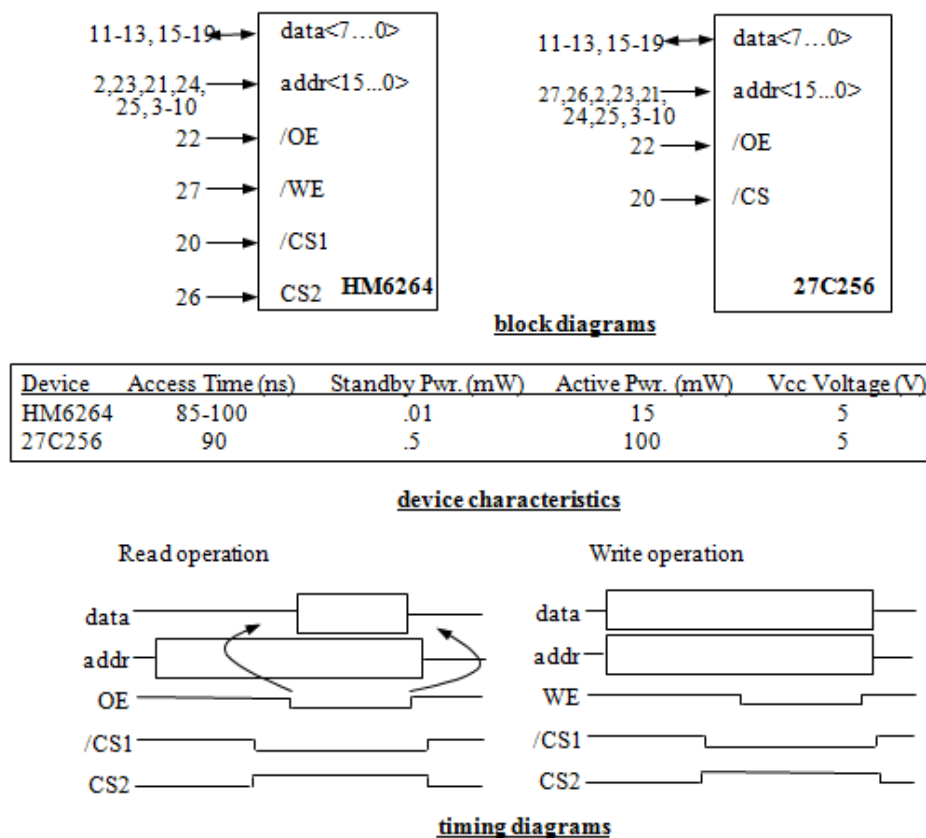


Figure 4.10: Example - HM6264 and 27C256 RAM/ROM Devices

4.4 Composing Memory

Composing Memory is needed when there is a need of particular-sized memory, which is not readily available. If the available memory is larger than required one, then we simply use the needed lower words of the memory and ignore the higher words which are not required. However, if the available

memory is smaller than needed, some more design procedures are needed to be followed. The various cases for composing memory have been discussed in the following paragraphs.

A. Case 1: To increase the width of words

When the number of words in the available memory is same to that of required one but the number of bits or width of word is not enough then the width must be increased. To do that, the available memories are connected side by side as shown in the given example.

Example 1: Compose 1K x 8 ROMs into a 1K x 32 ROM

Analysis: The available ROM 1K x 8 and required ROM of 1K x 32 have same number of words but width is different. The number of ROM to be placed side by side is given by n.

- $n = \text{width of required ROM} / \text{width of available ROM} = 32/8 = 4$
- Address line = 1K = 1024bytes = $2^{10} = 10$ address lines
- Data line = 8 lines

Hence, four 1K x 8 ROMs are placed side by side to compose 1K x 32.

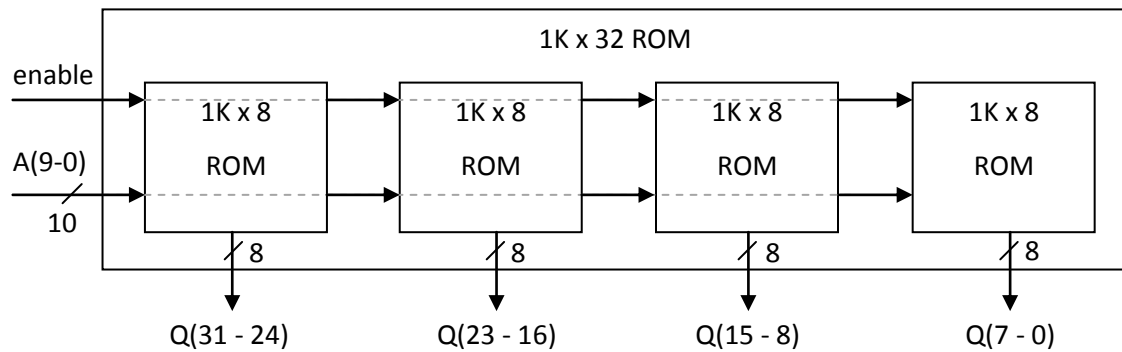


Figure 4.11: Composing 1K x 32 ROM from 1K x 8 ROM

B. Case 2: To increase the number of words

When the width of the word in the available memory and required memory is same but the number of words are different then the words must be increased. We connect the ROMs top to bottom and data line of each ROM is ORed. Since the number of words has to be increased, extra high-order address is required to select the particular ROM which can be implemented by using appropriate decoder.

Example 1: Compose 1K x 8 ROMs into a 4K x 8 ROM

Analysis: The available ROM 1K x 8 and required ROM 4K x 8 have same width of 8 bits but the number of words is different. Number of ROMs and the size of decoder can be determined as

- $N = \text{number of words in required ROM} / \text{number of words in available ROM} = 4K/1K = 4$
- Decoder: It must be able to select 4 ROM, so 2 x 4 decoder must be used.
- Higher address bits = $\log_2(4K) - \log_2(1K) = \log_2(2^{12}) - \log_2(2^{10}) = 12 - 10 = 2$ bits or lines
- Total Address line: $4K = 2^{12} = 12$ address lines, and 10 lines (A_9 to A_0) are connected to each ROM. 2 higher address is represented by inputs of decoder.

Hence, four ROM must be connected top to bottom and data line of each ROM is ORed. Decoder of 2 x 4 is used to select a particular ROM.

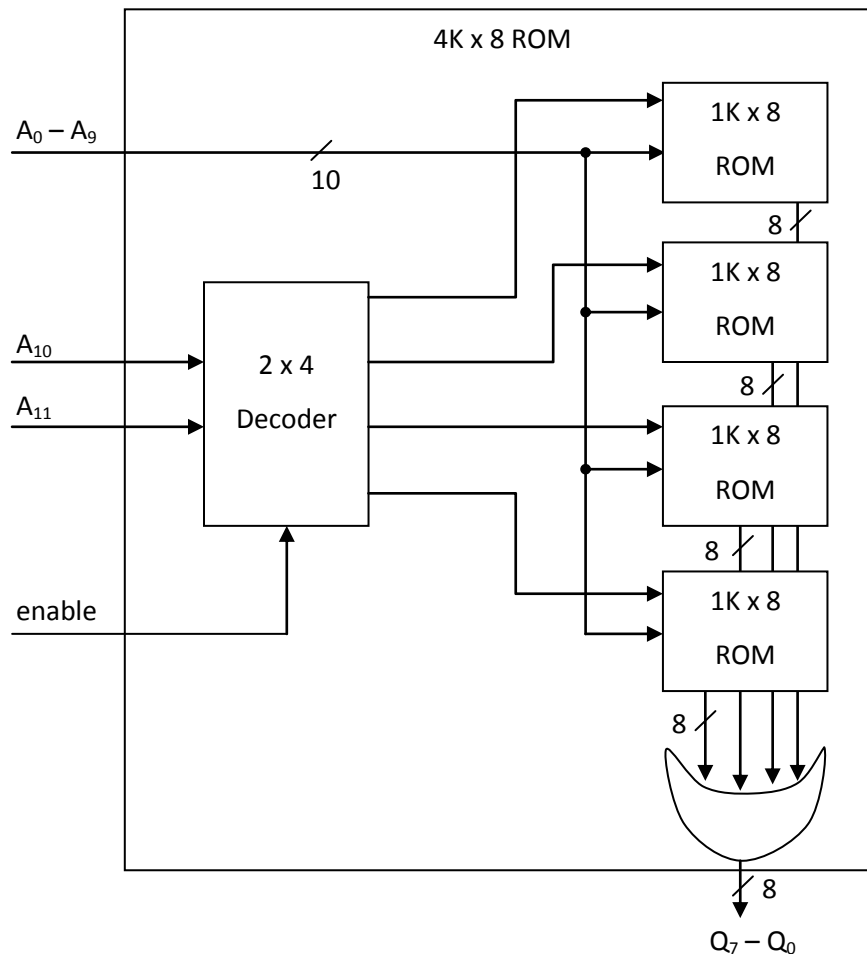


Figure 4.12: Composing 4K x 8 ROM using 1K x 8 ROM

C. Case 3: To increase both, number of words and word width

When the width of the word as well as the number of words in the available memory and required memory are different then the technique used in case 1 and case 2 must be combined. Initially, the number of words is increased and then the top-bottom set of ROMs with ORed data lines are placed side by side to increase the word width.

Example 1: Compose 1K x 8 ROMs into a 4K x 16 ROM

Analysis: The available ROM 1K x 8 and required ROM 4K x 16 differ in number of words as well as word width.

- Increase number of words: $4k/1k = 4$, Four ROMs are required with 2 x 4 decoder. 4K represents 12 address lines, 10 lines connected to every ROM and 2 lines represented by inputs of decoder.
- Increase word width: $16/8 = 2$, Four set of ROMs are repeated two times and placed side by side.

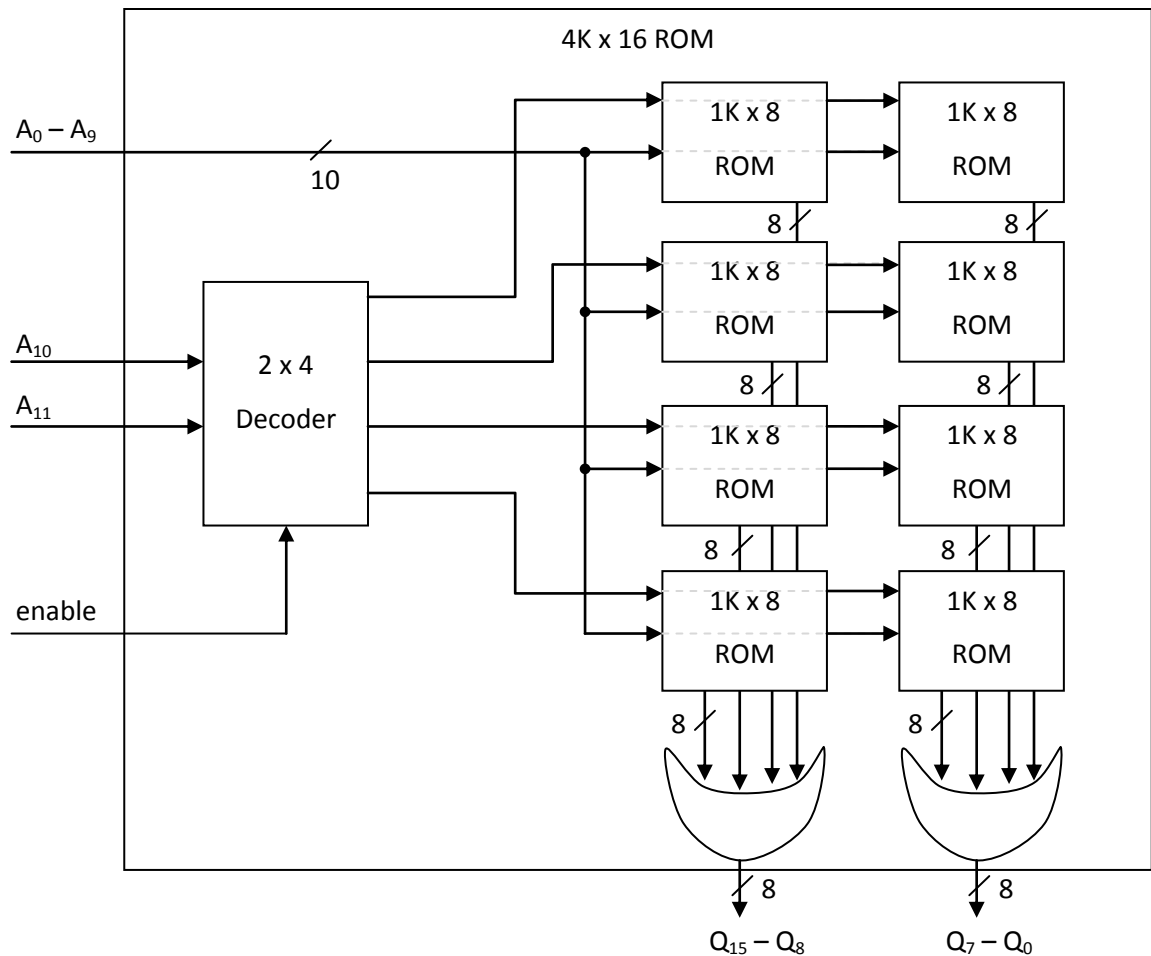


Figure 4.13: Composing 4K x 16 using 1K x 8 ROM

4.5 Memory Hierarchy and Cache

Memory Hierarchy

A system cannot be implemented with only fast memory as it makes the system very expensive. Also the use of only slow and low cost memory will make system very inefficient. So, the concept of

memory hierarchy comes into action in which a system is more likely to implement slow but high capacity memory for storage along with fast but small memory for high speed processing. Memory hierarchy defines the level of memory based on cost per bit, capacity and access time. As we move down the hierarchy, capacity increases, access time increases and cost per bit decreases.

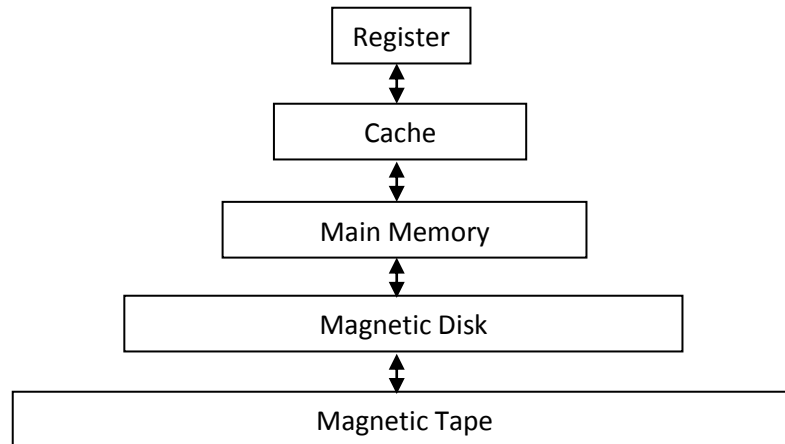


Figure 4.14: Memory Hierarchy

Cache Memory

Cache is a small but fast memory which contains a copy of portions of main memory to expedite operations of the system. Cache is designed using static RAM which makes it faster as compared to main memory. The access time for cache can get as low as one clock cycle while main memory access requires several cycles. So, the instructions and data which are supposed to get accessed frequently are placed in cache memory. Hence, the average access time is reduced resulting in improved performance.

During cache operation, the processor first checks the required word in cache. If it is available (cache hit), the word is delivered to the processor. But, however, if the word is not available (cache miss) in cache then the corresponding block of main memory is read into cache. And finally the word is made available to the processor. This operation leads to various cache design issues which are discussed in the following paragraph.

A. Cache Mapping Techniques

Cache memory is very small as compared to main memory. And all blocks of main memory cannot be assigned to cache memory at once. So, cache mapping techniques are required to assign particular block of main memory to the appropriate line in cache memory. There are basically three types of mapping techniques which are discussed below.

Direct Mapping

In this technique, main memory block is assigned to a fixed cache line. The cache stores the content of main memory, the tag and the valid bit. Here, the memory address is divided into the tag, the index and the offset. The index, which is defined by the cache size, represents the cache address. Index is used to select the particular cache line. The tag from main memory address is compared to the tag stored in cache. In case the tag matches, the data from the cache line is accessed. However, a single cache line can store few blocks of main memory. So to select a particular block, the offset part of main memory address is used. The valid bit in cache is used to indicate the validity of data stored in the cache slot.

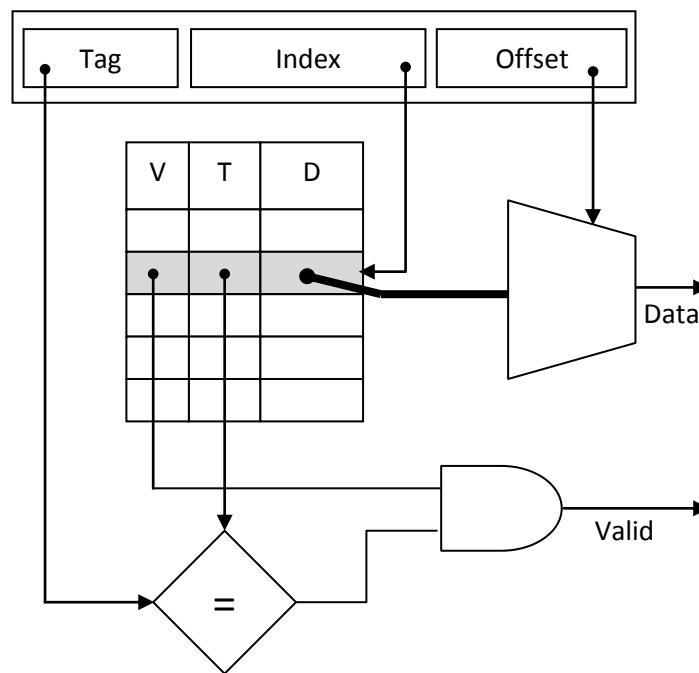


Figure 4.15: Direct cache mapping

Direct cache mapping is easy and simple in implementation. However, when two blocks of main memory which are assigned to a particular cache line are to be accessed frequently, then cache miss occurs repeatedly. This problem is commonly referred as thrashing. Also, replacement algorithm cannot be used, since main memory blocks are mapped to a fixed cache line.

Fully Associative mapping

In this mapping, main memory block can be assigned to any slot of cache line. The main memory address is divided into tag field and offset field. The tag from main memory is compared to each tag in the cache line. After the tag matches the offset is used to select a particular word in cache line.

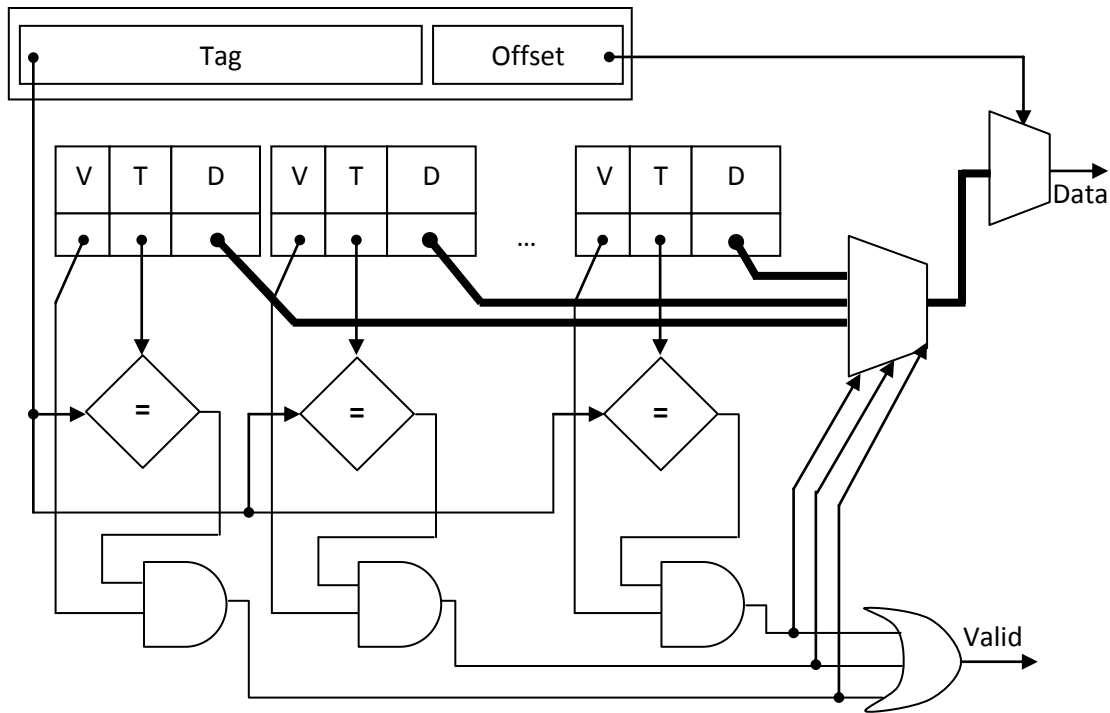


Figure 4.16: Fully associative cache mapping

Fully associative mapping provides high flexibility as block of main memory can be assigned to any cache line. However, the comparison logic is required for each cache line which makes this mapping method complex and expensive to implement. Miss rate can increase if frequently required block is replaced, so appropriate replacement algorithm must be utilized for efficient cache implementation.

Set Associative Mapping

It is a compromise mapping which, somehow, follows both direct and fully associative mapping. The cache is divided into sets, each with number of cache lines. A cache with a set of size N is called an N-way set associative cache. Each block of main memory can be mapped to particular line of any sets (fixed line but varying sets) or any lines of particular set (fixed set but varying lines). Taking former case into consideration, the main memory address is divided into tag, index and offset. The index field is used to select the fixed cache line, and the tag field of main memory is compared to tag of each sets. When the particular set is selected, the offset is used to select the particular word from the set in which the tag matches.

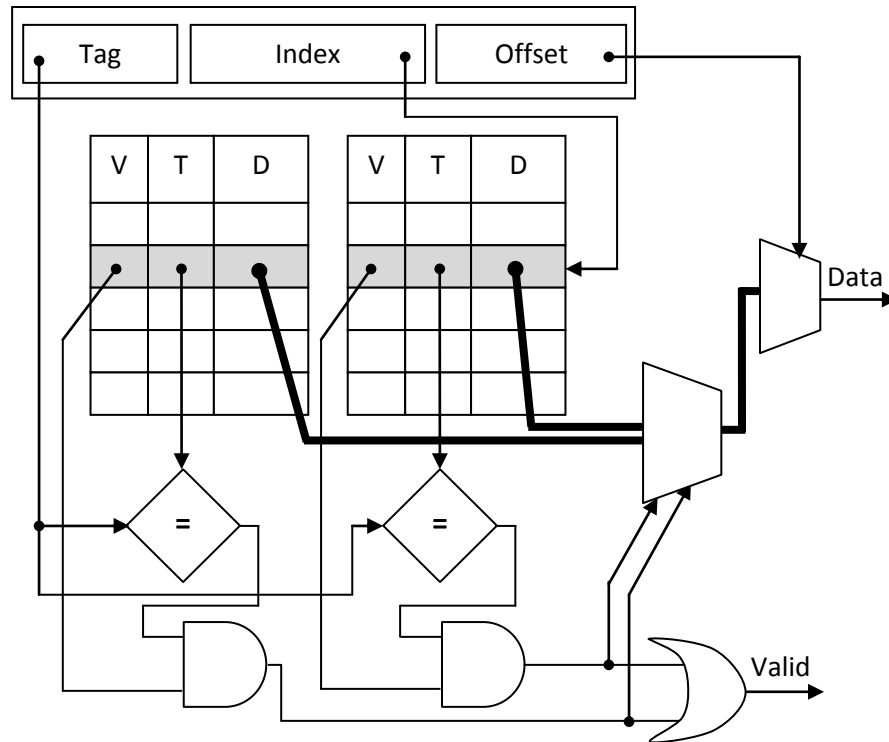


Figure 4.17: Two-way set associative

Set associative cache mapping is more flexible and can reduce cache misses as compared to direct mapping. Though the block of main memory is assigned to fixed cache line, block can be assigned to any sets of cache line. And proper implementation of cache replacement can be used to increase cache hit rate. Also the comparison logic is not required for every cache line rather is required for only available sets which reduce the complexity and expense for implementing comparison logic.

B. Cache-Replacement Policy

When cache is full and new main memory block is to be assigned to the cache then certain technique must be used to choose which cache line should be replaced. This mechanism of replacing the existing block by new set of blocks is referred as cache-replacement policy. In direct mapping the main memory block always maps to the fixed cache line, so replacement is fixed. But fully associative and set associative can follow various replacement algorithms. Least Recently Used (LRU), First In First Out (FIFO), Least-Frequently Used (LFU) and Random are few commonly used replacement techniques.

- **Random replacement** replaces the block randomly without following any specific algorithm.
- **Least Recently Used (LRU)** algorithm is based on time in which the block not accessed for longest time is replaced by the new block.

- **First In First Out (FIFO)** method uses queue mechanism to replace the first entered block. Each block is pushed into the queue when accessed. And when replacement is required the blocks are popped out from the queue.
- **Least Frequently Used (LFU)** technique is based on number of time the block is accessed. The block which is accessed less number of times is replaced.

C. Cache Write Techniques

A mechanism is required when content of cache is changed by the processor and the change must be updated to the corresponding main memory block. This technique of updating the main memory after change in cache is referred as cache write policy. There are two common cache write policy; write-through and write-back.

Write-through is a technique in which the main memory is updated immediately after the content in cache is changed. This technique is easier to implement but the processor has to wait for slower main memory frequent access. Also there are chances of unnecessary writes resulting in substantial memory traffic. For example when a particular value is changed four times, the last updated value must only be updated in the main memory. But the memory is updated four times for every change causing unnecessary memory access.

Write-back policy allows main memory to be updated only when cache line is to be replaced. Extra bit is associated with each cache line to represent whether the content of cache line is changed or not. Based on that extra bit the corresponding main memory block is updated when cache line is about to be replaced. Extra bit and update checking increase system complexity; however, it reduces number of slow main memory access and avoids memory congestion.

D. Cache Impact on System Performance

The performance of system is directly related to design and configuration of caches. The total size of cache, degree of associativity, and the data block size are important parameters that have direct impact on performance.

Cache size is the total number of bytes that the cache can store. The tags and extra bits, which do not contribute to the size of the cache, are also stored in cache along with the data of main memory block. Increasing the size of cache results in lower miss rates, however the access of data from the cache will be slower. So, larger cache size does not necessarily mean better performance.

Degree of associativity is related to number of sets used in set associative cache implementation. Increasing the number of sets will improve the hit rate. However, additional logic requirement will increase the access time latency.

Cache line size represents the size of each block in cache that holds the block of data of main memory. When line size is increased, the main memory access time is, obviously, reduced but only at the expense of more complex multiplexing circuitry which increases the access latency.

Example: Effect of cache size on system performance

Case I: Cache size = 2Kbytes, miss rate = 15%, hit cost = 2 cycles, miss cost = 20 cycles

$$\text{Average cost of memory access} = (0.85 \times 2) + (0.15 \times 20) = 4.7 \text{ cycles}$$

Case II: Cache size = 4Kbytes, miss rate = 6.5%, hit cost = 3 cycles, miss cost = 20 cycles

$$\text{Average cost of memory access} = (0.935 \times 3) + (0.065 \times 20) = 4.105 \text{ cycles}$$

Case III: Cache size = 8Kbytes, miss rate = 5.565%, hit cost = 5 cycles, miss cost = 20 cycles

$$\text{Average cost of memory access} = (0.94435 \times 5) + (0.05565 \times 20) = 4.8904 \text{ cycles}$$

In case II, increase in cache size, certainly, improved the performance as average cost of memory access is decreased. However, in case III, increase in cache size added more cycles for memory access in average.

Advanced RAM

A. Fast Page Mode Dram (FPM DRAM)

FPM DRAM is asynchronously controlled which is designed with some improvements on the basic DRAM architecture. In this design, each row of the memory bit-array is viewed as a page which contains multiple words. Each word is addressed by a unique column address. In its operation, first the row or page address is sent and then the corresponding column address must be sent to read a particular word. In each memory cycle, three data words can be read consecutively by providing their corresponding column address. Hence, it eliminates the requirement of extra cycle as three cycles would have been required to read three words.

B. Extended Data Out DRAM (EDO DRAM)

EDO DRAM is similar to FPM DRAM with additional feature that reduces the read/write latency. Here, new access cycle can be started while keeping the data output of previous cycle active. In simple words, new column address can be sent while reading previously selected word from the

memory. This results in overlapping of the operation which reduces the latency of memory access. However, extra output latch must be introduced in the architecture.

C. Synchronous DRAM (SDRAM)

In SDRAM, the information is latched to and from the controller on the active edge of the clock signal. The time required to detect the strobe signals in asynchronous DRAM is eliminated by SDRAM. This DRAM architecture can have additional column address counter which holds the starting address of the data to be accessed. This counter is incremented internally to provide new data in each clock cycle as long as the data required are consecutive memory locations. The enhanced synchronous DRAM (ESDRAM), is the improved version of the SDRAM. ESDRAM provides faster clocking and lower latency in reading and writing data.

D. Rambus DRAM (RDRAM)

Rambus represents the bus interface architecture which uses multiplexed address/data lines to connect the processor to the RDRAM device. RDRAM may be further divided into number of banks with each remain open for access. Multiple open page scheme and fast bus I/O can result in high throughput. However, as compared to other standards, Rambus showed increase in latency, heat output, complexity, and cost. Requirement of heat-spreaders along with packet demultiplexors makes it more complex while manufacturing. More complex interface circuitry and more number of memory banks increased the size and resulted to become expensive.

E. Double Data Rate SDRAM (DDR SDRAM)

The DDR SDRAM is capable of making higher transfer rates with more strict control of the timing of the data and clock signals. The interface transfers data on both the rising and falling edges of the clock signal to double the data bus bandwidth. DDR SDRAM also known as DDR1 was replaced by DDR2 which operated on same principle but for higher clock frequency and produced double throughput as compared to DDR1. Similarly, DDR3 and DDR4 offered better performance for increased bus speed and new features.

Memory Management Unit (MMU)

Memory Management Unit is a processor which translates the logical address to physical memory address. MMU has important role in handling DRAM refresh, bus interface and arbitration used in memory. In addition, it takes care of memory sharing among multiple processors. Contemporary CPUs have built-in MMU as a part of processor.

- **Communication Basics**
- **Microprocessor Interfacing**
 - I/O Addressing
 - Interrupts
 - Direct Memory Access
- **Arbitration**
- **Multilevel Bus Architectures**
- **Advanced Communication Principles**
- **Serial, Parallel and Wireless Protocols**

5.1 Communication Basics

Basic terminology

Wires are the connecting lines of two terminals in communication system. It may be uni-directional or bi-directional. A single line can be used to represent multiples wires with the help of small angled line drawn through it.

Bus refers to the set of wires with a single function. Address bus for address, data bus for data are two examples of single functioned buses. Bus can also be the entire collection of wires. System bus, for instance, consists of address, data and control lines.

Port is the actual conducting device on periphery which connects bus to processor or memory or other devices. Port is a medium through which a signal is input to output from the processor. Port is also referred as pin which extends from the IC package and that can be plugged into a socket (IC base) on a printed circuit board. Metallic balls instead of pins may be present. However, metal pads are more common these days.

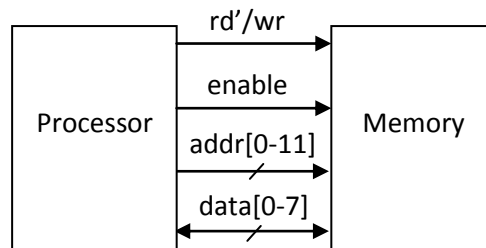


Figure 5.1: A simple bus example

Timing diagram is a diagrammatic representation for describing hardware protocol. In the diagram, time proceeds to the right along x-axis. It represents state of control lines or data lines. The control lines may be either low or high, whereas the data lines – address or data -- can be valid or not valid. Active high means that a one on the line makes it active while active low means that a zero on the line makes it active. Asserting a line means making it active and de-asserting the line deactivates the line. A protocol may have several sub-protocols which are also called bus cycle or transaction. A bus cycle may consist of several clock cycles.

Example: Timing diagram for read protocol

The timing diagram of memory read protocol gives the following information to the designer

- The processor must set the *rd'/wr* line low for a read operation

- Address of memory must be placed on *addr* line for atleast t_{setup} time before setting the *enable* line high.
- Setting *enable* line high will cause memory to place the data on the *data* line after at time t_{read} .

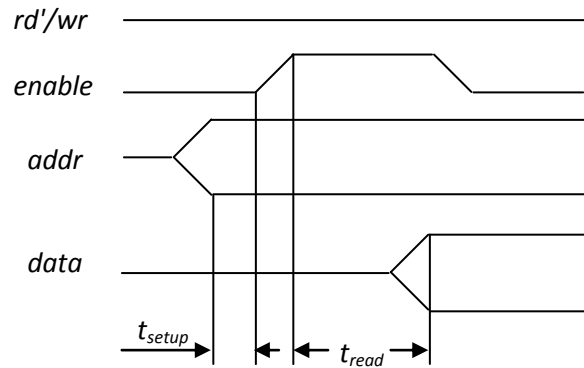


Figure 5.2: Timing Diagram: Memory read protocol

Basic Protocol Concepts

An **actor** is a device that can be processor or memory which takes part in data transfer. Actor can be a master or a slave. A master initiates the data transfer whereas a slave responds to the initiation request.

Data direction represents movement of data among actors. The direction of data is independent of type of actor. Either master or slave can send or/and receive data.

Addresses represent a special kind of data which specify a location in memory, a peripheral, or a register within a peripheral. A protocol often includes both an address and data. In every memory access protocol, the address specifies the location where the data should be read from or written to in the memory.

Time multiplexing represents a technique in which the multiple sets of data are sent one at a time over the shared line. Number of wires requirement can be reduced to a single line at the expense of time. The following figures show the examples of time multiplexing. In both cases, single bus is used to send multiple data at different time instant.

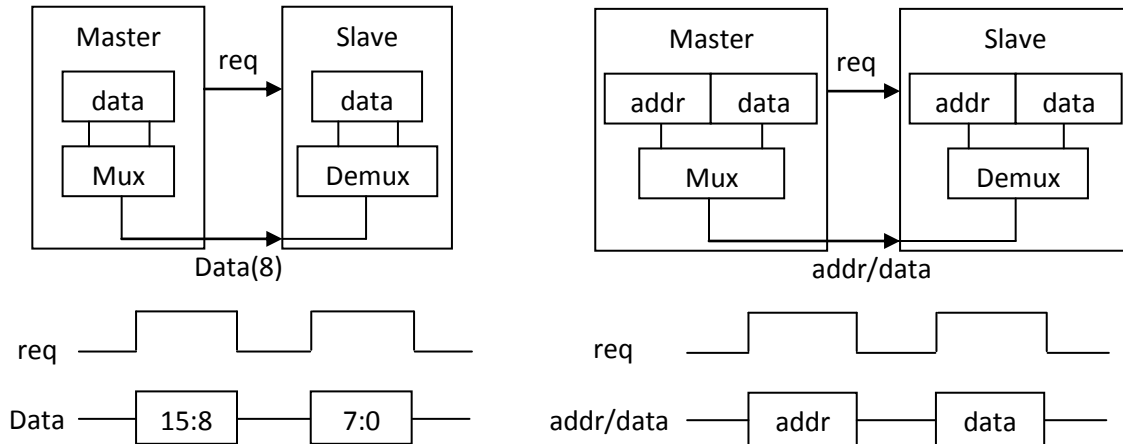


Figure 5.3: Time-multiplexed data transfer: data serialization and address/data muxing

Control methods are schemes for initiating and ending the data transfer. Strobe and handshake are two common control methods.

Strobe Protocol

In strobe protocol, master uses a control line and activates it to initiate the data transfer. Then the slave has certain time to put data on data bus. Assuming data to be valid, master reads the data from data bus and deactivates its control line. And both actors are ready for next data transfer. The main disadvantage of strobe protocol is that the master that initiates the transfer has no way of knowing whether the slave has received the data or not.

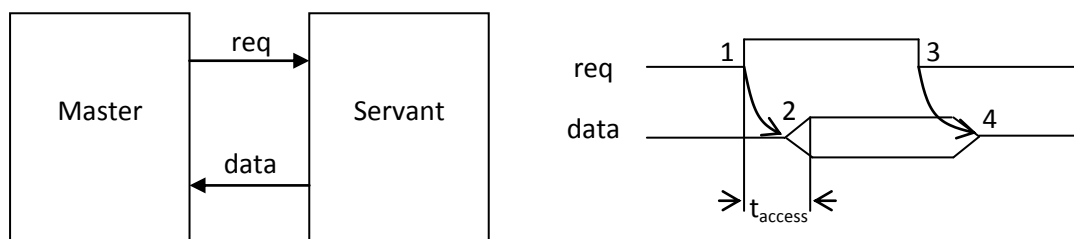


Figure 5.4: Strobe Protocol

The flow in timing diagram can be explained as:

1. Master asserts *req* to receive data
2. Servant puts data on *data* line within time t_{access}
3. Master receives data and deasserts *req*
4. Servant ready for next request

Handshake Protocol

In this protocol, servant uses extra line to acknowledge that the data is ready. Initially, master asserts request line to start the transfer. Then the servant, taking its time to put data on data line, asserts acknowledge line to inform the master that the data is ready. Next, the master reads the data from the data line and deasserts the request line which is followed by slave deasserting acknowledge line. Finally the transfer is complete and both actors are ready for next transfer. Though the protocol is somewhat complex, it is more reliable compared to strobe protocol as data availability is confirmed by the sending device.

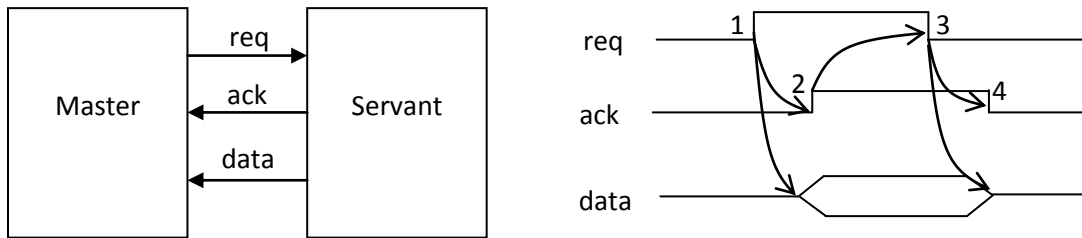


Figure 5.5: Handshake Protocol

The flow, as indicated by numbers, in timing diagram can be summarized as:

1. Master asserts *req* line to receive data
2. Servant puts data on *data* line and asserts *ack*
3. Master receives data from data bus and deasserts *req*
4. Servant ready for next transfer

Strobe/Handshake Compromise

A compromise protocol can be used to achieve the speed of strobe protocol and varying response time tolerance of handshake protocol. As represented in figure 5.6, servant can use wait line, if it is not ready to put data on data line.

- If the servant can put data within time t_{access} then it follows strobe protocol representing fast response. And wait line remains unused in this protocol.
- If the servant can't put the data within t_{access} time then it asks master to wait longer by asserting wait line. After the data is ready, the wait line is deasserted by servant and master receives the data. And it represents slow response as master has to wait for certain time for the data transfer

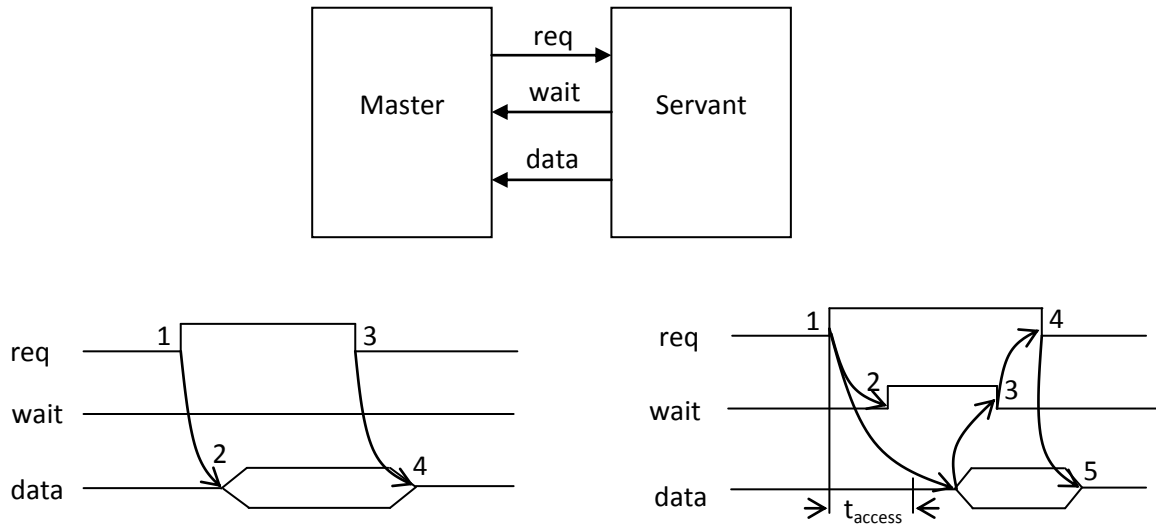


Figure 5.6: A strobe/handshake compromise: fast and slow response

The flow as indicated by numbers in timing diagram can be summarized as:

For fast response

1. Master asserts *req* line to receive data
2. Servant puts data on data bus within time t_{access} , *wait* line remains unused
3. Master receives data and deasserts *req*
4. Servant ready for next request

For slow response

1. Master asserts *req* to receive data
2. Servant can't put data within t_{access} , asserts *wait* line
3. Servant puts data on bus and deasserts *wait*
4. Master receives data and deasserts *req*
5. Servant ready for next request

Example: The ISA Bus Protocol – Memory Access

The Industry Standard Architecture bus protocol is common in systems using an 80x86 microprocessors. The processor uses 20-bit memory address and follows compromise strobe/handshake protocol. If the memory is not ready then the processor inserts wait cycles. Four cycles is default for the operation to complete. For the read operation, in the first clock cycle the processor puts address on the address line and asserts address latch enable signal. During second and third clock cycle, the processor asserts memory read signal. After third clock cycle, the data is

available on data lines. Finally all signals are deasserted at fourth clock cycle. The timing diagram for memory read operation and memory write operation is shown in the figure below.

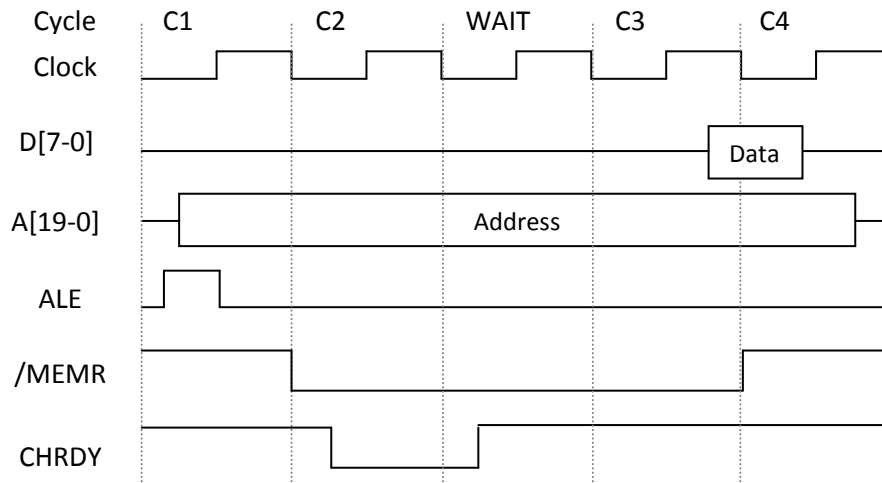


Figure 5.7: ISA bus protocol – read bus cycle

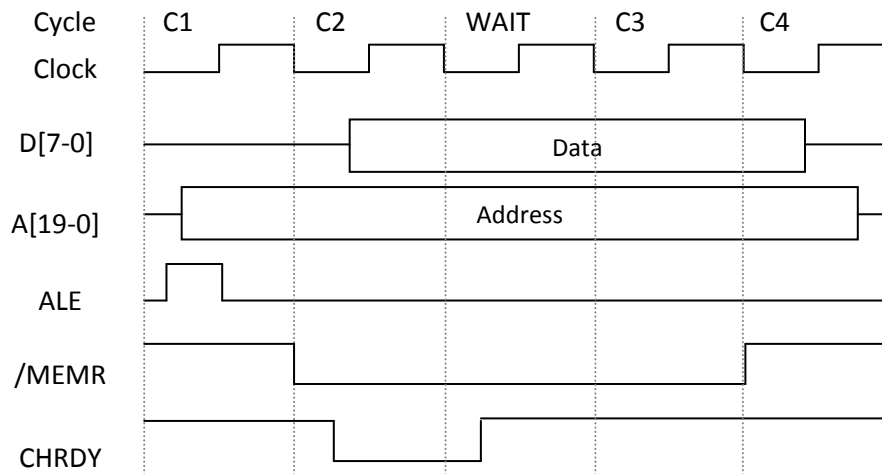


Figure 5.8: ISA bus protocol – write bus cycle

For Write Operation

- In C1, processor puts 20 bit address memory address on the address line and asserts ALE signal.
- During C2 and C3, the processor puts the data on the data line and asserts MEMW signal to enable write operation. However, if the memory, when not ready, deasserts CHRDY signal in C2 then processor inserts wait cycles until CHRDY is reasserted.
- In cycle C4, all signals are deasserted.

5.2 Microprocessor Interfacing

A. I/O Addressing

Port based I/O

In port based I/O, a port can be directly read from or written into with the help of processor instructions. It is also referred as parallel I/O. Generally the devices may be provided with one or more N-bit ports to facilitate port based I/O and each port is bit addressable. For example, 8051, AVR microcontrollers have 8 bit I/O ports. In 8051, $P1 = 0xF7$ statement will write into Port 1 of microcontroller. Also, for bit addressing, $P2.4 = 1$ will set the pin 4 of Port 2.

The port based I/O can be extended using appropriate peripheral which extends the number of available ports from four to six. Each port on peripheral is associated with a register that can be read or written into by the processor.

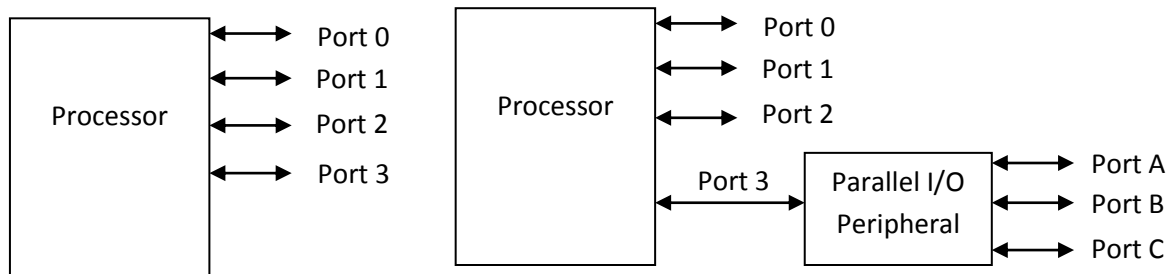


Figure 5.9: Port Based I/O and Extended Parallel I/O

Bus based I/O

In bus based I/O, the processor has address, data and control lines for I/O addressing. The communication protocol is built into the processor. A single instruction is available which causes the hardware to write or read data. If a system with bus based I/O requires parallel I/O then parallel I/O peripheral can be connected to the system bus.

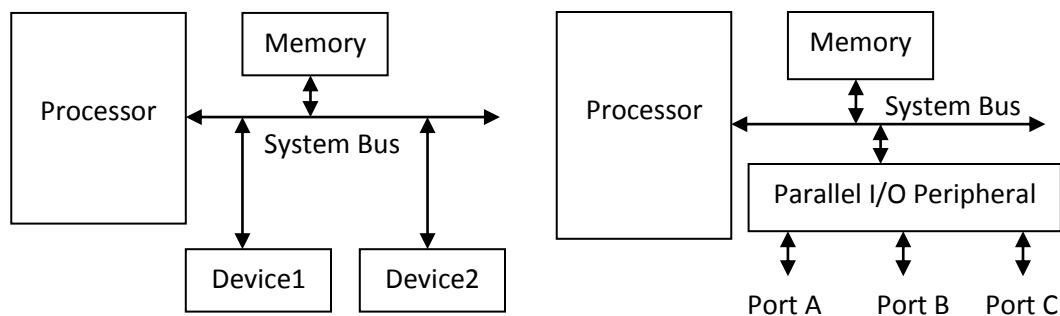


Figure 5.10: Bus-based I/O and Extended bus-based I/O with parallel I/O peripheral

Memory-Mapped I/O

Memory-mapped I/O is a type of bus-based I/O addressing for a processor to communicate with peripherals in which peripherals are addressed using the specific existing address space. The total address space is divided into memory address and peripheral address. Hence, there is loss of memory addresses to peripherals. Also, no special instructions for peripherals are required for data transfer, since instructions like MOV used for memory will also work with peripherals.

Example: A bus with 16 bit have total of 65536 addresses. So, lower 32768 addresses may correspond to memory address while upper 32768 correspond to I/O addresses.

Standard I/O

Standard I/O is a type of bus-based I/O addressing in which extra control line (M/IO) is used to indicate whether the address represents memory location or peripheral. Memory locations and peripherals use all sets of address for addressing, so there is no loss of memory addresses. This addressing, however, requires special instructions. MOV, LOAD instructions for memory while IN, OUT for peripherals. Also the address decoding logic for peripherals is simple as the high order address bits can be ignore when the number of peripherals is less.

Example: A bus with 16 bit have total of 65536 addresses. All 65536 addresses can be used to address memory and peripheral. The M/IO line is used to select either memory or peripheral. If M/IO is zero then the address on the address bus corresponds to a memory address.

Example: The ISA Bus Protocol – Standard I/O

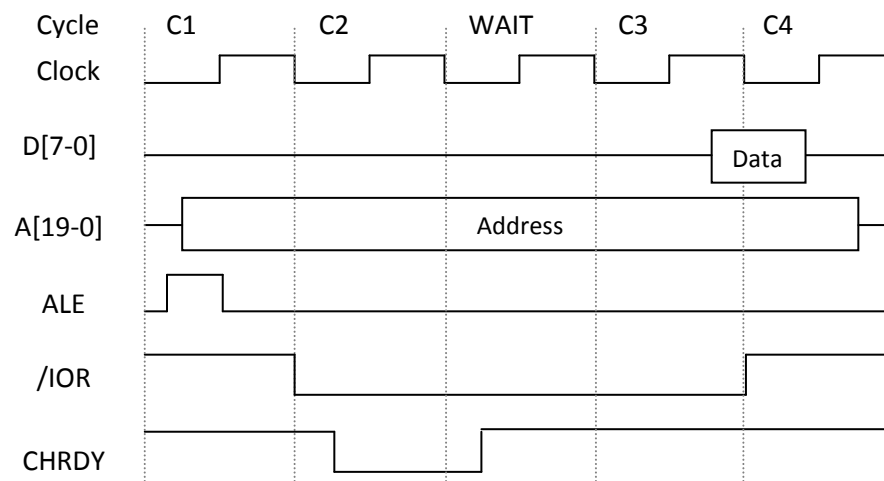


Figure 5.11: ISA bus protocol for standard I/O

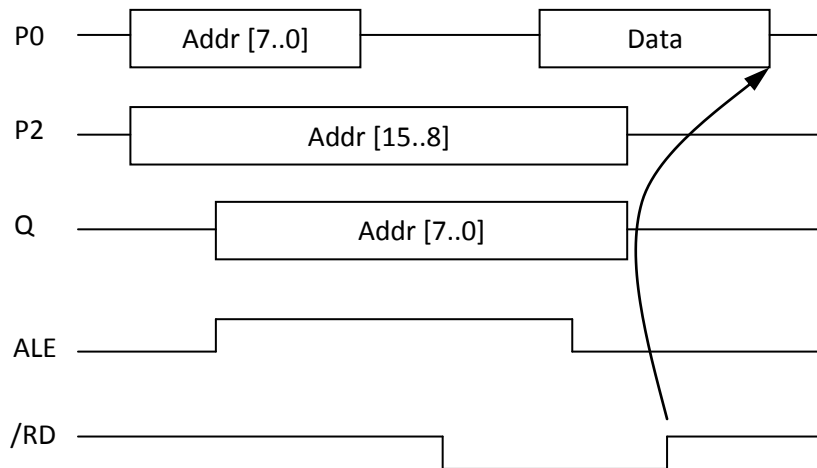
Example: A Basic Memory Protocol

Figure 5.12: A Basic Memory Protocol for 8051 microcontroller: Timing diagram for read operation.

Read Operation in 8051 microcontroller

- Microcontroller places source address on ports P2 and P0. Port 2 holds 8 most significant address bits and retains its value throughout the read operation. Port 0 holds the eight least-significant address bits which is stored using a 8-bit latch.
- The ALE signal is used to trigger the latching of port P0. And controller asserts high impedance on P0 to allow memory device to drive it with requested data. The memory outputs valid data as long as RD signal is asserted. The microcontroller reads the data and deasserts its control can port signals.

B. Interrupts – Interrupt Driven I/O

The peripherals may require service from the processor which is very much unpredictable. So there is an issue on how to serve the peripherals by the processor as it remains busy on its own task. Polling and interrupts are two basic methods to address that issue.

1. **Polling** is a method in which the processor checks for service requirement of every peripheral. This method, though, is easier and simple to implement, the repeated checking, however, wastes many clock cycles which could have been used to do certain useful work.
2. **Interrupt** is a feature of the processor through which the peripherals can request for service even when processor is busy in its own task. For external interrupt, there is always a pin available to implement interrupt feature. Whenever the interrupt pin is asserted, the processor jump to a particular address at which the routine for the interrupt is stored. Interrupt

overcomes the limitations of polling, but interrupt, in itself, is the type of polling. The pin is checked after the execution of every instruction, so it does not require extra clock cycles.

Interrupt address vector represents the address in which the interrupt service routine (ISR) resides. Fixed interrupt and vectored interrupt are two common methods by which the processor obtains the address of ISR.

In **fixed interrupt**, the address of subroutine is built into microprocessor and remains fixed. Programmer simply has to store the ISR at that location or can put jump instructions to move to actual location of ISR where programmer has saved. Suppose a data from sensor (peripheral1) is to be read, processed and then a motor (peripheral2) is controlled based on calculated data. The flow of actions can be summarized as:

- Peripheral1 has data in its register; meanwhile the processor is executing its main program.
- Peripheral1 asserts INT to request service from the processor.
- After execution of each instruction, the processor checks INT pin. So processor detects the service requirement. It saves its present context and sets the PC to the fixed ISR location.
- The ISR is executed which reads data from peripheral1, modifies it and sends the resulting data to peripheral2. At the same time, peripheral1 deasserts INT after data is read from it.
- The processor retrieves its state and resumes its work.

In **vectored interrupt**, peripheral must provide the address to the processor. In this method, along with INT pin, INTA pin is also required to acknowledge that the interrupt has been detected and the peripheral can provide the address of relevant ISR using system bus. The peripheral provides the address through the data bus which is read by microprocessor.

The flow of actions can be summarized as:

- Peripheral1 has data in its register; at the same time the processor is executing its main program.
- Peripheral1 asserts INT to request service from the processor.
- After execution of each instruction, the processor checks INT pin. So processor detects the service requirement and it asserts INTA.
- Peripheral1 detects INTA and puts interrupt address vector on the data bus.

- Processor jumps to the address read from data bus and executes its corresponding ISR. It reads data from peripheral1 processes it and sends the result to the peripheral2. Meanwhile peripheral1 deasserts INT after data is read from it.
- The processor retrieves its state and resumes its work.

In **interrupt address table**, which is a compromise between fixed and vectored interrupt methods, a table with ISR addresses is stored in memory of processor. A peripheral instead provides the number, rather than the address of ISR, corresponding to an entry in the table. One major advantage is that the bit requirement to address the table is very less compared to number of bits of real address of ISR. Also it provides the flexibility to assign and change the location of ISR.

Additional Interrupt Issues

External interrupts may be **maskable** or **nonmaskable**. In maskable interrupt, the programmer can use specific instruction to disable the interrupt by configuring certain bits of interrupt register. It is important when more critical works need to be executed first. Nonmaskable interrupt cannot be disabled by the programmer. It requires a separate pin for drastic situations. For instance, if power fails, the nonmaskable interrupt can cause a jump to a subroutine that stores critical data in non-volatile memory, before power is completely gone.

Another issue regarding the interrupt is jump to ISR in which the microprocessor either saves complete context or partial state before jumping to ISR. Some processors save PC, registers which consumes many cycles, while others save the content of PC only. The ISR, however, must not modify registers if its content is not saved.

C. Direct Memory Access – DMA controller

INTRODUCTION

When the communication between memory and peripherals involves microprocessor then there will, somehow, always be waste of processor's time. Since the speed of the processor and peripherals may not match, data must be stored temporarily before processing which is referred as buffering. Buffering will, certainly, impact on system performance. Also while using interrupt feature, the storing and restoring of state of processor is an inefficient process, since this process requires many clock cycles. And, the regular program stalls during transfer of data causing more problems in the performance of the system. So, a separate single-purpose processor called a DMA

controller is required which relieves processor from all data transfers involving memory and peripherals.

DMA controller is specifically used to transfer data between memories and peripherals. The peripherals request the service from DMA controller which then requests control of the system bus from processor. After that, processor relinquishes the system bus. Finally, the data transfer between memory and peripheral is initiated by DMA controller without the involvement of processor. Hence, the overhead required for storing and restoring the state is eliminated. Also, the processor can continue its regular task unless it requires the system bus or the particular data being transferred.

BLOCK DIAGRAM

The simple block diagram of system involving DMA controller is shown in the figure below:

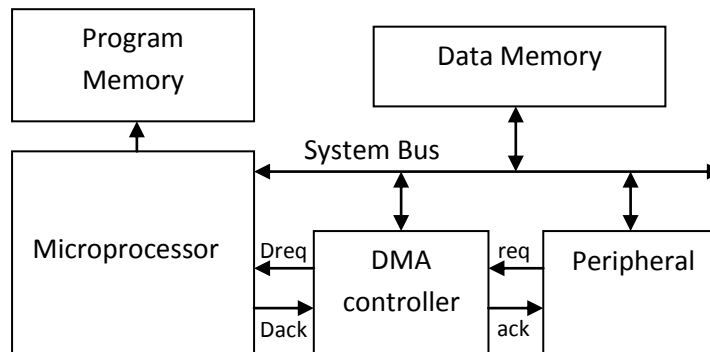


Figure 5.13: Simple system with DMA Controller

OPERATION

The flow of action for the transfer of data between peripheral and memory using DMA can be summarized as:

- Initially, processor is busy executing its main program.
- After peripheral has data within its register it asserts request line for service from DMA.
- DMA asserts request signal to request the system bus from processor.
- Processor releases the system bus after seeing the request from DMA, and acknowledges about it to DMA.
- DMA asserts acknowledge signal to peripheral, and starts transfer of data as requested.
- After the completion of transfer, all control lines are deasserted and processor retakes the control of the system bus.

5.3 Arbitration

Arbitration is the mechanism through which a service or shared resource is provided to particular requesting device, out of many contenting devices for service.

A. Priority Arbiter

INTRODUCTION

Priority arbiter is a single purpose processor which is used to arbitrate among various requests from peripherals. Each of the peripherals, which are connected to the arbiter, can make request for the service. Using certain priority mechanism, arbiter selects a peripheral to permit the required service. The figure 5.14 shows the priority arbiter connected with peripherals which use vectored interrupt to request service and processor which provide service to the peripherals. Arbiter is connected to system bus for configurations only. The configurations may include setting priorities of the peripherals.

The main advantage of this arbitration is that it can support advanced priority schemes. Also, failure of single peripheral does not have any impact on the operation of whole system. The system, however, must be redesigned if new peripherals are to be added. So, this method is less flexible if new peripherals are required to be added or removed.

BLOCK DIAGRAM

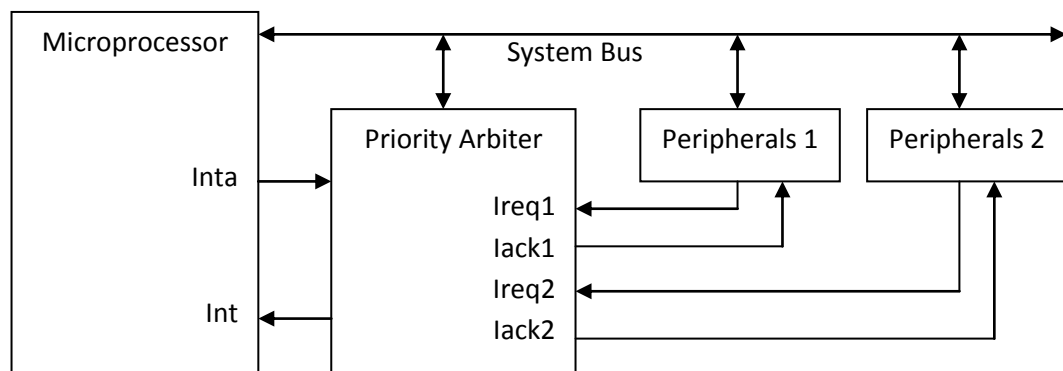


Figure 5.14: Arbitration using a priority arbiter

OPERATION

The stepwise operation of arbitration using priority arbiter is listed below:

- Initially, microprocessor is busy in its own operation.

- Both peripherals can assert request to priority arbiter which interrupts processor when at least one request is available from peripherals.
- Processor stops its current operation, stores its state and asserts interrupt acknowledgement signal.
- After acknowledged by processor, priority arbiter asserts acknowledge signal to any one peripheral based on priority.
- The selected peripheral puts its interrupt address vector on the system bus.
- Microprocessor reads ISR address from data bus and jumps into its, executes the ISR.
- After execution of requested ISR, processor retrieves its state and resumes its operation.

TYPES OF PRIORITY ARBITER

The priority among peripherals can be determined based on, basically, two schemes; fixed priority or rotating priority.

Fixed Priority

Each peripheral is assigned a unique rank. If two peripherals simultaneously request for service then the arbiter chooses the one with the higher rank. Such method is efficient when there is a clear distinction in priority among peripherals. But it can cause high-ranked peripherals to get much more servicing than other peripherals.

Rotating priority or Round – robin priority

In this method, each peripheral gets almost equal time for service from the arbiter. This priority method is efficient when there is not much difference in priority among peripherals. The priority of peripherals changes based on the history of servicing of those peripherals, so the arbiter can get more complex in rotating priority.

B. Daisy-Chain Arbitration

INTRODUCTION

In daisy-chain arbitration, peripherals are connected to each other in daisy-chain manner. The arbitration is build within the peripherals with each having a request and acknowledge signals as shown in the figure 5.15. The request signal and acknowledge signals flow through the peripherals: peripheral's request signal flows downstream to processor and processor's acknowledge signal flows upstream to requesting peripheral. The peripheral connected first to the processor has the highest priority while the peripheral at the end of chain has lowest priority.

The main advantage of this arbitration method is that one can easily add or remove peripherals from the system without the requirement of system redesign. This method, however, does not support rotating priority. Also, if one peripheral is damaged in the chain, other peripherals beyond that broken point will remain inaccessible as signal cannot pass through the chain.

BLOCK DIAGRAM

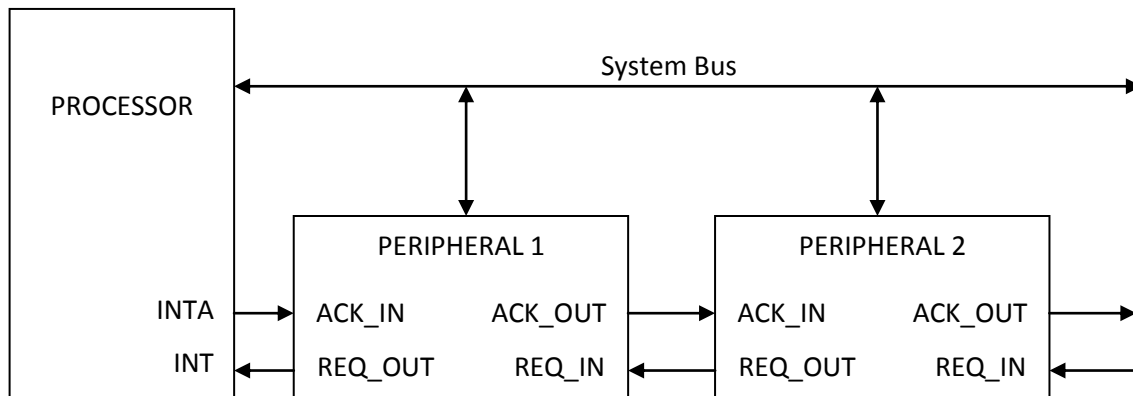


Figure 5.15: Daisy Chain Configuration

OPERATION

Suppose peripheral 2 requires service from the processor then the operation can be summarized as:

- Microprocessor is busy in executing its own program.
- The request signal from peripheral 2 is send to processor through the peripheral 1 and interrupt pin is asserted.
- Processor stops its current work, stores its state, and asserts acknowledgement signal.
- The acknowledgement signal reaches to peripheral 2 through peripheral 1. Since the request is not generated by peripheral 1, it passes the acknowledge signal to peripheral 2.
- Peripheral 2 puts its interrupt address vector on the system bus.
- Microprocessor reads ISR address from data bus and jumps into its, executes the ISR.
- After execution of requested ISR, processor retrieves its state and resumes its operation.

Daisy Chain aware peripherals

Generally, peripherals have acknowledge input and request out lines but daisy chain aware peripherals must have additional acknowledge output and request input lines. However, if the peripherals do not contain acknowledge output and request input lines then they will not be daisy chain aware peripherals. But they can be made daisy chain aware by certain logic whose complexity

may increase based on complexity of system. One simple example for making a peripheral daisy chain aware is shown in the figure below.

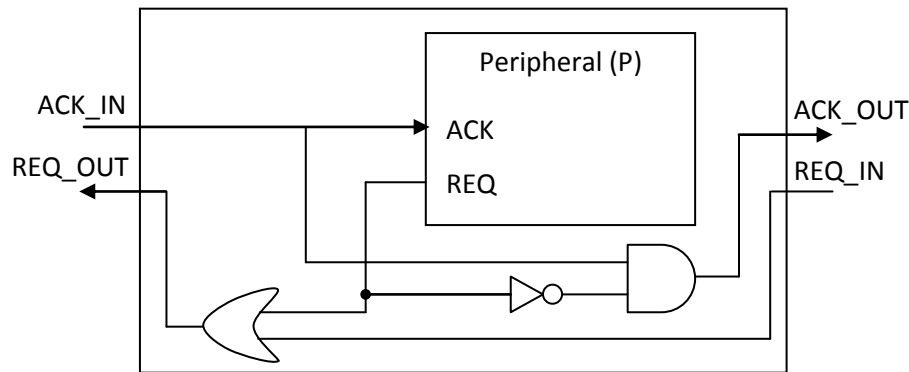


Figure 5.16: Simple Logic to make Daisy Chain Aware

Case 1: When request is from downstream peripherals

- Peripheral (P) does not participate in the flow of signal

Case 2: When request is from upstream peripherals beyond peripheral (P)

- REQ_IN = 1 but REQ = 0, resulting in REQ_OUT = 1
- ACK_IN = 1 and REQ = 0, resulting in ACK_OUT = 1

Case 3: When request is from peripheral (P)

- REQ = 1, REQ_IN = X (don't care), resulting in REQ_OUT = 1
- ACK_IN = 1 and REQ = 1 resulting in ACK = 1 and ACK_OUT = 0

C. Network-Oriented Arbitration

In network oriented arbitration, arbitration is done for multiple microprocessors sharing a common to form a network. Arbitration is build into the bus protocol, as bus is the only the medium that connects multiple processors. However, multiple processors may try to access the bus simultaneously resulting in data collision. The protocol must be designed in such a way that the contending processors don't start sending the data at the same time. Also some statistical methods can be used so as to make chances of data collision very rare, if not eliminate it. Some protocols use efficient address encoding schemes in which higher priority address will override the lower-priority one.

5.4 Multilevel Bus Architectures

Multilevel bus architectures are implemented in the system to improve the overall performance of the system. One can easily presume a single high-speed bus would be enough for all the communications in the system. But, however, there are various drawbacks of using single high speed bus. Few of them are discussed in the following paragraph.

Inefficient interface

For a single high speed bus, each peripheral requires a high-speed bus interface. But the peripherals may not need such high-speed transfer resulting in extra power consumption, increase in number of gates, and high cost. Also, the high-speed bus can be very processor specific which can lead the interface of a peripheral to be non portable.

Slower bus

When many peripherals are connected to a single bus, all peripherals may not get the access to the bus when required. This condition results to slow down the speed of transfer, hence it can create a performance lag.

Two Level Bus Systems

Generally, two level bus systems consist of a high-speed processor local bus, a lower-speed peripheral bus and a bridge to connect two buses.

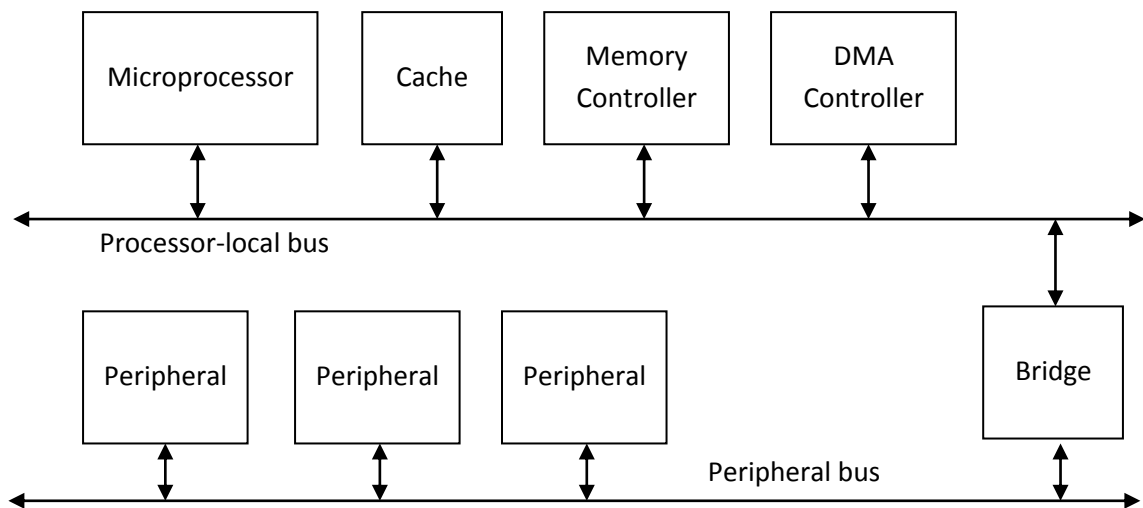


Figure 5.17: A two-level bus architecture

- The **processor local bus** connects very high speed devices such as microprocessor, cache, memory controllers, and certain high-speed coprocessors. These buses are wide, as wide as a memory word and frequent communication takes place through it.

- On the other hand, the **peripheral bus** connects to those peripherals which do not have access to processor-local bus. It emphasizes on portability, low power, or low gate count. It is often narrower and slower than a processor local bus. The frequency of communication through peripheral bus is also less as compared to that of processor local bus. So the interface for peripheral bus is comparatively efficient one in terms of number of pins, gates and power consumption.
- A **bridge** is a single purpose processor that connects the two buses of the system and also makes the various conversions required. Speed synchronization is another important function of bridge. Data speed and data formats of processor-local bus is different to that of peripheral bus, such problem is resolved by bridge using various mechanisms.

Three level bus hierarchy

Three level bus systems consist of processor local bus, system bus and peripheral bus. A local bus connects the processor to a cache and may support one or more local devices. The system bus, acting as high-speed bus, offloads much of the traffic from the processor local bus. And the peripheral bus is used to connect various peripherals in the system.

5.5 Advanced Communication Principles

Parallel Communication

In parallel communication, the physical layer carries multiple bits of data at a time. With each wire carrying a single bit, the bus consists of data wires along with control and power lines.

Advantages

- High data throughput: Many bits are transferred at a time.
- Less complexity: Easily implemented in hardware requiring only a latch to copy data onto a data bus.

Disadvantages

- Long parallel wires can result in Ferranti Effect. And according to this effect, there is a voltage build up due to capacitance and voltage at receiving end becomes more than that of sending end.
- Little variation in wire length can cause data misalignment as the bits at the receiving ends with reach at different time.
- It is more costly to construct and can make system bulky. The cost further increases if insulation of wires to prevent the interference is considered.

Usage

- It is used to connect devices which reside on the same circuit board or same IC.

Serial Communication

In serial communication, the physical layer carries one bit of data at a time. With all bits of data passing through the single wire, the bus is composed of single data wire along with control and power lines.

Advantages

- Significant reduction in the size, the complexity of the connectors and the associated costs.
- Throughput can be better for two distant devices as compared to parallel communication of two distant devices.
- It does not exhibit Ferranti effect and data misalignments.

Disadvantages

- Complex interfacing logic and communication protocols; the data are decomposed into bits at sending end, which must be assembled properly at receiving end.
- For short distance communication, its throughput is very less as compared to that of parallel communication.

Usage

- It is used to connect distant devices. But it doesn't mean that it cannot be used to connect devices at short distance. However, it is more efficient for distant communications

Wireless Communication

In wireless communication, the devices do not need to be connected physically for data transfer. Infrared and radio frequency channels are used as a physical layer.

Infrared wave

Infrared wave, which cannot be seen by naked eye, uses electromagnetic wave frequencies that are below the visible light spectrum. Infrared waves are generated using infrared diode whereas infrared transistors are used to detect the infrared emitted by infrared diode. Such infrared transistors conduct when exposed to infrared wave. One advantage of infrared communication is that it is cheap to build transmitters and receivers. But the main disadvantage of this sort of communication is that it requires line of sight between the two devices participating in

communication. Also the range of communication is low, which makes it an inefficient method of communication for distant devices.

Radio Frequency

Radio frequency uses electromagnetic wave frequencies in the radio spectrum. For such communication, analog circuitry as well as an antenna is required at communicating devices. The main advantage of this type of communication is that the line of sight is not required. Also the longer distance communication is possible. The range of communication is dependent on the transmission power. But building transmitters and receivers can be complex and costly in Radio Frequency communication.

Layering

Layering the communications process means breaking down the communication process into smaller and easier to handle interdependent categories, with each solving an important and somehow distinct aspect of data exchange process. Layering can also be viewed as a hierarchical organization of a communication protocol where lower levels of the protocol provide services to the higher levels. The main objective is to break the complexity of a communication protocol into simple levels which ensures easier handling and simplified design. The physical layer provides the lower level services of sending and receiving bits or words of data, whereas the application layer provides the high level service to the user.

Error Detection and Correction

Error detection is the process of detecting errors that may occur during the transmission of data in any communication process. Errors can be bit error or burst of bit errors. In bit error, single bit in the transmitted data is invalid. But in case of burst of bit errors, more than one bit gets changed.

Error correction is the process of correcting the bits that were detected during communication process. Parity and checksum are two basic error correction methods.

In **parity** check, extra bit is send along with data to provide additional information about the data. If extra bit makes an odd number of 1s in data word bits plus parity bit then it is referred as odd parity otherwise it will be even parity. The parity of data must be check at both ends of communication and the parity of the data sent must be same to that of parity of data received. This type of checking

method is efficient for single bit error but can create problems for burst of bit errors as it is not able to detect change in even number of bits.

Example:

Data of 7 bits – 0011010

Transmitted data with even parity – 00110101

Received data with parity – 10110101, it indicates error as it should have an even parity.

Received data with parity – 10010101, change in two bits and error not detected.

In **checksum** error checking, multiple words of data in packets are checked for error. The extra word which represents the XOR sum of all data words in a packet is transmitted along with packet. Though it can be implemented for burst of bits error, it does not account for all error combinations. The transmitter sends the packet of data along with the checksum word which is checked at receiving end on reception. If the checksum word is correct then it represents successful transmission. However, few error combinations can generate the checksum word same as received, in such case the error checking fails.

Example:

Data words of packet to be transmitted: 010101, 011101, 110011, 101100

Checksum word at transmitter: 010111 (XOR of all data words)

Received checksum word: 010111

Received data words of packet: 110101, 010101, 110011, 101100, error exists

Calculate Checksum word at receiver: 111111, checksum does not match, error check success.

Received data words of packet: 010101, 011101, 110011, 101100, error exists

Calculate Checksum word at receiver: 010111, checksum match, error check fails.

5.6 Serial, Parallel and Wireless Protocols

A. Serial Protocol

Inter-IC or I2C or I²C

I2C is a serial protocol for two-wire interface to connect low-speed devices like microcontrollers, EEPROMs, A/D and D/A converters, I/O interfaces and other similar peripherals in embedded systems. The I2C has 7-bit or 10 bit address space. Seven bit addressing allows a total of 128 devices to communication over a shared I2C bus. The common speed of I2C bus is 100 Kbit/s in standard

mode and 10 Kbit/s in low-speed mode. Recent revisions of I2C can host more nodes and run at faster speeds: 400 Kbit/s in Fast mode, 1 Mbit/s in Fast mode plus and 3.4 Mbit/s in High speed mode. I2C uses only two wires: SCL (serial clock) and SDA (serial data). Length of the wires is not limited as long as the total bus capacitance is less than 400pf.

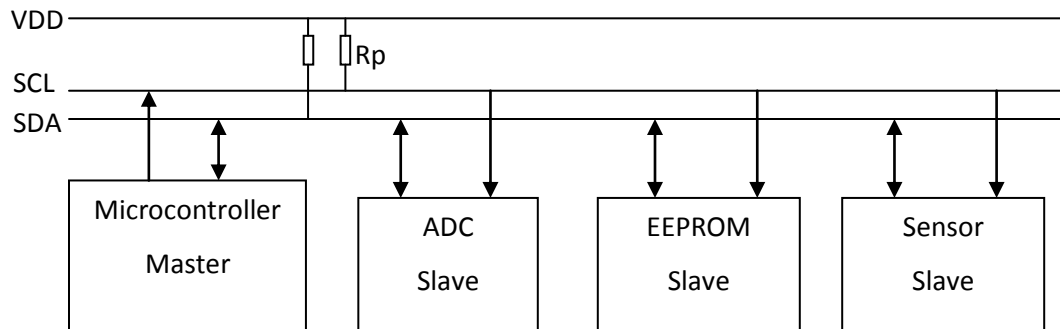


Figure 5.18: I2C bus structure

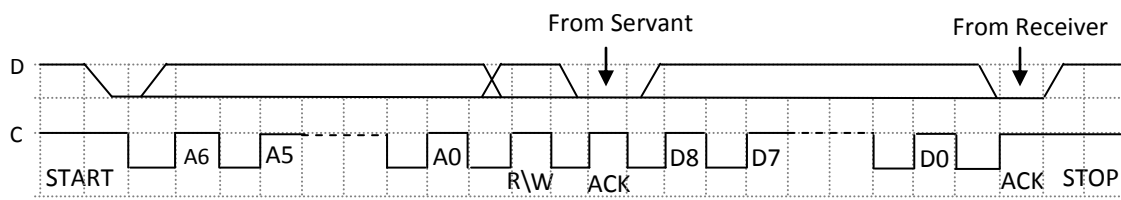


Figure 5.19: Timing diagram of a typical read/write cycle

A typical I2C byte write cycle operates as follows:

- The master initiates the transfer with a start condition. Start condition is represented by a high to low transition of SDA line while the SCL is held high.
- Then, the address of the device to which the data is to be written is sent with most significant bit down to the least significant bit.
- For write operation, the master sends a zero after sending the address. And the slave acknowledges the transmission by holding the SDA line low during first ACK clock cycle.
- Next, the master transmits a byte of data with most significant bit first.
- The slave acknowledges the reception of data by holding the SDA line low during second ACK clock cycle.
- Finally, master terminates the transfer by generating a stop condition. Stop condition is represented by a low to high transition of SDA line while the SCL is held high.

Serial Peripheral Interface (SPI)

The serial peripheral Interface bus is a synchronous serial communication interface specification used for short distance communication. It is used to send data between processors/controllers and small peripherals. It uses separate clock and data lines along with a select line to choose the device that should be communicated with. The SPI bus consists of four logic signals: SCLK – serial clock, MOSI – Master Output Slave Input, MISO – Master Input Slave Output, and SS – Slave select. The SPI bus can operate with a single master device and with one or more slave devices. Full duplex communication, higher throughput, simple software and hardware implementation, etc are few characteristics of SPI protocol.

Control Area Network (CAN)

CAN is an International Standardization Organization (ISO) defined serial communications bus originally developed for the automotive industry to replace the complex wiring harness with a two wire bus. The specification calls for high immunity to electrical interference and the ability to self-diagnose and repair data errors. These features have led to CAN's popularity in a variety of industries including building automation, medical, and manufacturing.

Some of the characteristics of the CAN protocol includes high-integrity serial data communications, real-time support, data rates up to 1 Mbit/s, error detection and confinement capabilities. Balanced differential signaling in CAN protocol not only reduces noise coupling but also allows high signaling rates over twisted pair cable. The CAN protocol incorporates five methods of error checking which forces transmitting node to resent the message until it is received correctly. But if the error limit is reached then the faulty node is deprived of transmit capability. Faulty nodes are automatically dropped from the bus, which prevents any single node from bringing a network down. This error containment also allows nodes to be added to a bus while the system is in operation, otherwise known as hot-plugging. It implements a non destructive, bit-wise arbitration in which the node winning arbitration continues with the message without being corrupted by another node. The high speed ISO 11898 standard specifications are given for a maximum signaling rate of 1 Mbps with a bus length of 40m with a maximum of 30 nodes.

To summarize, the protocol defines data packet format and transmission rules to prioritize messages, guarantee latency times, allow for multiple masters, handles transmission errors,

retransmit corrupted messages, and distinguish between a permanent failures of a node versus temporary errors

FireWire

FireWire is a serial bus protocol for high-speed data transfer. It was initiated by Apple and developed by IEEE P1394 group so may refer this protocol as IEEE 1394. It supports mass information transfer and allows peer to peer device communication without the involvement of system memory or CPU. Some of the characteristics of FireWire protocol include transfer rates up to 400 Mbit/s, plug and play and hot swapping, packet based layered design structure and provision of power through the cable. Also, the 64 bit addressing allows a local-area network to consist of 1023 sub-networks, each consisting of 63 nodes. FireWire devices are organized at the bus in a tree or daisy chain topology. In arbitration, the closest node requesting for the data transfer gets the high priority. It provides two types of data transfer: asynchronous and isochronous. In asynchronous, data transfer can be initiated as a given length of data arrives in a buffer. But, in isochronous data transfer, data flows at a pre-set rate.

Universal Serial Bus (USB)

The Universal Serial Bus (USB) protocol was designed to connect a wide range of peripherals to a computer, including pointing devices, displays, data storage, communication devices and other devices. It standardized the connection of computer peripherals to personal computers, both to communicate and to supply electric power. The original USB 1.0 specification defined data transfer rates of 1.5 Mbit/s for low data rate devices and 12 Mbit/s for high speed devices. The transfer rate for USB 2.0 is 480 Mbit/s while for USB 3.0 it can go up to 5 Gbit/s. USB On-The-Go is the special feature of USB in which two USB devices communicate with each other without requiring a separate USB host. USB uses a tiered star topology, which means some USB devices can serve as connection ports for other USB peripherals. USB hubs and standalone hubs can be used to provide handful of convenient USB ports. USB host controllers manage and control the driver software and bandwidth required by each peripheral connected to the bus.

B. Parallel Protocols

PCI Bus

The Peripheral Component Interconnect (PCI) bus is a high-performance bus for attaching hardware devices in a computer. It is synchronous bus architecture with all data transfers being performed relative to a system clock. The maximum clock rate can go up to 66MHz; however, use of 33MHz is

very common in personal computers. So the transfer rate can vary from 132 to 512 MB/s. PCI implements a 32-bit multiplexed address and Data bus which allows reduced pin count on the PCI connector resulting in lower cost and smaller package size. It supports rigorous auto configuration mechanisms which allow identification of the type of device and the company that produced it. In PCI, any device has the potential to take control of the bus and initiate transactions with any other device making multiple master implementations easier which otherwise had been difficult.

ARM Bus

ARM bus was designed to connect and manage different function blocks in a system on a chip (SoC) designs. It supports 32-bit data transfer and 32-bit addressing and is implemented using synchronous data transfer architecture. The transfer rate is the function of the clock speed used in a particular application. The ARM Advanced Microcontroller Bus Architecture is an open-standard for on chip interconnection.

C. Wireless Protocols

Infrared Data Association (IrDA)

The infrared Data Association (IrDA) is an international organization that creates and promotes infrared data interconnection standards. It provides specifications for a complete set of protocols for wireless infrared communications. IrDA has been implemented in portable devices like smart phones, laptops, cameras, etc. It is designed to support communication between two devices over point to point infrared at speeds between 9.6 kbps and 4 Mbps. Simplicity and low cost of IrDA hardware makes it an attractive option. Also, line of sight, very low bit error rate and physically secure data transfer are few important features of IrDA. Other wireless technologies with no requirement of direct line of sight have displaced IrDA. However, it is still applicable where interference makes radio based wireless technologies unusable.

Bluetooth

Bluetooth is a wireless technology standard for exchanging data over short distances from fixed and mobile devices. It operates at frequencies between 2402 and 2480 MHz which is the globally unlicensed Industrial, Scientific and Medical (ISM) 2.4 GHz short-range frequent band. Since Bluetooth uses a radio-based link, it does not require line of sight for communication. Bluetooth 4.0 may provide the transfer rate of up to 25Mbps. Bluetooth is a packet based protocol with a master-slave structure and one master may communicates up to maximum of seven slave devices. Low power consumption and short range based communication is the typical feature of Bluetooth.

Permitted transmission power and range of communication depend on the radios class. For class 3 radio, range is up to 1m with max permitted power of about 1mW. The range is about 10m and 2.5mW power is permitted in case of class 2 radios, and class 1 radios have a range of about 100m and 100mW of transmission power. Handsfree headset and wireless speakers are two, out of many, examples using Bluetooth.

IEEE 802.11

IEEE 802.11 is a set of media access control (MAC) and physical layer (PHY) specifications for implementing wireless local area network. IEEE 802.11, often termed as Wi-Fi, has the data transfer rate of around 1 - 2Mbps. IEEE 802.11 has a variety of standards, each with a letter suffix; 802.11a, 802.11b, 802.11g, 802.11n standards are quite common. All these 802.11 Wi-Fi standards operate within the ISM frequency bands. Generally, 2.4 GHz band is common which also makes the chips easier and cheaper to manufacture. The data rate can go up to 54 Mbps with some standard, while few latest standards may support up to 6.75 Gbit/s.

The PHY layer defines the means of transmitting bits over a physical link connecting network nodes. It provides an electrical, mechanical and procedural interface to the transmission medium. Modulation, line coding, synchronization are few functions, out of many, performed by the physical layer. The MAC layer provides the addressing and channel access control mechanism that ensures the communication of several nodes within a shared medium. Each device is assigned a unique serial number which is also known as MAC address. Unique MAC address makes it possible for data packets to be delivered to a destination within a sub-network. Multiple access protocol allows several stations connected to the same physical medium to share it. The most common multiple access protocol is the contention based Carrier Sense Multiple Access with Collision Avoidance (CSMA/CD) protocol.

- **Operating System basics**
- **Task Process and Threads**
- **Multiprocessing and Multitasking**
- **Task Scheduling**
- **Task Synchronization**
- **Device Drivers**

6.1 Operating System basics

The operating system acts as a bridge between the user applications/tasks and the underlying system resources through a set of system functionalities and services. The primary function of an operating system is

- Make the system convenient to use
- Organize and manage the system resources efficiently and correctly

The following figure shows the basic components of an operating system and their interfaces with rest of the world.

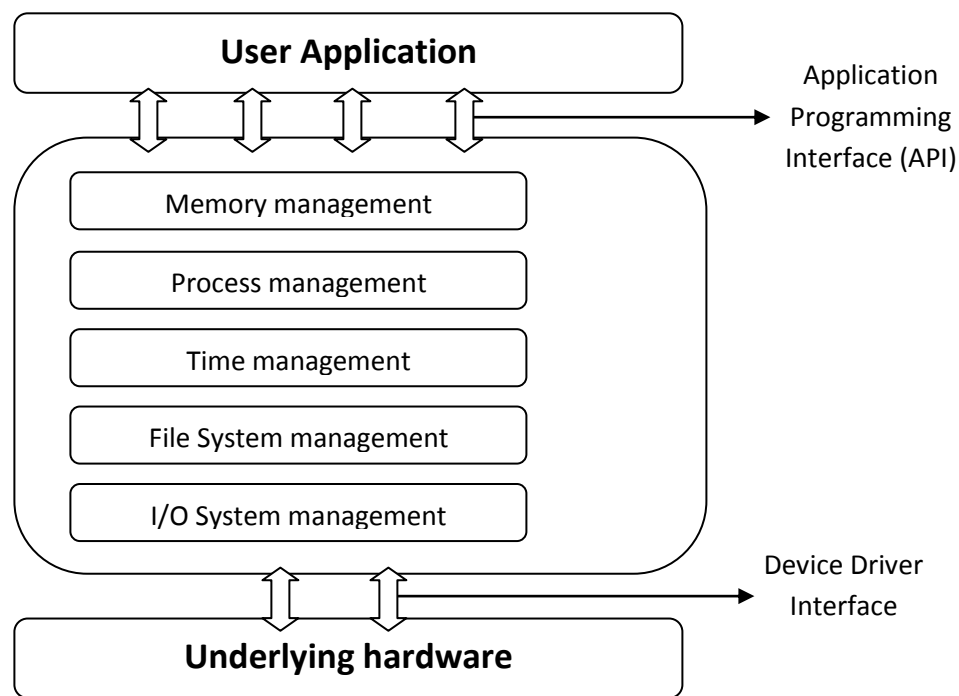


Figure 6.1: Operating System Architecture

Comparison of General Purpose OS (GPOS) with Real Time OS (RTOS)

General Purpose Operating System is software that manages all the system resources and provides common service to all programs running in the system. In case of real time operating system, along with management and services it performs certain function within a specified time constraint. However, both operating systems provide a number of services to application programs and users. Application Programming Interfaces (API) or system calls are the medium through which the services are accessed by the applications.

The differences between GPOS and RTOS can be clarified using following parameters.

- **Deterministic nature**

RTOS are deterministic in nature; the time required to execute the services is fixed. However, there may not be fixed time defined for any service in case of GPOS.

- **Task Scheduling**

RTOS uses priority based preemptive scheduling, while scheduling in GPOS is defined so as to achieve high throughput. In RTOS, high priority process execution will override the low priority ones. In GPOS, high priority process may be delayed to perform several low priority tasks.

- **Time Critical systems**

RTOS is used in time critical systems in which delay in processing can result in undesirable consequences. However, GPOS are implemented in non time critical systems.

- **Preemptive Kernel**

The kernel of an RTOS is preemptive where as a GPOS kernel is non preemptive. In preemptive kernel, the high priority user process can preempt a kernel call. In other words, the execution of low priority system process can be stopped by high priority user process.

- **Priority Inversion Problem**

Priority Inversion problem is seen in RTOS in which the high priority task has to wait for the shared resource occupied by low priority task. This results in execution of low priority task first rather than high priority task.

A. The Kernel

The kernel is the core of the operating system and is responsible for managing the system resources and the communication among the hardware and other system services. It acts as the abstraction layer between system resources and user applications. The kernel contains different services for handling the following.

Process Management

It includes setting up the memory space for the process, loading the process's code into the memory space, allocating system resources, scheduling and managing the execution of the process, setting up and managing the process control block (PCB), Inter Process Communication and Synchronization, process termination/deletion, etc.

Primary Memory Management

The term primary memory refers to the volatile memory (RAM) where processes are loaded and variables and shared data associated with each process are stored. The Memory Management Unit (MMU) of the kernel is responsible for

- Keeping track of which part of the memory area is currently used by which process
- Allocating and De- allocating memory space on a need basis (Dynamic memory allocation)

File System Management

File is a collection of related information. A file could be a program, text files, word documents, audio/video files, etc. Each of these files differs in the kind of information they hold and the way in which the information is stored. The file operation is a useful service provided by the OS. The file system management service of Kernel is responsible for

- The creation, deletion and alteration of files and directories
- Saving of files in the secondary storage memory
- Providing automatic allocation of file space based on the amount of free space available
- Providing flexible naming convention for the files

I/O System (Device) Management

Kernel is responsible for routing the I/O requests coming from different user applications to the appropriate I/O devices of the system. In a well-structured OS, the direct accessing of I/O devices are not allowed and the access to them are provided through a set of Application Programming Interfaces (APIs) exposed by the kernel. The kernel maintains a list of all I/O devices of the system. The list may be available in advance and recent kernel dynamically updates the list of available devices. The service 'Device Manager' of the kernel is responsible for handling all I/O device related operations. The kernel talks to the I/O device through a set of low level system calls, which are implemented in a service, called device drivers. The device manager is responsible for

- Loading and unloading of device drivers

- Exchanging information and the system specific control signals to and from the device

Secondary Storage Management

The secondary storage management deals with managing the secondary storage memory devices that are connected to the system. Secondary memory is used as backup medium for programs and data since the main memory is volatile. In most systems, the secondary storage is kept in disks (Hard Disk). The secondary storage management service of kernel deals with

- Disk storage allocation
- Disk scheduling (time interval at which the disk is activated to backup data)
- Free Disk space management

Protection Systems

Modern operating systems are designed in such a way to support multiple users with different levels of access permissions (For example: Administrator, Standard, Restricted, Guest, etc). Implementing security policies to restrict the access to both user and system resources by different applications or processes or users, one user may not be allowed to view or modify the whole/portions of another user's data or profile details. Some application may not be granted with permission to make use of some of the system resources.

Interrupt Handler

Kernel provides a mechanism to handle all external/internal interrupts generated by the system. Based upon the priority of the interrupt the process either runs in the foreground or background. Depending on the type of operating system, a kernel may contain lesser number of services or more number of services which may include network communication, network management, user-interface graphics, timer services (delays, timeouts, etc.), error handler, database management, etc.

B. Kernel Space and User Space

The program code corresponding to the kernel applications/services are kept in a contiguous area of primary memory and are protected from un-authorized access by user programs/applications. The memory space at which the kernel code is located is known as 'Kernel Space'. Similarly, all user applications are loaded to a specific area of primary memory and this memory area is referred as 'User Space'. User space is the memory area where user applications are loaded and executed. The partitioning of memory into kernel and user space is purely OS dependent.

C. Types of Kernel

Based on the kernel architecture/design, kernels can be classified into Monolithic and Micro.

Monolithic Kernel

In Monolithic Kernel architecture, all kernel services run in the kernel space. Here all kernel modules run within the same memory space under a single kernel thread. It runs all basic system services and provides powerful abstraction of the underlying hardware. Amount of context switches and messaging involved are greatly reduced which makes it run faster than microkernel. The major drawback of monolithic kernel is that any error or failure in any one of the kernel modules leads to the crashing of the entire kernel application. The inclusion of all basic services in kernel space leads to different drawbacks such as requirement of large kernel size, lacking extensibility, poor maintainability. LINUX, SOLARIS, MS-DOS kernels are examples of monolithic kernel.

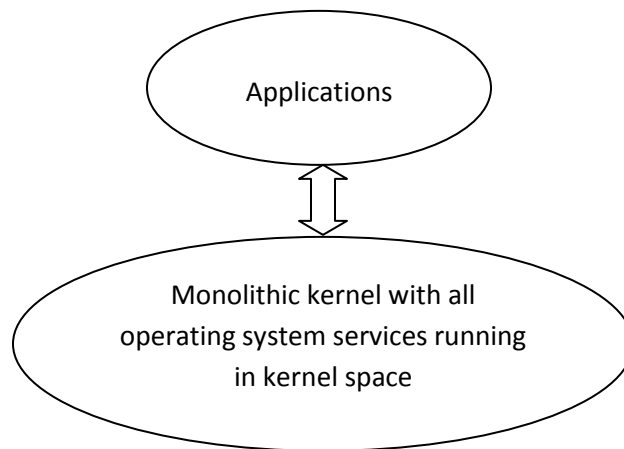


Figure 6.2: The Monolithic Kernel Model

Microkernel

The microkernel design incorporates only the essential set of operating system services such as communication and I/O control into the kernel. The rest of the operating system services are implemented in programs known as 'Servers' which runs in user space. It is more stable than monolithic as the kernel is unaffected even if the server fails. Memory Management, process Management, timer systems and interrupt handlers are the essential services, which forms the part of microkernel. Microkernel based design approach offers the following benefits.

- **Robustness:** If a problem is encountered in any of the service, which runs as 'Server' application, the same can be reconfigured and re-stated without the need for re-starting the entire OS.

- **Configurability:** services can be changed, updated without corrupting the essential services residing within the microkernel.

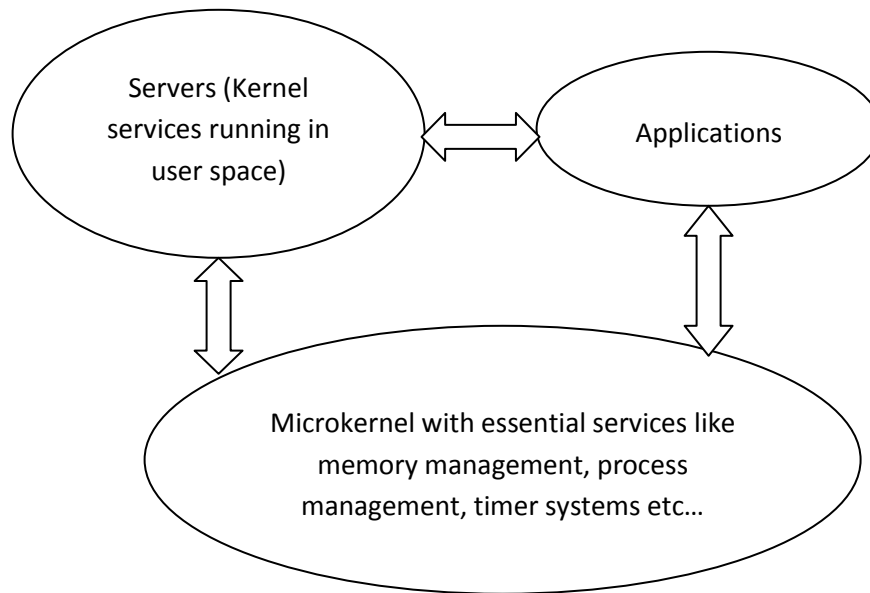


Figure 6.3: The Microkernel Model

6.2 Task Process and Threads

A task is defined as a program in execution and related information maintained by OS for that program. Task is also known as 'Job' in the operating system context. A program or part of it in execution is also called a 'Process'. The terms 'Task', 'Job' and 'Process' refer to the same entity in the operating system context and most often they are used interchangeably.

A. Process

A process is an instance of a program or part of program in execution. A process requires various system resources such as the CPU for executing the process, memory for storing the code corresponding to the process and associated variables, I/O devices for information exchange etc. A program by itself is not a process; a program is a passive entity, such as a file containing a list of instructions stored on the disk (executable file). A process is an active entity. A program becomes a process when an executable file is loaded into memory.

Structure of a process

A process holds a set of registers, process status, a Program Counter (PC) to point to the next executable instruction of the process, a stack for holding the local variables associated with the

process and the code corresponding to the process. From a memory perspective, the memory occupied by the process is separated into three regions, stack memory, data memory and code memory.

The stack memory holds all temporary data such as variables local to the process. Data memory holds all global data for the process. The code memory contains the program code (instructions) corresponding to the process.

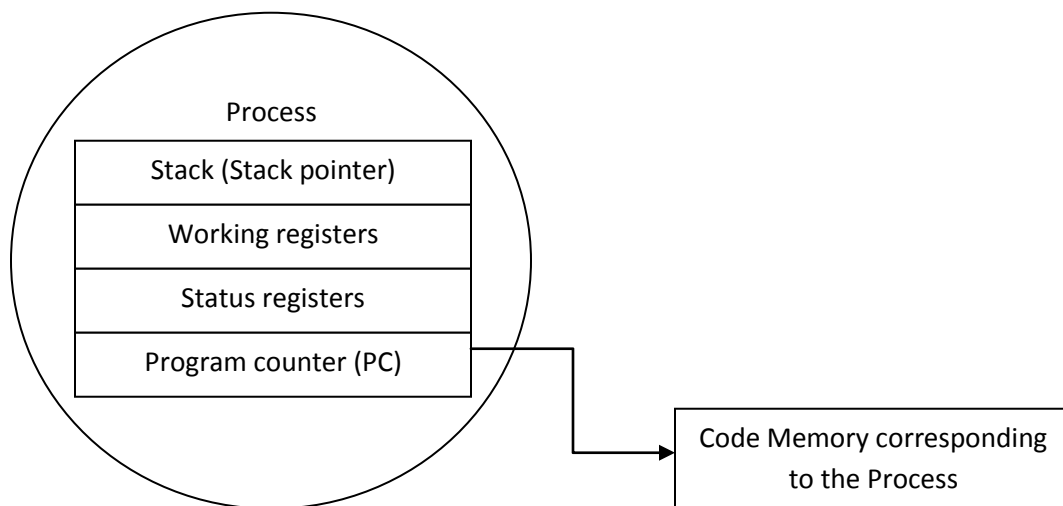


Figure 6.4: Structure of a Process

Process States and State Transition

The process traverses through a series of states during its transition from the newly created state to the terminated state. The cycle through which a process changes its state from 'newly created' to 'execution completed' is known as 'Process Life Cycle'.

- **Created State:** it is the state at which a process is being created. The operating system recognizes a process in the created state but no resources are allocated to the process.
- **Ready State:** It is the state, where a process is loaded into the memory and awaiting the processor time for execution. The process is placed in the ready list queue maintained by the OS.
- **Running State:** It is the state where the source code instructions corresponding to the process are being executed. The process execution happens in this state.

- **Blocked/Waiting State:** it refers to a state at where a running process is temporarily suspended from execution and does not have immediate access to resources. The blocked state might be invoked by various conditions like: the process enters a wait state for an event to occur or waiting for getting access to a shared resource.
- **Terminated/Completed State:** It is a state where the process completes its execution.

Different OS kernel can have different name for the state associated with a task. Created state may be stated as dormant state, waiting state may be restated as Pending state and so on.

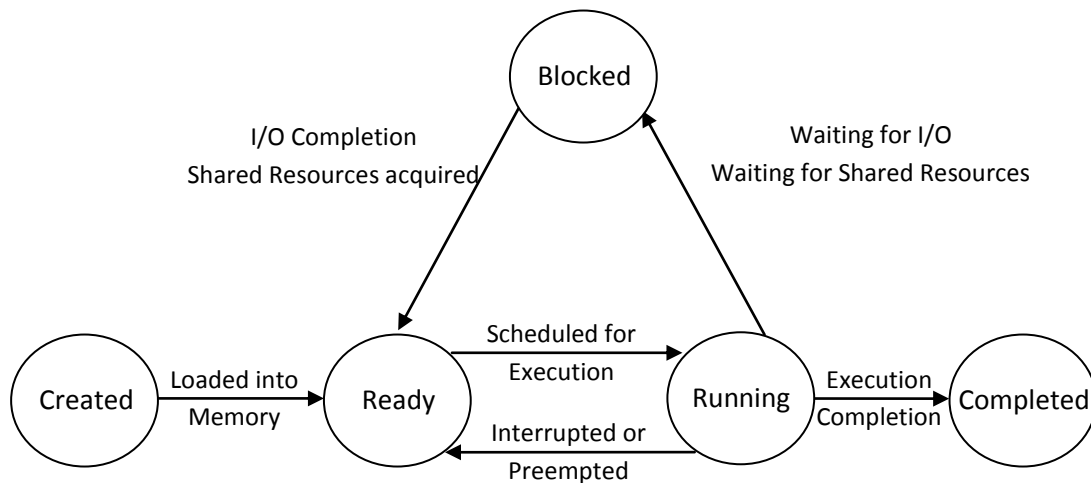


Figure 6.5: Process states and state transition representation

Process Control Block (PCB)

Each process is represented in the OS by a process control block. A PCB serves as a repository for any information that may vary from process to process. A PCB contains many pieces of information associated with a specific process.

- **Process state:** The state may be new, ready, running, waiting/blocked/pending or completed.
- **Program counter:** It indicates the address of next instruction to be executed for current process.
- **CPU registers:** They include accumulators index registers, stack pointers, general purpose registers along with any status registers. The content of PC along with the state information of a process must be saved when an interrupt occurs.

- CPU Scheduling Information: This information includes the process priority and the pointers to the scheduling queues
- Memory management Information: This information includes the value of the base registers, page tables depending upon the memory system used by the OS.
- Accounting information: This information includes the amount of CPU time, time limits and process numbers.
- I/O status information: It includes the list of I/O devices allocated to a process.

B. Threads

A thread, basic unit of CPU utilization, is a single sequential flow of control within a process. A process can have many threads of execution. Different threads which are part of process share the data memory, code memory and the heap memory.

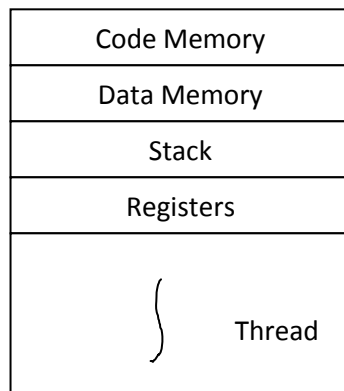


Figure 6.6: Single-Threaded Process

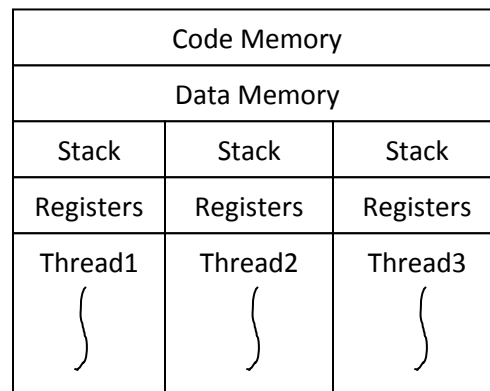


Figure 6.7: Multi-Threaded Process

However, the threads maintain their own thread status (CPU register value), Program Counter (PC) and stack. If a process has multiple threads of control, it can perform more than one task at a time. It is called a multi threaded process. If a process has a single thread of control it can perform a single task and is called single threaded process.

Concept of Multithreading

A process contain various sub-operations like getting input from I/O devices connected to the processor, performing some internal calculations/operations, updating some I/O devices etc. If all the sub-functions of a task are executed in sequence, the CPU utilization may not be efficient. For example, if the process is waiting for a user input, the CPU enters the wait state for the event, and for

the process execution also enters a wait state. If a process is split into different threads carrying out the different sub-functionalities of the process, the CPU can be effectively utilized and when the thread corresponding to the I/O operation enters the wait state, another thread which do not require the I/O event for their operation can be switched into execution. This leads to more speedy execution of the process and the efficient utilization of the processor time and resources.

The benefits of multi-threaded can be broken down into the following major categories:

- **Responsiveness:** Multi-threading on interactive application may allow a program to continue running even if part of it is blocked or is performing a lengthy operation, thereby increasing responsiveness to the user.
- **Economical:** Process creation is costly in terms of allocating memory and resources. Multiple thread creation within a process is economical because threads share the resources of the process to which they belong (code, data, heap memory). Creation of threads and context-switch of threads is economical.
- **Utilization of multiprocessor architecture:** The benefits of multithreading can be greatly increased in a multi processor architecture, where threads may be running in parallel in different processors. A single threaded process can only run on one processor, no matter how many processors are available. Multi threading on a multiprocessor machine increases concurrency.
- **Efficient CPU utilization:** CPU is engaged all the time. Since a process is split into different threads, when a thread enters a wait/block state, the CPU can be utilized by other threads of the process. This speeds up the execution of a process.

C. User level & Kernel level threads

User Level Threads

The user level threads don't have kernel/OS support and they exist only in a running process. Even if a process contains multiple user level threads, the OS treats it as a single thread. It is the responsibility of the process to schedule each thread as and when ever required. User level threads of a process are non-preemptive at the thread level from the OS perspective.

Kernel Level Threads

These are individual units of execution, which the OS treats as separate threads. The OS interrupts the execution of the currently running kernel thread and switches the execution of another kernel

thread based on the scheduling policies implemented by the OS. Kernel level threads are pre-emptive.

Relationship between User level thread and Kernel level thread

There are many ways for binding/connecting user level threads with kernel level threads.

Many to One model: Many user level threads are mapped to a single kernel thread. The kernel treats all user level threads as single thread and the execution switching among the user level threads happens when a currently executing user level thread voluntarily blocks itself or relinquishes the CPU.

One to One model: Each user level thread is bonded to a kernel/system level thread. It provides more concurrency than the many to one model by allowing another thread to run when a thread makes a blocking system call. It allows multiple threads to run in parallel on multiprocessor. Creating a user level thread requires creating a corresponding kernel level thread.

Many to many model: It multiplexes many user level threads to a smaller or equal no of kernel level threads. Developers can create as many user level threads as necessary and the corresponding kernel level threads can run in parallel on a multiprocessor. When a thread performs a blocking system call, the kernel can schedule another thread for execution.

D. Thread Libraries

A thread library provides the programmer with an API for creating and managing threads. There are two primary ways of implementing thread library.

The first approach is to provide a library entirely by the user space with no kernel/OS support. All code and data structure for library exists in user space. This means that invoking a function in the library results in a local function call in user space and not a system call.

The second approach is to implement a kernel level library supported directly by the OS. In this case, code and data structure for the library exists in the kernel space. Invoking a function in the API for the library, results in a system call to the kernel.

There are three main thread libraries that are used today.

- **POSIX threads:** POSIX stands for Portable OS Interface. The POSIX standard for defining API, for thread creation and management, is pthreads. Pthreads library defines the set of POSIX thread creation and management functions in C language. Pthread may be provided as either a user level or a kernel level library.

Thread Call	Description
pthread_create()	Creates a new thread
pthread_exit()	Terminates the calling thread
pthread_join()	Blocks the current thread and waits until the completion of the thread pointed by it.
pthread_yield()	Releases the CPU to let another thread run
pthread_attr_init()	Create and initialize a thread's attributes
pthread_attr_destroy()	Releases a thread's attributes

- **Win32 threads:** Win32 threads are supported by various flavors of the windows OS. The win32 API libraries provide a standard set of win32 thread creation and management function. Win32 thread library is a kernel level library.

Thread Call	Description
CreateThread()	Creates a new thread
SuspendThread()	Temporarily suspends thread execution
ResumeThread()	Wakes up a suspended thread
ExitThread()	It terminates a thread and allocates the thread stack resources along with other resources that were held by it.

- **Java threads:** Java threads are the threads supported by Java programming language. The java thread class 'Thread' is defined in the package 'java.lang'. The java thread API allows thread creation and management directly in the java programs. Since a java virtual machine runs on the top of host operating system, the JAVA thread API on the top of a host OS, the JAVA thread API typically implemented using a thread library available on the host system. This means that on windows system, java threads are typically implemented using the win32 API. UNIX and LINUS systems user pthreads.

Thread Call	Description
Start()	Allocates memory and initializes a new thread in JAVA
Yield()	A running thread enters the ready state
Sleep()	A thread enters the suspend state
Wait()	A thread enters a blocked state
Stop()	Terminates a thread and de-allocates resources

E. Difference between Thread and Process

Thread	Process
It is a single unit of execution and is a part of the process	A process is a program in execution and combines one or more threads
A thread shares the code, data, heap memory with other threads of the same process	A process has its own code, data and stack memory
A thread cannot live independently	A process contains at least one thread
Threads are very inexpensive to create	Processes are expensive to create. Involves many OS overhead
Context switching is inexpensive and fast	Context switching is complex and involves lot of OS overhead and is comparatively slower.
If a thread expires, its stack is reclaimed by the process.	If a process dies, the resources allocated to it are reclaimed by the OS and all the associated threads of the process also dies.

6.3 Multiprocessing and Multitasking

Multiprocessing describes the ability to execute multiple processes simultaneously. Systems which are capable of performing multiprocessing are called multiprocessor system. Multiprocessor systems possess multiple CPUs/processors and can execute multiple processes simultaneously. The ability of an OS to have multiple programs in memory, which are ready for execution, is referred as multiprogramming.

In a uniprocessor system, it is not possible to execute multiple processes simultaneously. However, it is possible for a uniprocessor system to achieve some degree of pseudo parallelism in the execution of multiple processes by switching the execution among different processes. The ability of an operating system to hold multiple processes in memory and switch the processor from executing

one process to another process is known as multitasking. Multitasking creates the illusion of multiple tasks executing in parallel. Multitasking involves 'Context switching', 'Context saving' and 'Content retrieval'.

A. Context Switching

Each task may exist in any one of the different states (running, ready, blocked, etc). During the execution of an application program, individual tasks are continuously changing from one state to another. At any point of the execution, only one task is in running mode. During the process of state change, CPU control changes from one task to another, context of the to-be-suspended task will be saved while context of the to-be-executed task will be retrieved.

The process of saving the context of a task being suspended and restoring the context of a task being resumed is called context switching.

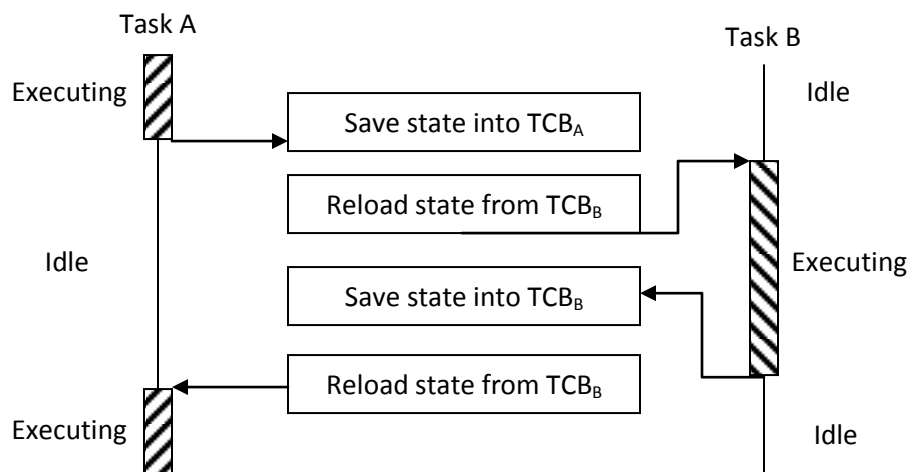


Figure 6.8: Simple Context Switching Diagram

Context Saving is the act of saving the current contents which contains the context details (register details, memory details, system resource usage details, etc for the currently running process at the time of CPU switching. Context retrieval is the process of retrieving the saved context details for a process which is going to be executed due to CPU switching. Context switch time is pure overhead because the system does no useful work while switching.

B. Types of Multitasking

Multitasking involves the switching of execution among multiple tasks. Depending on how the switching act is implemented, multitasking can be classified into different types.

Co-operative Multitasking: It is the most primitive form of multitasking in which a task/process gets a chance to execute only when the currently executing task/process voluntarily relinquishes the CPU. Any task/process can hold the CPU as much time as it wants. If the currently executing task is non-cooperative, the other tasks may have to wait for a long time to get the CPU.

Preemptive Multitasking: It ensures that every task/process gets a chance to execute. When and how much time a process gets is dependent on the implementation of the preemptive scheduling. The currently running task/process is preempted to give a chance to other tasks/process to execute. The preemption of task may be based on time slots or task/process priority.

Non-preemptive Multitasking: In non-preemptive multitasking, the process/task, which is currently given the CPU time, is allowed to execute until it terminates or enters the 'Blocked/Wait' state, waiting for an I/O or system resource. In co-operative multitasking, the currently executing process/task need not relinquish the CPU when it enters the 'Blocked/Wait' state, whereas in non-preemptive multitasking the currently executing task relinquishes the CPU when it waits for an I/O or system resource or an event to occur.

6.4 Task Scheduling

Multitasking involves the execution switching among the different tasks. Determining which task/process is to be executed at a given point of time is known as task/process scheduling. Scheduling policies forms the guidelines for determining which task is to be executed when. The scheduling policies are implemented in an algorithm and it is run by the kernel as a service. The process scheduling decision may take place when a process switches its state to

- Ready state form Running state
- Blocked/Wait state from Running state
- Ready state from Blocked/Wait state
- Completed state

The selection of a scheduling criterion should consider the following factors

- **CPU utilization:** the scheduling criterion should always make the CPU utilization high. CPU utilization is a direct measure of how much percentage of the CPU is being utilized.
- **Throughput:** This gives an indication of the number of processes executed per unit of time. The throughput for a good scheduler should always be higher.

- **Turnaround time:** It is the amount of time taken by a process for completing its execution. It includes the time spent by the process for waiting for the main memory, time spent in the ready queue, time spent on completing the I/O operations, and the time spent in execution. The turnaround time should be a minimal for a good scheduling algorithm.
- **Waiting Time:** It is the amount of time spent by a process in the 'Ready' queue waiting to get the CPU time for execution. The waiting time should be minimal for a good scheduling algorithm.
- **Response time:** It is the time elapsed between the submission of a process and the first response, for a good scheduling algorithm, the response time should be as least as possible.

The operating system maintains various queues in connection with the CPU scheduling, and a process passes through these queues during the course of its admittance to execution completion.

The various queues maintained by OS in association with CPU scheduling are:

- **Job Queue:** Job queue contains all the processes of the system.
- **Ready Queue:** contains all the processes, which are ready for execution and waiting for CPU to get their turn for execution
- **Device Queue:** contains the set of processes, which are waiting for an I/O device.

The scheduling algorithm can be classified as:

A. Non-preemptive Scheduling

It is employed in systems which implements non-preemptive multitasking model. In this scheduling type, the currently executing task/process is allowed to run until it terminates or enters the wait state waiting for an I/O or system resources. Various types of non-preemptive scheduling are listed below.

- **First Come First Served (FCFS) / FIFO Scheduling:** The FCFS scheduling algorithm allocates CPU time to the processes based on the order in which they enter the ready queue. The first entered process is serviced first. E.g. ticketing reservation system where people need to stand to a queue and the first person standing in the queue is serviced first.
- **Last Come First Served (LCFS) / LIFO scheduling:** The LCFS scheduling algorithm also allocates CPU time to the Processes based on the order in which they are entered in the ready queue. The last entered process is serviced first.
- **Shortest Job First (SJF) scheduling:** SJF scheduling algorithm sorts the ready queue each time a process relinquishes the CPU to pick the process with shortest estimated completion

time. The process with the shortest estimated run time is scheduled first, followed by the next shortest process, and so on.

- **Priority Based Scheduling:** This scheduling algorithm ensures that a process with high priority is serviced at the earliest compared to other low priority processes in the ready queue. The SJF algorithm can be viewed as a priority based scheduling where each task is prioritized in the order of the time required to complete the task. Another way of priority assigning is associating a priority to the task/process at the time of creation of the task/process. The priority is the number ranging from 0 to the maximum priority supported by the OS. For windows CE operating system a priority number 0 indicates the highest priority.

B. Preemptive Scheduling

Preemptive scheduling is employed in systems, which implements preemptive multitasking model. In this scheduling, every task in the ready queue gets a chance to execute. When and how often each process gets a chance to execute is dependent on the type of preemptive scheduling algorithm. In this scheduling method, the scheduler can preempt (stop temporarily) the currently executing process and select another task from the ready queue for execution. The task which is preempted by the scheduler is moved to the ready queue. The act of moving a running process into a ready queue by the scheduler, without the processes requesting for it is known as preemption. The different types of preemptive scheduling adopted in process scheduling are explained below.

- **Preemptive SJF Scheduling / Shortest Remaining Time (SRT):** The preemptive SJF scheduling algorithm sorts the ready queue when a new process enters the ready queue and checks whether the execution time of the new process is shorter than the remaining of the total estimated time for the currently executing process. If the execution time of the new process is less, the currently executing process is preempted and the new process is scheduled for execution. Preemptive SJF scheduling is also known as Shortest Remaining Time (SRT) scheduling.
- **Round Robin Scheduling:** In this scheduling method, each process in the ready queue is executed for a pre-defined time slot. The execution starts with picking up the first process in the ready queue. It is executed for a pre defined time slice and when the pre-defined time elapses or the process completes before the pre-defined time slice, the next process in the ready queue is selected for execution. Once each process in the ready queue is executed for

the pre-defined time period, the scheduler picks the first process in the ready queue again for execution and the sequence is repeated. So, the round robin scheduling is similar to the FCFS scheduling but time slice preemption is added to switch the execution between the processes in the ready queue.

- **Priority Based Scheduling:** Priority based preemptive scheduling algorithm is same as that of the non-preemptive priority based scheduling except for the switching of execution between processes. In preemptive scheduling, any high priority process entering the ready queue is immediately scheduled for execution whereas in the non-preemptive scheduling any high priority process entering the ready queue is scheduled only after the currently executing process completes its execution or only when it voluntarily relinquishes the CPU.

6.5 Task Synchronization

In a multitasking environment, multiple processes run concurrently and share the system resources. When two processes try to access display hardware connected to the system or two processes try to access a shared memory area where one process tries to write to a memory location while the other process is trying to read from this. Then, an issue will arise and hence each process must be made aware of the access of the shared resources. The act of making processes aware of the access of the shared resources by each process to avoid conflicts is known as task/process synchronization. Various synchronization issues may arise if processes are not synchronized properly.

Task Communication/Synchronization Issues

A. Racing

Racing or Race condition is the situation in which multiple processes compete each other to access and manipulate shared data concurrently. In a race condition the final value of the shared data depends on the process which acted on the data finally.

Suppose that two processes A and B have access to a shared variable Count:

Process A: $\text{Count} = \text{Count} + 5$

Process B: $\text{Count} = \text{Count} + 10$

Assume that process A and process B are executing concurrently in a time-shared, multi-programmed system.

Each statement requires several machine level instructions such as

For $\text{Count} = \text{Count} + 5$

A1: Load R_a, Count

A2: Add Ra, 05

A3: Store Count, Ra

For Count = Count + 10

B1: Load Rb, Count

B2: Add Rb, 10

B3: Store Count, Rb

In a time-shared or multi-processing system the exact instruction execution order cannot be predicted.

Scenario 1	Scenario 2
A1: Load Ra, Count	A1: Load Ra, Count
A2: Add Ra, 05	A2: Add Ra, 05
A3: Store Count, Ra	Context Switch
Context Switch	B1: Load Rb, Count
B1: Load Rb, Count	B2: Add Rb, 10
B2: Add Rb, 10	B3: Store Count, Rb
B3: Store Count, Rb	Context Switch
	A3: Store Count, Ra
Count is increased by 15	Count is increased by 5

B. Deadlock

A race condition produces incorrect results whereas a deadlock condition creates a situation where none of the processes are able to make any progress in their execution, resulting in a set of deadlocked processes. In its simplest form, 'deadlock' is the condition in which a process is waiting for a resource held by another process which is waiting for a resource held by the first process. For instance, process A holds a resource x and it wants a resource y held by process B. Process B is currently holding resource y and it wants the resources x which is currently held by process A. None of the competing process will be able to access the resources held by other processes since they are locked by the respective processes.

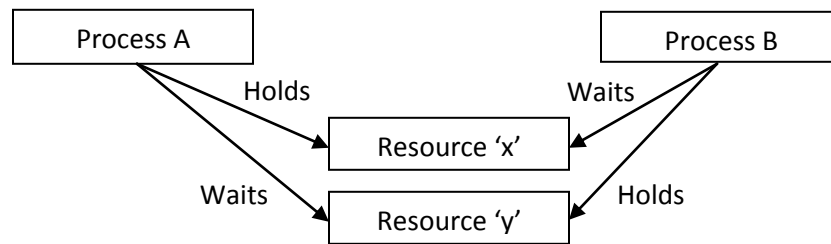


Figure 6.9: Scenarios leading to Deadlock

Coffman conditions: The different conditions favoring a deadlock situation are listed below

- **Mutual Exclusion:** The criteria that only one process can hold a resource at a time. Processes should access shared resources with mutual exclusion.
- **Hold and Wait:** The condition in which a process holds a shared resource by acquiring the lock controlling the shared access and waiting for additional resources held by other processes.
- **No Resource Preemptive:** The criteria that operating system cannot take back a resource from a process which is currently holding it and the resource can only be released voluntarily by the process holding it.
- **Circular Wait:** A process is waiting for a resource which is currently held by another process which in turn is waiting for a resource held by the first process. In general, there exists a set of waiting process $P_0, P_1 \dots P_n$ with P_0 is waiting for a resource held by P_1 and P_1 is waiting for a resource held by P_0 , ..., P_n is waiting for a resource held by P_0 and P_0 is waiting for a resource held by P_n and so on... This forms a circular wait queue.

Deadlock Handling

A smart OS may foresee the deadlock condition and will act proactively to avoid such a situation.

The OS may adopt any of the following techniques to detect and prevent deadlock conditions.

- **Ignore Deadlocks:** Always assume that the system design is deadlock free. This is acceptable for the reason the cost of removing a deadlock is large compared to the chance of deadlock to occur.
- **Detect and Recover:** This approach suggests the detection of a deadlock situation and recovery from it. OS keeps a resource graph in their memory. The resource graph is updated on each resource request and release. A deadlock condition can be detected by analyzing the resource graph by graph analyzer algorithms. Once a deadlock condition is detected, the system can terminate a process or preempt the resource to break the deadlocking cycle.

- **Avoid Deadlocks:** Deadlock is avoided by the careful resource allocation techniques by the operating system.
- **Prevent Deadlocks:** Prevent the deadlock condition by negating one of the four conditions favoring the deadlock situation.
 - Ensure that a process does not hold any other resources when it requires a resource.
 - Ensure resources preemption.

C. Livelock

In a livelock condition, a process changes its state with time but is unable to make any progress in the execution completion. While in deadlock a process enters a wait state for a response and continues in that state forever without making any progress in the execution. For example, two people attempting to cross each other in a narrow corridor. Both the person moves towards each side of the corridor to allow the opposite person to cross. Since the corridor is narrow, none of them are able to cross each other. Here both of the persons perform some action but still they are unable to achieve their target.

D. Starvation

In the multitasking context, starvation is the condition in which a process does not get the resources required to continue its execution for a long time. As time progresses the process starves on resource. Starvation may arise due to various conditions like byproduct of preventive measures of deadlock, scheduling policies favoring high priority tasks and tasks with shortest execution time, etc.

Task Synchronization techniques

Task synchronization is essential for:

- Avoiding conflicts in resource access in a multitasking environment.
- Ensuring proper sequence of operation across processes.
- Communicating between processes.

The code memory area which holds the program instructions for accessing a shared resource, shared variables is known as critical section. In order to synchronize the access to shared resources, the access to the critical section should be exclusive.

Consider two processes Process A and Process B running on a multitasking system. Process A is currently running and it enters its critical section. Before Process A completes its operation in the critical section, the scheduler preempts process A and schedules Process B for execution. Process B also contains the access to the critical section which is already in use by Process A. If process B continues its execution and enters the critical section which is already in use by Process A, a racing condition will be resulted. A mutual exclusion policy enforces mutually exclusive access of critical sections.

A. Mutual Exclusion through Busy Waiting/Spin Lock

The Busy Waiting technique uses a lock variable for implementing mutual exclusion. Each process/thread checks this lock variable before entering the critical section. The lock is set to 1 by a process/thread if the process/thread is already in its critical section; otherwise the lock is set to 0.

The major challenge in implementing the lock variable based synchronization is the non-availability of a single atomic instruction which combines the reading, comparing and setting of the lock variable. Most often the three different operations related to the locks, the operation of reading the lock variable, checking its present value and setting it are achieved with multiple low level instructions. The low level implementations of these operations are dependent on the underlying processor instruction set and the compiler in use.

Consider a situation where process 1 reads the lock variable and tests it and found that the lock is available and it is about to set the lock for acquiring the critical section. But just before process 1 sets the lock variable, process 2 preempts process 1 and starts executing. Process 2 contains a critical section code and it tests the lock variable for its availability. Since process 1 was unable to set the lock variable, its state is still 0 and process 2 sets it and acquires the critical section. Remember, process 1 was preempted at a point just before setting the lock variable. Now process 1 sets the lock variable and enters the critical section. It violates the mutual exclusion policy and may produce unpredictable results.

The above issue can be effectively tackled by combining the actions of reading the lock variable, testing its state and setting the lock into a single step. This can be achieved with the combined hardware and software support. Most of the processors support a single instruction Test and Set Lock (TSL) for testing and setting the lock variable. The TSL instruction call copies the value of the lock variable and sets it to a nonzero value.

The lock based mutual exclusion implementation always checks the state of a lock and waits till the lock is available. This keeps the processes always busy and forces the processes to wait for the availability of the lock for proceeding further. Hence this synchronization mechanism is known as Busy Waiting. This method is useful in handling scenarios where the processes are likely to be blocked for a shorter period of time on waiting the lock, as they avoid OS overheads on context saving and process re-scheduling. The drawback of Spin Lock based synchronization is that if the lock is being held for a long time by a process and if it is preempted by the OS, the other threads waiting for this lock may have to spin a longer time for getting it. The busy waiting mechanism keeps the process always active, performing a task which is not useful and leads to the wastage of processor time and high power consumption.

B. Mutual Exclusion through Sleep & Wakeup

The Busy waiting mutual exclusion enforcement mechanism used by processes makes the CPU always busy by checking the lock to see whether they can proceed. This results in the wastage of CPU time and leads to high power consumption. This is not affordable in embedded systems powered on battery. In sleep and wakeup mechanism, when a process is not allowed to access the critical section, which is currently being locked by another process, the process undergoes Sleep and enters the blocked state. The process which is blocked on waiting for access to the critical section is awakened by the process which currently owns the critical section. The process which owns the critical section sends a wakeup message to the process, which is sleeping as a result of waiting for the access to the critical section, when the process leaves the critical section. The sleep and wakeup policy for mutual exclusion can be implemented in different ways.

- **Semaphore:** It is a sleep and wakeup based mutual exclusion implementation for share resource access. Semaphore is a system resource and the process which wants to access the share resource can first acquire this system object to indicate the other processes which wants the shared resource that the shared resource is currently acquired by it. The resources which are shared among a process can be either for exclusive use by a process or for using by a number of processes in a time. The display device of an embedded system is a typical example for the shared resource which needs exclusive access by a process. The hard disk of a system is a typical example for sharing the resource among a limited number of multiple processes.

Binary Semaphore (Mutex): The binary semaphore provides exclusive access to shared resource by allocating the resource to a single process at a time and not allowing the other processes to access it when it is being owned by a process. Mutex is a synchronization object provided by OS for process synchronization. Any process can create a mutex object and other processes of the system can use this mutex object at a time. The state of a mutex object is set to signaled when it is not owned by any process, and set to non-signaled when it is owned by any process.

Counting Semaphore: The counting semaphore limit the access of resources to fixed number of processes or threads. It maintains a count between zero and a value. It limits the usage of the resource to the maximum value of the count supported by it. The state of the counting semaphore object is set to signaled when the count of the object is greater than zero. The count associated with a semaphore object is decremented by one when a process acquires it and the count is incremented by one when a process releases the semaphore object. The state of the semaphore object is set to non-signaled when the semaphore is acquired by the maximum number of processes that the semaphore can support.

- **Events:** Event object is a synchronization technique which uses the notification mechanism for synchronization. In concurrent execution we may come across situations which demand the processes to wait for a particular sequence for its operations. A thread/process can wait for an event and another thread/process can set this even for processing by the waiting thread/process. The creating and handling event objects for notification is OS kernel dependent.

Priority Inversion

Priority inversion is the byproduct of the combination of blocking based process synchronization and pre-emptive priority scheduling. It is the condition in which a high priority task needs to wait for a low priority task to release a resource which is shared between the high priority task and the low priority task, and a medium priority task which doesn't require the shared resource continue its execution by preempting the low priority task. Priority based preemptive scheduling technique ensures that a high priority task is always executed first, where as the lock based process synchronization mechanism ensures that a process will not access a shared resource, which is currently in use by another process. The synchronization technique is only interested in avoiding

conflicts that may arise due to concurrent access of the shared resources and not at all bother about the priority of the process which tries to access the shared resource.

Consider a three process A, B, C with priorities High, Medium and Low respectively. Process A and C share a variable X and the access to this variable is synchronized through mutual exclusion. Process C is ready and is picked up for execution by the scheduler and process C tries to access the shared variable X and acquires the semaphore to indicate the other processes that it is accessing the shared variable X. At the same time, process B enters the ready state with higher priority compared to C, so Process C gets preempted and B starts executing. Now if Process A enters the ready state at this point. Process B is preempted and process A is scheduled for execution. Process A involves access of shared variable X which is currently being accessed by process C. So process A is put into blocked state and process B gets the CPU and it continues its execution until it relinquishes the CPU voluntarily or enters a wait state or preempted by another high priority task. The high priority A has to wait till Process C gets a chance to execute and release the semaphore. This produces unwanted delay in the execution of the high priority task which is supposed to be executed immediately when it was ready.

The commonly adopted priority inversion workarounds are:

A. Priority Inheritance

A low priority task that is currently accessing a shared resource requested by a high priority task temporarily inherits the priority of that high priority task, from the moment the high priority task raises the request. Boosting the priority of the low priority task to that of the priority of the task which requested the shared resource holding by the low priority task eliminates the preemption of the low priority task by other tasks whose priority are below that of the task requested the shared resource and thereby reduces the delay in waiting to get the resource requested by the high priority task. The priority of the low priority task which is temporarily boosted to high is brought to the original value when it releases the shared resources. Priority inheritance handles priority inversion at the cost of run time overhead at scheduler. It imposes the overhead of checking the priorities of all tasks which tries to access shared resources and adjust the priority dynamically.

B. Priority Ceiling

In Priority Ceiling, a priority is associated with each shared resource. The priority associated to each resource is the priority of the highest priority task which uses this shared resource. This priority level

is called ceiling priority. Whenever a task accesses a shared resource, the scheduler elevates the priority of the task to that of the ceiling priority of the resource. If the task which accesses the shared resource is a low priority task, its priority is temporarily boosted to the priority of the highest priority task to which the resource is also shared. This eliminates the preemption of the task by other medium priority tasks leading to priority inversion. The priority of the task is brought back to the original level once the task completes the accessing of the shared resource. Priority Ceiling brings the added advantage of sharing resources without the need for synchronization techniques like locks. The priority of the task accessing shared resources is boosted to the highest priority of the task among which the resource is shared; the concurrent access of shared resource is automatically handled. Another advantage is that all the overheads are at compile time instead of run-time.

PRIORITY LIST	Process C acquires shared variables 'X'	Process B preempts C	Process A preempts B	Process A requires shared variable 'X', Priority of C is increased to High	Process C releases the shared resource, so A starts executing with that resource and the priority of C is lowered to its original value.	Process A completes its execution. B starts executing	Process B completes its execution and C starts its execution.
Process A			Running	Waiting	Running		
Process B		Running	Waiting			Running	
Process C	Running	Waiting		Running	Waiting		Running

Figure 6.10: Illustration of Priority Inheritance

6.6 Device Drivers

It is a piece of software that acts as a bridge between the OS and the hardware. The architecture of OS kernel will now allow direct device access from the user application. All devices related access should flow through OS kernel, and the OS kernel routes it to the concerned hardware peripherals. Device drivers are responsible for initiating and managing the communication with hardware

peripherals. They are responsible for establishing connectivity, initializing hardware (setting up various CPU registers) and transferring data.

Device drives which are part of OS are called built in drivers or on-board drivers. These drivers are loaded by OS at the time of booting the device and are kept in RAM. Device drivers which need to be installed for accessing a device are called installable drivers. Whenever the device is connected, the OS loads the corresponding driver into memory. Driver files are usually in the form of '.dll' files. Drivers can run either in user space or in kernel space. Device drivers which run in user space are called user mode driver and the driver which run in kernel space are called kernel mode drivers.

A device driver implements the following:

Device initialization and interrupt configuration: The driver configures the different registers of the device. The interrupt configuration part deals with configuring the interrupts that needs to be associated with the hardware. The basic interrupt configuration involves:

- Set the interrupt type (Edge triggered or Level triggered), enable the interrupts and set the interrupt priorities.
- Bind the interrupt with an interrupt request (IRQ). The processor identifies an interrupt through IRQ. These IRQs are generated by the Interrupt Controller. In order to identify and interrupt the interrupt needs to be bonded to an IRQ.
- Register an Interrupt Service Routine (ISR) with an IRQ. ISR is the handler for an interrupt. In order to service an interrupt, an ISR should be associated with an IRQ.

Interrupt handling and processing: An interrupt is served based on its priority, and the corresponding ISR is invoked. The processing part of an interrupt is handled in an ISR. The whole interrupt processing can be done by the ISR itself or by invoking an Interrupt Service Thread (IST). The IST performs interrupt processing on behalf of the ISR. Since interrupt processing happens at kernel level, user application may not have direct access to the drivers to pass and receive data.

Client Interfacing: The client interfacing implementation makes use of the Inter Process Communication mechanisms supported by the embedded OS for communicating and synchronizing with user applications and drivers. For example, to inform a user application that an interrupt is occurred and the data received from the device is placed in a shared buffer, the client interfacing code can signal an event.

NUMERICAL EXAMPLES

Example 1: Three processes with process IDs P1, P2, P3 with priorities 2, 3, 0 and estimated completion time 10, 5, 7 milliseconds respectively enter the ready queue together in the order P1, P2, P3. Calculate the Waiting Time and Turn Around Time for each process and also the Average Waiting Time and Average Turn Around Time. Assume there is no I/O waiting for the process. Use the following non-preemptive scheduling algorithms.

- First Come First Serve Scheduling
- Priority Based Scheduling
- Shortest Job First Scheduling

Solution:

Given information from the question are tabulated as shown below:

Process	Entry Time	Completion Time	Priority	Entered
P1	0	10	2	1 st
P2	0	5	3	2 nd
P3	0	7	0	3 rd

A. First Come First Served Scheduling

P1	P2	P3
0	10	15
		22

Execution Sequence of Processes

Waiting Time calculation

$P1 = (0-0) = 0\text{ms}$

$P2 = (10-0) = 10\text{ms}$

$P3 = (15-0) = 15\text{ms}$

Average Waiting Time

$= (0+10+15)/3$

$= 8.33\text{ms}$

Turn Around Time calculation

$P1 = (10-0) = 10\text{ms}$

$P2 = (15-0) = 15\text{ms}$

$P3 = (22-0) = 22\text{ms}$

Average Turn Around Time

$= (10 + 15 + 22)/3$

$= 15.67\text{ms}$

Waiting Time = Execution Start Point – Entry Point

Turn Around Time = Completion Point – Entry Point

B. Priority Based Scheduling

P3	P1	P2	
0	7	17	22

Execution Sequence of Processes

Waiting Time = Execution Start Point – Entry Point

Turn Around Time = Completion Point – Entry Point

Waiting Time calculation

P3 = (0 - 0) = 0ms

P1 = (7 - 0) = 7ms

P2 = (17 - 0) = 17ms

Average Waiting Time

= (0 + 7 + 17)/3

= 8ms

Turn Around Time calculation

P3 = (7 - 0) = 7ms

P1 = (17 - 0) = 17ms

P2 = (22 - 0) = 22ms

Average Turn Around Time

= (7 + 17 + 22)/3

= 15.33ms

C. Shortest Job First

P2	P3	P1
0	5	12
		22

Execution Sequence of Processes

Waiting Time calculation

$P2 = (0 - 0) = 0\text{ms}$

$P3 = (5 - 0) = 5\text{ms}$

$P1 = (12 - 0) = 12\text{ms}$

Average Waiting Time

$= (0 + 5 + 12)/3$

$= 5.67\text{ms}$

Turn Around Time calculation

$P2 = (5 - 0) = 5\text{ms}$

$P3 = (12 - 0) = 12\text{ms}$

$P1 = (22 - 0) = 22\text{ms}$

Average Turn Around Time

$= (5 + 12 + 22)/3$

$= 13\text{ms}$

Waiting Time = Execution Start Point – Entry Point

Turn Around Time = Completion Point – Entry Point

Example 2: Three processes with process IDs P1, P2, P3 with priorities 0, 1, 3 and estimated completion time 6, 9, 3 milliseconds respectively enter the ready queue together. If a new process P4 (priority 2) with estimated completion time 2ms enters the ready queue after 3ms of execution of P1. Calculate the Waiting Time and Turn around Time for each process and also the Average Waiting Time and Average Turn Around Time. Make use of following non-preemptive scheduling algorithm to solve the problem.

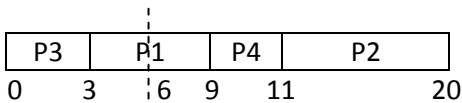
- Shortest Job First (SJF) Scheduling
- Priority Based Scheduling

Solution:

A. Non - Preemptive SJF Scheduling

Given information from the question are tabulated as shown below

Process	Entry Time	Completion Time	Priority
P1	0	6	0
P2	0	9	1
P3	0	3	3
P4	3ms after P1 starts	2	2



Execution Sequence of Processes

Waiting Time = (Execution Starting Point – Entry Point)

Turn Around Time = (Completion Point – Entry Point)

Waiting Time calculation

$$P3 = (0 - 0) = 0\text{ms}$$

$$P1 = (3 - 0) = 3\text{ms}$$

$$P4 = (9 - 6) = 3\text{ms}$$

$$P2 = (11 - 0) = 11\text{ms}$$

Average Waiting Time

$$= (0 + 3 + 3 + 11)/4$$

$$= 4.25\text{ms}$$

Turn Around Time calculation

$$P3 = (3 - 0) = 3\text{ms}$$

$$P1 = (9 - 0) = 9\text{ms}$$

$$P4 = (11 - 6) = 5\text{ms}$$

$$P2 = (20 - 0) = 20\text{ms}$$

Average Turn Around Time

$$= (3 + 9 + 5 + 20)/4$$

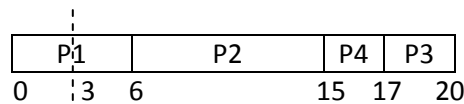
$$= 9.25\text{ms}$$

Explanation: Entry point for three processes P1, P2 and P3 is same at 0ms but the process P4 enters only after 3ms of execution of P1. So, the entry point for P4 will be at 6ms. Regardless of the shortest completion time of P4, P4 will not halt the execution of P1 as the algorithm is non pre-emptive. However, after the execution of P1, there remain two processes P2 and P4 with completion time 9ms and 2ms respectively. Hence, P4 will start to execute after completion of P1 according to shortest job first scheduling.

B. Non - Preemptive Priority Based Scheduling

Given information from the question are tabulated as shown below

Process	Entry Time	Completion Time	Priority
P1	0	6	0
P2	0	9	1
P3	0	3	3
P4	3ms after P1 starts	2	2



Execution Sequence of Processes

Waiting Time = (Execution Starting Point – Entry Point)

Turn Around Time = (Completion Point – Entry Point)

Waiting Time calculation

$$P1 = (0 - 0) = 0\text{ms}$$

$$P2 = (6 - 0) = 6\text{ms}$$

$$P4 = (15 - 3) = 12\text{ms}$$

$$P3 = (17 - 0) = 17\text{ms}$$

Average Waiting Time

$$= (0 + 6 + 12 + 17)/4$$

$$= 8.75\text{ms}$$

Turn Around Time calculation

$$P1 = (6 - 0) = 6\text{ms}$$

$$P2 = (15 - 0) = 15\text{ms}$$

$$P4 = (17 - 3) = 14\text{ms}$$

$$P3 = (20 - 0) = 20\text{ms}$$

Average Turn Around Time

$$= (6 + 15 + 14 + 20)/4$$

$$= 13.75\text{ms}$$

Example 3: Three processes P1, P2, P3 with estimated completion time 9, 4, 6 ms and priorities 1, 3, 2 respectively enters the ready queue together. A new process P4 with estimated completion time 4ms and priority 0 enters the ready queue after 2 ms of start of execution of P1. Calculate the Waiting Time and Turn Around Time for each process. Also Calculate the Average Waiting Time and Average Turn Around Time, using the Preemptive Shortest Job First Scheduling and Priority Based Scheduling.

Solution:

A. Preemptive SJF Scheduling

Given information from the question are tabulated as shown below

Process	Entry Time	Completion Time	Priority
P1	0	9	1
P2	0	4	3
P3	0	6	2
P4	2ms after P1 starts	4	0

P2	P3	P1	P4	P1	
0	4	10	12	16	23

Execution Sequence of Processes

Waiting Time = (Execution Starting Point – Entry Point) + Halted time

Turn Around Time = (Completion Point – Entry Point)

Waiting Time calculation

$$P2 = (0 - 0) = 0\text{ms}$$

$$P3 = (4 - 0) = 4\text{ms}$$

$$P4 = (12 - 12) = 0\text{ms}$$

$$P1 = (10 - 0) + (16 - 12) = 14\text{ms}$$

Average Waiting Time

$$= (0 + 4 + 0 + 14)/4$$

$$= 4.5\text{ms}$$

Turn Around Time calculation

$$P2 = (4 - 0) = 4\text{ms}$$

$$P3 = (10 - 0) = 10\text{ms}$$

$$P4 = (16 - 12) = 4\text{ms}$$

$$P1 = (23 - 0) = 23\text{ms}$$

Average Turn Around Time

$$= (4 + 10 + 4 + 23)/4$$

$$= 10.25\text{ms}$$

Explanation: Entry point for three processes P1, P2 and P3 is same at 0ms but the process P4 enters only after 2ms of execution of P1. So, the entry point for P4 will be at 12ms. At 12ms, there are two processes remaining; P1 with 7ms left to execute and P4 with 4ms. Since P4 is shorter compared to remaining part of P1, P4 will halt the execution of P1 at 12ms and starts its own execution. After P4 completes its execution at 16ms, P1 resumes.

B. Preemptive Priority based scheduling

Given information from the question are tabulated as shown below

Process	Entry Time	Completion Time	Priority
P1	0	9	1
P2	0	4	3
P3	0	6	2
P4	2ms after P1 starts	4	0

P1	P4	P1	P3	P2	
0	2	6	13	19	23

Execution Sequence of Processes

Waiting Time = (Execution Starting Point – Entry Point) + Halted time

Turn Around Time = (Completion Point – Entry Point)

Waiting Time calculation

$$P1 = (0 - 0) + (6 - 2) = 4\text{ms}$$

$$P4 = (2 - 2) = 0\text{ms}$$

$$P3 = (13 - 0) = 13\text{ms}$$

$$P2 = (19 - 0) = 19\text{ms}$$

Average Waiting Time

$$= (4 + 0 + 13 + 19)/4$$

$$= 9\text{ms}$$

Turn Around Time calculation

$$P1 = (13 - 0) = 13\text{ms}$$

$$P4 = (6 - 2) = 4\text{ms}$$

$$P3 = (19 - 0) = 19\text{ms}$$

$$P2 = (23 - 0) = 23\text{ms}$$

Average Turn Around Time

$$= (13 + 4 + 19 + 23)/4$$

$$= 14.75\text{ms}$$

Points to Remember

- When a process entering at the middle of execution does not halt the executing process, then its entry point and start of execution will never be at same point. Hence, its WT is never 0 and TAT is always greater than Completion Time.
- When a process entering at the middle of execution halts the executing process, then its entry point and start of execution will be same. Hence, it's $WT = 0$ and $TAT = \text{Completion Time}$.

- **Introduction**
- **Open-Loop and Closed-Loop Control Systems Overview**
- **General Control Systems and PID Controllers**
- **Software Coding of PID Controller**
- **PID Tuning**
- **Practical Issues Related to Computer-Based Control**
- **Benefits of Computer-Based Control Implementation**

7.1 Introduction

Control systems, a class of embedded systems, focus on tracking the reference input that is provided to the system. Initially the reference input is set and the output is more likely to track the same input regardless of the different external factors involved. The tracking can get difficult with the presence of disturbances. However, the system must be able to adjust to external factors for optimum performance. The objective of a control system is to track the reference output. The following figures represent good tracking and bad tracking respectively.

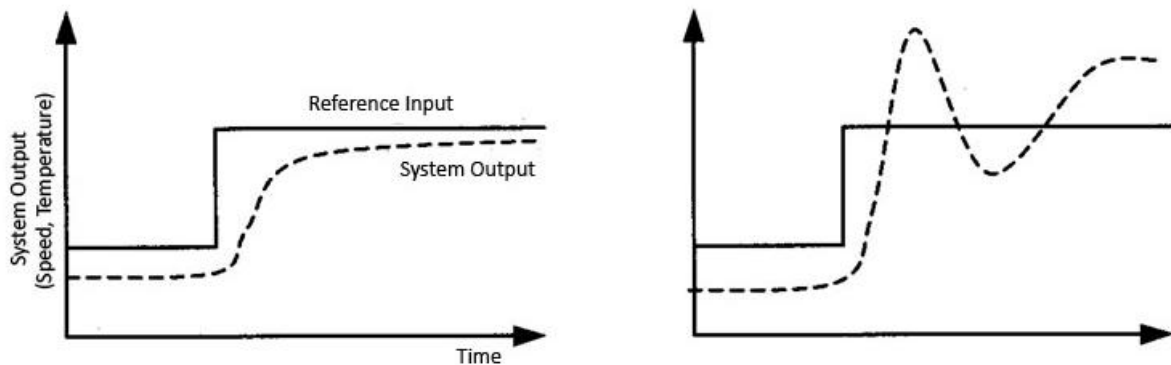


Figure 7.1: Good Tracking and Bad Tracking

7.2 Open-Loop and Closed-Loop Control Systems Overview

Open-Loop Control Systems are those systems in which the output has no influence on the control action of the input signal. It is also referred as feed-forward system or non-feedback system since the output is not fed back for comparison with the reference input. Also the controller is not aware about the tracking of reference input, so optimization is not possible. These systems are best utilized in case of predictable systems whose model is accurate and disturbance effect is minimal. In general, the open-loop control systems consist of following:

- **Plant**, which is also referred as a process, is the physical system to be controlled. Automobiles, fan, heater, disk etc are few examples.
- **Output** is the aspect or attribute of the physical system that we are about to control. Speed, temperature can be taken as examples.
- **Reference** input is the desired value that is required to be observed as an output of the physical system. Desired speed, temperature set by the user represents a reference input.
- **Actuator** is the device that is used to control the input to the plant. Motor can be taken as an example of an actuator.

- **Controller** is the main processing part of the system which computes the input to the plant such that desired output is achieved based on given reference input.
- **Disturbance** is an undesirable input to the system that may cause the output to deviate from the desired reference input.

The general block diagram of Open-Loop Control Systems is shown in the figure below.

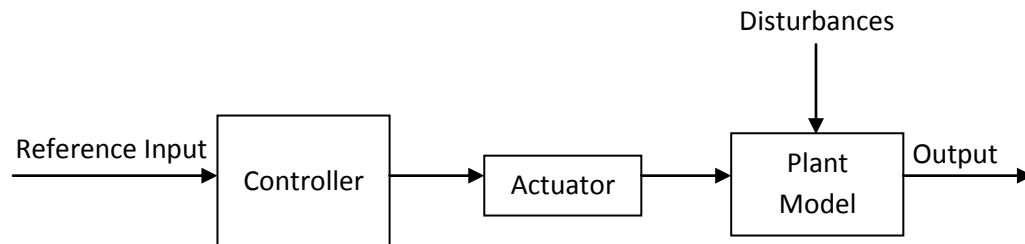


Figure 7.2: Block Diagram of Open-Loop Control System

Closed-Loop Control Systems are the systems operating on feedback principle. In such system the output is fed back, compared with the reference input and error signal is produced. The controller processes the error signal and reduces the error to obtain the desired output. Since the controller is aware about the output variations, optimization can be done and optimum performance can be obtained by minimizing the error. Apart from the plant, output, reference, controller, actuator, and disturbances, closed-loop control system contains additional components as sensor and error detector.

- **Sensor** is used to sense the output of the system and is fed to the input where error is calculated.
- **Error Detector** determines the error being produced in the system. Error is calculated by determining the difference between the output of the system and the reference input.

The general block diagram of Closed-Loop Control Systems is shown in the figure below.

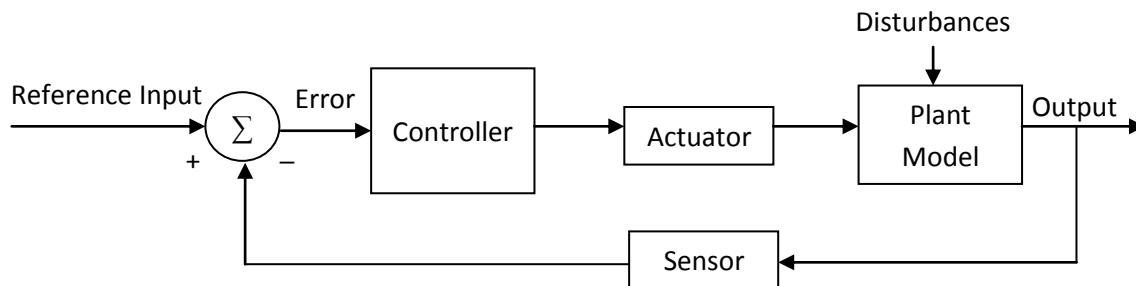


Figure 7.3: Block Diagram of Closed-Loop Control System

Comparison of Open-Loop and Closed Loop Control Systems

SN	Open-Loop Control System	Closed-Loop Control System
1.	Feed Forward System: Output is not fed back	Feed Back System: Output is fed back and compared with input
2.	It is simple and economical.	It is complex and expensive
3.	Good calibration can lead to good accuracy but optimization is not possible	Feedback principle reduces error, increases accuracy and supports optimization
4.	It is slow and unreliable but stable	It is fast and more reliable but unstable

7.3 General Control Systems and PID Controllers

Control Objectives

The main objective of control system design is to make output track the reference input even in the presence of measurement noise, model error and disturbances. The objective fulfillment can be analyzed and assessed through various metrics.

- **Stability:** For the system to be stable, all variables in the system remain bounded
- **Performance:** It describes how well the output is tracking the change in the reference input.

The various aspects of performance is shown in the figure below:

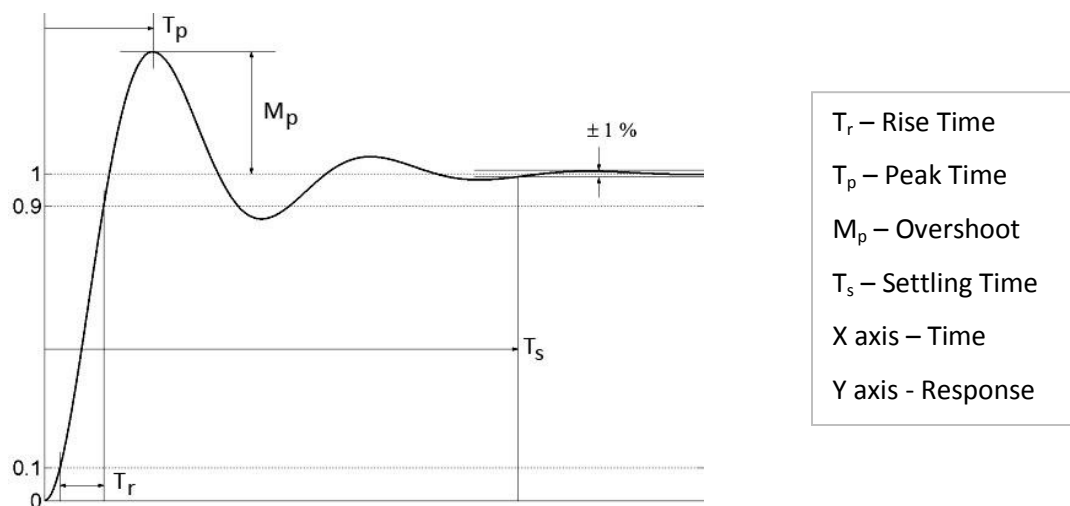


Figure 7.4: Aspect of performance metrics in Control system response.

The different aspects of performance are discussed below:

- **Rise Time (T_r)** is the time required to change from 10% to 90% of its final value. It is a measure of the ability of a system to fast input signals.
- **Peak Time (T_p)** is the time required to reach the first peak of the response.
- **Overshoot (M_p)** refers to an output exceeding its final, steady-state value. it is the percentage amount by which the peak of the response exceeds the final value.
- **Settling Time (T_s)** is the time required for the system to settle down to within 1% of its final value.
- **Disturbance rejection:** Disturbances are the undesired effects which cannot be eliminated but its impact can be minimized.
- **Robustness:** The system to be designed must be able to tolerate the modeling error of the plant. The stability and performance of the system should not be significantly affected by the presence of model errors.

Transient Response and Steady State Response of Control System

Transient response occurs just after the system starts and when any undesired conditions occur. The system's response during the settling time is transient response. Whereas the Steady state occurs after the system becomes settled. Steady State Error is defined as the difference between the actual output and the desired output when system reaches steady state.

Modeling Real Physical Systems

The accurate modeling of the behavior of the plant is an essential factor in control system design. Since the controller will be designed based on the plant model, the plant model must be accurate as far as possible. The key features of real systems are:

- **Continuous in nature:** It responds as continuous variables and as continuous function of time. Since real physical systems are continuously reacting, the plant model is represented by differential equations. Though continuous in nature, equivalent discrete time model can be determined. But the sampling period, however, must be selected much smaller than the reaction time of the system. Such sampling ensures system does not change much between sampling instants.

- **Complexity:** It is much more complex than any model we generally assume in our design. Our model may not include nonlinear effects, all system states, or all system interactions. Generally assumed model is a linear model which is sufficient when the variables of the model have a small operating range.

Controller Design

Proportional Control

A Controller that multiplies the tracking error by a constant is referred as proportional control. The form of proportional control is:

$$u(t) = P * e(t)$$

Where, $u(t)$ is the output of the controller, P is the proportional Constant, $e(t)$ is the measured error and is the difference between reference input and output of the system.

Proportional Constant affects transient response, steady state tracking error and disturbance rejection. High value of proportional constant can cause system to become unstable by resulting in high overshoot and oscillation, whereas low value of P will cause the system to be less response or less sensitive, since rise time will be high for low value of P . Also the steady state error will be high for low value of P . The following figure shows the response of an arbitrary control system for different values of P .

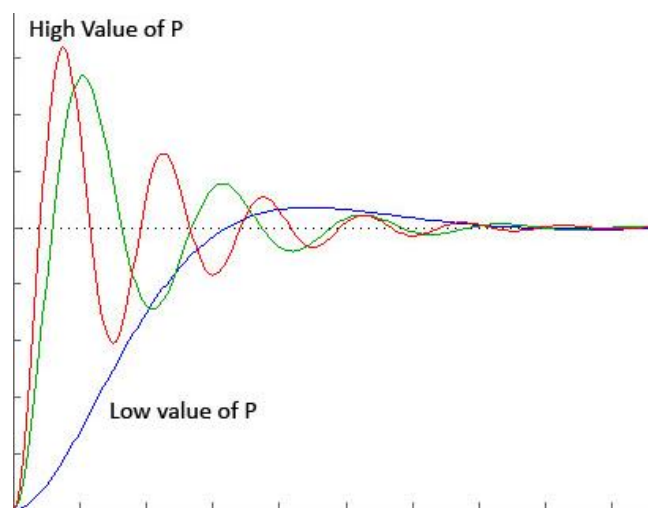


Figure 7.5: Response of system for different values of proportional constant

Proportional and Derivative (PD) Control

Derivative action predicts system behavior and improves settling time and stability of the system. Derivative term allows the transient response to be optimized without affecting the steady state response and disturbance rejection characteristic. Hence, transient response and the steady state error independently can be adjusted by using appropriate values of P and D in PD controller. The form of PD control is:

$$u(t) = P * e(t) + D * (e(t) - e(t-1))$$

Characteristics of PD control:

- Rise time reduces, improves damping, overshoot reduces, response is stable

The following figure shows the response of an arbitrary system for PD control action.

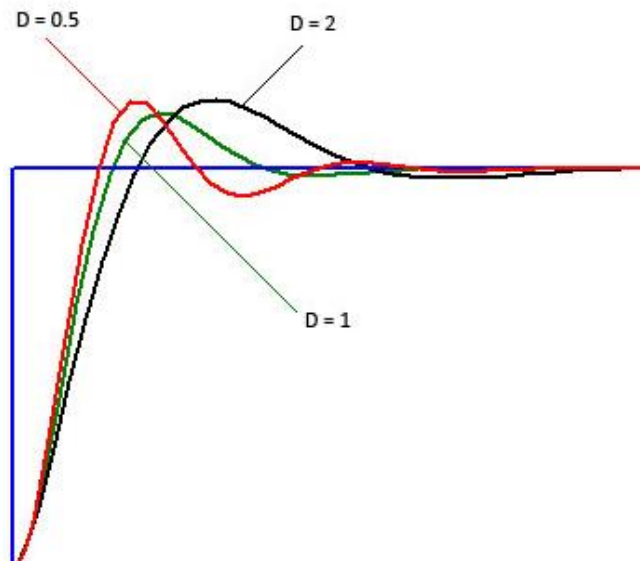


Figure 7.6: Effect of Derivative term

Proportional and Integral (PI) Control

A PI controller is a special case of the PID controller in which the derivative of the error is not used. The integral term in PI control is the sum of the instantaneous error over time and the accumulated error is multiplied by integral constant. Its output is given by

$$u(t) = P * e(t) + I * (e(0) + e(1) + e(2) + \dots + e(t))$$

PI controller is used to eliminate the steady state error resulting from P controller. However, it has undesirable impact on speed and stability of the system.

Characteristics of PI control

- Steady state accuracy improves, rise time increases, response is oscillatory

The following figure shows effect of different values of integral constant in the response of an arbitrary system for PD control action.

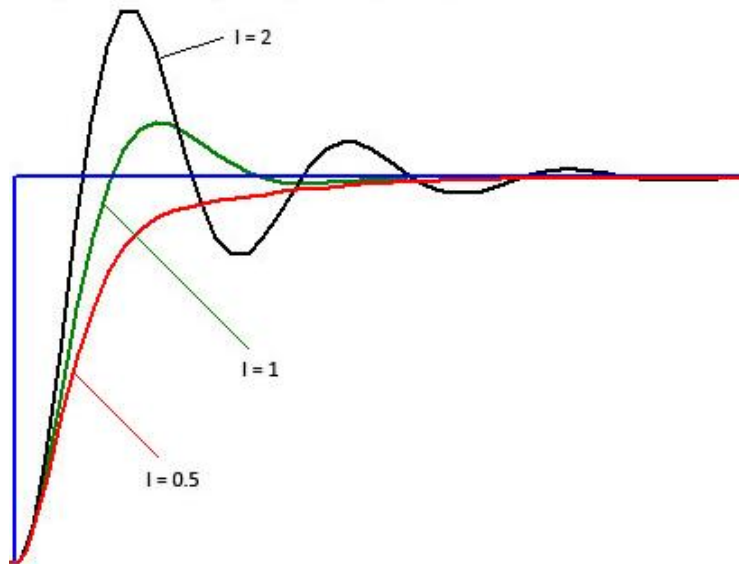


Figure 7.7: Effect of Integral Term in system response

Proportional Integral and Derivative (PID) Control

PID controller is a feedback controller that helps to attain a set point irrespective of disturbances or any variation in characteristics of the plant of any form. It calculates its output based on the measured error and the three controller gains; proportional gain P , integral gain K , and derivative gain D .

- The proportional gain simply multiplies the error by a factor P . It reduces steady state errors while minimizes the effect of external disturbances.
- The integral term is a multiplication of the integral gain and the sum of the recent errors. The integral term helps in getting rid of the steady state error and causes the system to catch up with the desired set point.
- The derivative controller determines the reaction to the rate of which the error has been changing and it increases damping and improves stability but has almost no effect on steady state error.

Its output is given by

$$u(t) = P * e(t) + I * (e(0) + e(1) + e(2) + \dots + e(t)) + D * ((e(1) - e(0)) + (e(2) - e(1)) + \dots + (e(t) - e(t-1)))$$

The general block diagram of PID controller is shown in the figure below:

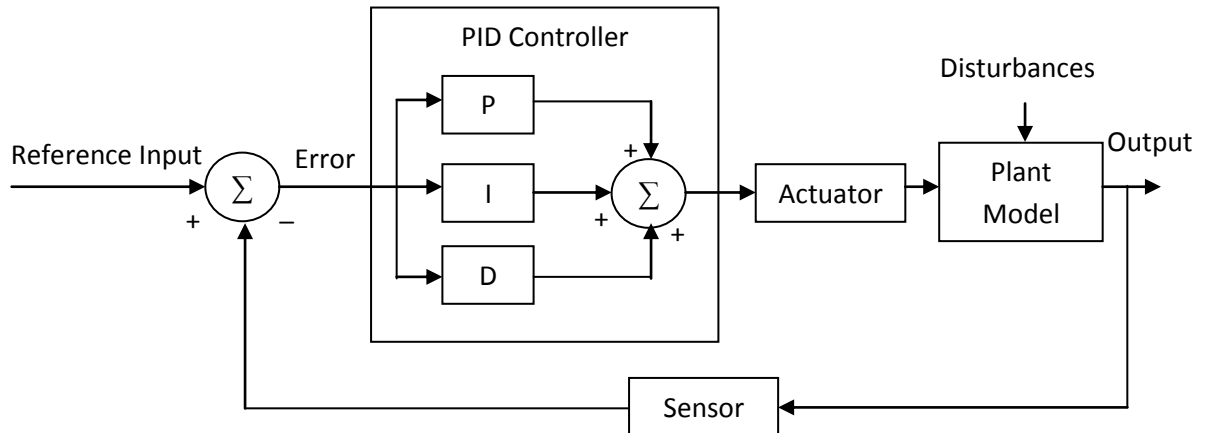


Figure 7.8: General Block Diagram of PID Controller

The following figure shows effect of different values of P, I, D in the response of an arbitrary system for PID control action.

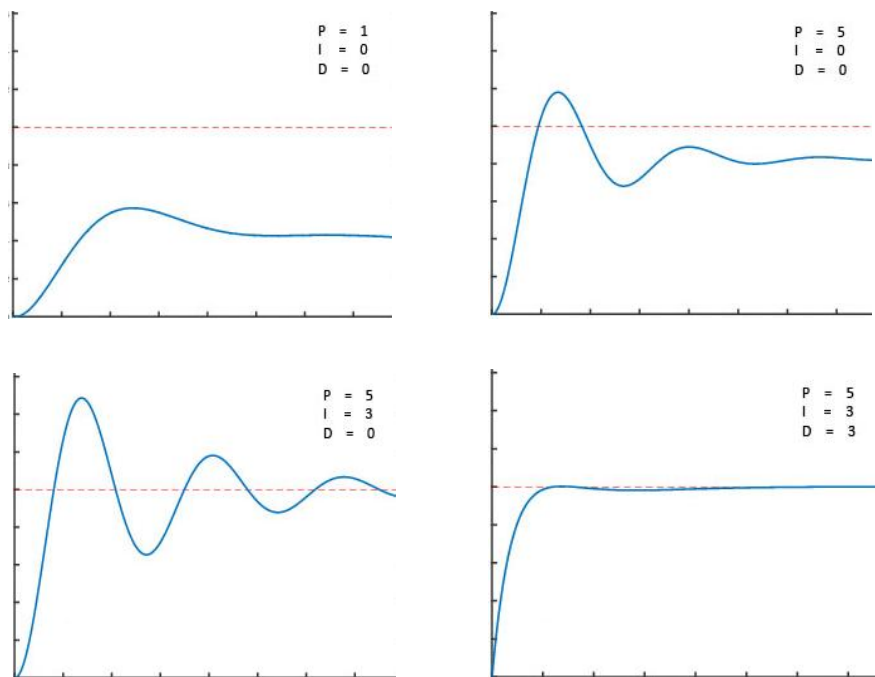


Figure 7.9: Effects of different values of P, I, D in the response of an arbitrary system

Summary of PID control action

Type	Rise Time	Maximum Overshoot	Settling Time	Steady-state error	Stability
P	Decrease	Increase	Small Change	Decrease	Degrade
I	Decrease	Increase	Increase	Eliminate	Degrade
D	Small Change	Decrease	Decrease	No/Small Change	Improve

(*Note: In above table, the effect is considered based on optimal value rather than increasing or decreasing the value of Proportional, Integral and Derivative constant)

7.4 Software Coding of PID Controller

A PID controller can be implemented using software. At first, required initialization is done which is followed by reading reference value and sensor value. Then, after that error can be calculated which further is used to compute the output of PID controller. The refined output is fed to the actuator which in turn controls the plant based on the value of proportional, integral and derivative constant defined in the program. The pseudo code for the PID controller can be written as:

- ➔ Set values for Pgain, Igain, Dgain
- ➔ Initialize prior_error = 0 and integral = 0
- ➔ Repeat following steps
 - sensorValue = getValueFromSensor()
 - refValue = getReferenceValue()
 - error = refValue – sensorValue
 - integral = integral + error*iterationTime
 - derivative = (error – prior_error)/iterationTime
 - output = Pgain * error + Igain * integral + Dgain * derivative
 - setActuator(output)
 - prior_error = error
 - wait(iterationTime)

7.5 PID Tuning

PID tuning is the adjustment of its control parameters to the optimum values for the desired control response. Quantitative analysis can be used to determine the values of P, I, and D.

However, quantitative analysis is not necessary when safety and cost of using plant is not a concern. There are various methods for PID tuning, one of which is ad hoc tuning process. The steps for ad hoc tuning process are

- Start with small value of P gain, D and I gains as 0
- Increase value of D gain until oscillation is seen, and then D gain is decremented by a factor of 2 to 4.
- Then, increase value of P gain until oscillation or excessive overshoot is observed, and then P gain is reduced by a factor of 2 to 4.
- Next, increase the value of I gain and reduce it slightly when oscillation or excessive overshoot is seen.
- Above steps are repeated until satisfactory performance is achieved.

7.6 Practical Issues Related to Computer-Based Control

The various practical issues related to computer-based control are explained in the following paragraphs.

a. Quantization and Overflow Error

Quantization error occurs when machine number is altered to fit the constraints of the computer memory.

- Case I: when arithmetic results require more precisions than original values. For example, in operation $0.50 \times 0.25 = 0.125$, the final result requires more precision.
- Case II: when analog signals from sensors are quantized by analog to digital converter it can create quantization error. In quantization process limited set of discrete values are defined and if the signal or value from the sensors doesn't match the defined quantized discrete values then rounding or truncation will occur which results in quantization error. For example: When 4 levels are defined between -1.5 and 1.5 as -1.5, -0.75, 0, 0.75 and 1.5 then the value 1.3 will be taken as 1.5.

Overflow error occurs when the system attempts to operate on or results a number that does not lie within the defined range of the system. For example, let us consider a case of signed binary numbers where five bits are used to represent the magnitude while sixth or MSB is used to represent sign. Using such representation, when two binary numbers 010010 (+18) and 010101 (+21) are added then it results in 100111 which is (-25) rather than (+39). Since the first bit is used for sign representation, the undesirable output resulted due to

overflow error. The situation can get more complex if we consider multiplication operation and floating point numbers.

b. Aliasing

Aliasing is the consequence of improper sampling process. It arises when a signal is discretely sampled at a rate that is insufficient to capture the changes in the signal. In simple term, aliasing causes the reconstructed signal to be different from original signal. It causes different signals to become indistinguishable. Let us consider an example in which the sampling is done at a period of 0.4 second which results a sampling frequency of 2.5 Hz. Then the following signals will be indistinguishable

$$y(t) = 1.0 * \sin(6\pi t), \text{ frequency } 3 \text{ Hz}$$

$$y(t) = 1.0 * \sin(\pi t), \text{ frequency } 0.5 \text{ Hz}$$

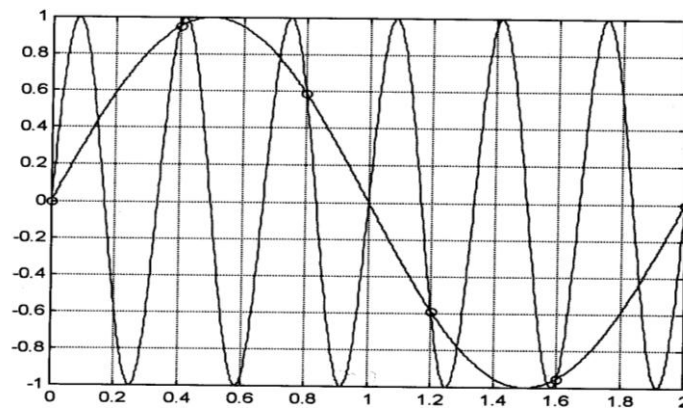


Figure 7.10: Aliasing Illustration

For a sampling rate of 2.5 Hz, sine wave with frequency of 0.5 Hz is indistinguishable from sine waves at 3 Hz, 5.5 Hz, 8 Hz and so on. Also, it can only correctly sample signal below Nyquist Frequency, which is half the value of sampling rate.

c. Computation Delay

Delay results in control signal being applied later than desired time. Computation delay is the attribute of many digital systems but too much delay results in performance degradation. The effect of delay can be accurately analyzed and we need to characterize implementation delay to ensure its effect is negligible in system's performance. Synchronous Design makes hardware delay to be characterized easily. Software delay, however, is harder to predict. So, code should be organized carefully to make delay

predictable. Also code can be written with predictable timing behavior, such that the effect of delay can be minimized to acceptable level.

7.7 Benefits of Computer-Based Control Implementation

The following are the benefits of computer-based control implementation.

a. Repeatability

Analog systems are more prone to aging, temperature and manufacturing tolerance effects which cause results to vary with time. However, the digital systems can produce identical results for longer time

b. Stability

Since digital systems are less prone to different sorts of degradations and optimizations can be implemented efficiently, systems can become more stable.

c. Programmability

Advanced features can be easily implemented in digital systems but that would be very complex in analog implementations. Few features include: control mode and gain switching, on-line performance evaluation, data storage, performance parameter estimation, and adaptive behavior.

d. Flexibility

Computer based control can be easily re-configured based on requirement which allows periodic upgrade and enhancement of the system. It permits modification of the sequencing and control procedures for different products and for frequent change in product specifications.

- **Introduction**
- **Full-Custom (VLSI) IC Technology**
- **Semi-Custom (ASIC) IC Technology**
- **Programmable Logic Device (PLD) IC Technology**

8.1 Introduction

A structural representation of the system generally deals with the various components and their interconnections to implement system's functionality. IC technology is more about mapping the structural representation to a physical implementation. The physical implementation can be done using various methods, out of which full-custom, semi-custom and programmable technologies are few common methods. As CMOS transistor is the core of every component, let us take a look at CMOS transistor and different layers required for its physical implementations.

CMOS Transistor

CMOS transistor consists of three terminals: the source, drain, and gate. Source and drain are created by implanting ions on the surface of silicon. Gate is formed using poly-silicon, and lies between source and drain. Gate is placed on top of silicon and is isolated from silicon with the help of silicon dioxide insulating layer. Gate voltage controls the current flowing from source to the drain. In case of nMOS transistor, a high voltage at gate will attract electrons from silicon substrate towards it resulting in formation of conducting channel between source and drain. For a low voltage at gate, the conducting channel is not formed.

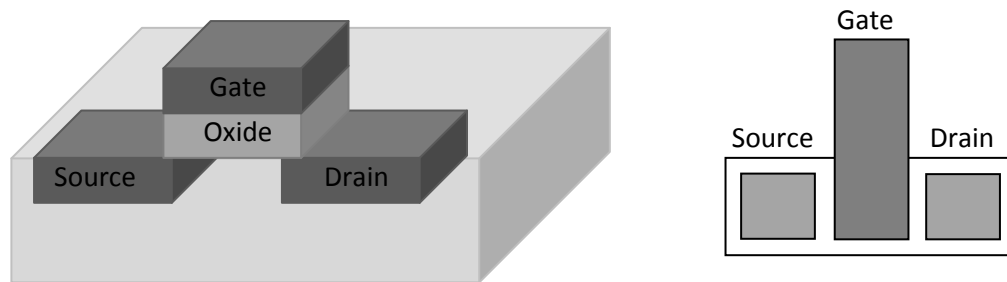


Figure 8.1: CMOS Transistor and Its Top-Down View

Layers in Physical Implementation

The transistor basically has three layers: diffusion layer for source and drain, oxide layer for insulation, and poly-silicon layer for gate. For circuits, there will be number of transistors connected together to represent particular functionality. These connections are represented by metal layers. There can be number of metal layers based on complexity of circuit implemented. Each metal layer is insulated from another layer using oxide layer. Hence, there exists number of oxide layer.

Metal2 Layer		
Oxide Layer		
Metal1 Layer		
Oxide Layer		
Poly-silicon Layer		
Oxide Layer		
pdiff		ndiff
Silicon Substrate		

Figure 8.2: Basic Layers in Physical Implementation

Example 1: Draw the transistor level circuit schematic and top-down view for a NAND gate

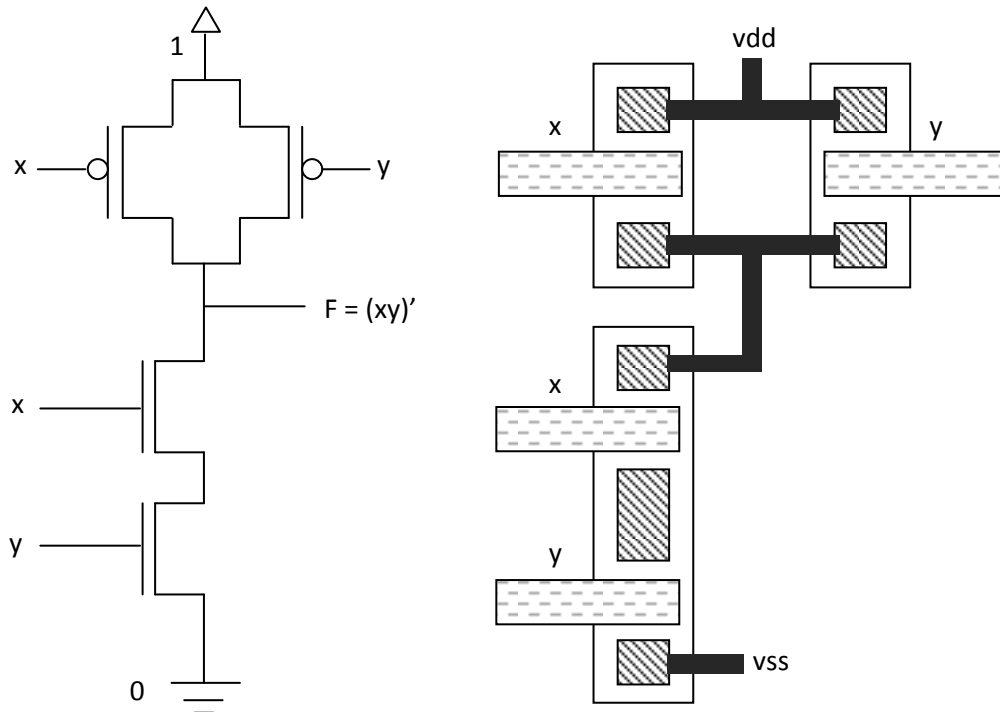


Figure 8.3: Circuit Schematic and top-down view of NAND gate

IC Manufacturing Process

Basically, IC manufacturing process can be divided into two phases: design phase and manufacturing phase. In design phase, structural design and layout design is done, whereas manufacturing phase includes various steps from mask creation to final IC packaging.

A. Design Phase

In design phase, the structural description along with the layout of the system is developed. Initially, the behavioral description of the system is implemented using hardware description language. The high-level HDL describes the circuit at the Register Transfer Level. The first step in the synthesis process is compilation which converts high-level VHDL language into a netlist at the gate level. The second process is speed and area optimization which is performed on gate-level netlist. Finally, the physical layout of the system is generated with the help of place-and-route software. The layout specifies the placement of every transistor and every wire connecting those transistors. Several EDA (Electronic Design Automation) tools are available for circuit synthesis, implementation, and simulation.

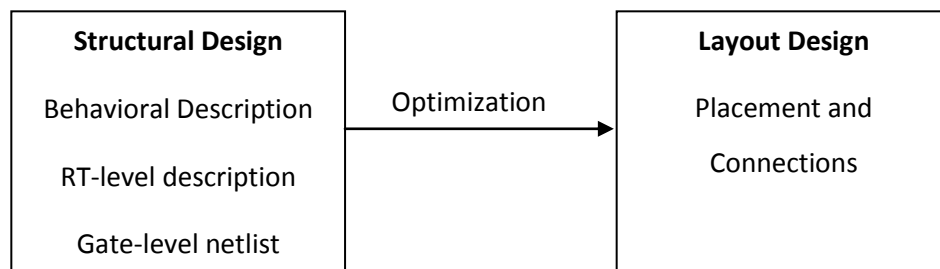


Figure 8.4: Design phase in IC manufacturing process

B. Manufacturing Phase

Manufacturing consists of several steps which are shown in the figure below and later each step is explained briefly.

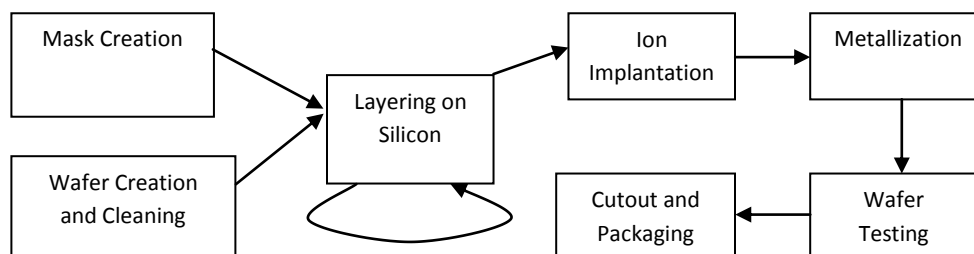


Figure 8.5: Manufacturing phase in IC manufacturing process

Mask Creation: The layout design of the system is translated into masks. The number of masks requirement may vary based on number of layers defined by the systems complexity. Masks for different layers – such as oxide layer, metal layers, etc – are generated. Generally, masks contain number of identical regions, so that number of IC's can be produced at once.

Silicon Wafer Creation and its cleaning: In a crucible, high purity silicon is melted. Donor impurity atoms can be added to dope the crystal. A seed crystal is dipped into molten silicon and pulled upwards rotating it. And cylindrical ingot is extracted by controlling temperature gradients, rate of pulling and speed of rotation. Finally, the ingot is sliced with a wafer saw and polished to form wafers.

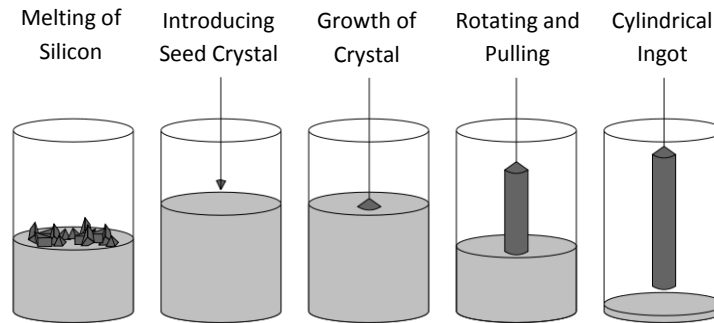


Figure 8.6: Silicon Wafer Creation

Wafer must be cleaned before any layer is deposited on it. Various cleaning methods can be used. Chemical cleaning methods are commonly used. First method of chemical cleaning is by using piranha solution in which wafer is immersed in hot mixture of hydrogen peroxide and sulfuric acid. Another method is using sonic waves in cleaning solution which is known as megasonic cleaning process. After the wafer is cleaned with chemical, it must be rinsed with De-ionized (DI) water. Finally, the wafer is dried using either nitrogen gun or by baking. Also spun dry method can be used to make the wafer dry after cleansing process.

Layering on Silicon: Various layers are developed on the silicon surface. Layer for masks can be created using different layering techniques. Photolithography, which uses optical radiation to create patterns, is very common method in layering process. In this process, the layer required, for example silicon dioxide, is built onto the silicon surface which is overlapped by photoresist. Positive photoresist becomes soluble when UV rays are exposed on it. Using proper alignment, the UV rays are passed through the masks which cast a shadow on the photoresist wherever the layer of silicon dioxide is required. Then the soluble photoresist is washed using appropriate solvent. Finally, the exposed silicon dioxide is etched away using chemicals and the remaining photoresist is removed to expose the regions of silicon dioxide that we required in our layer. The whole process is repeated for each layer.

Ion Implantation: Ions are accelerated at a very high energy and impinged on the target. The ion energy ranges from several KeV to MeV. The main purpose of ion implantation is doping in which impurities are added into wafer. This process is finalized with annealing process that repairs the lattice damage inflicted by high energy ions.

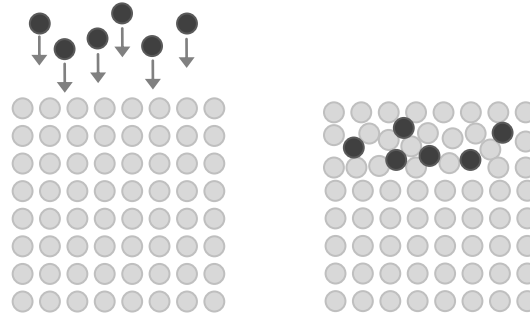


Figure 8.7: Ion Implantation and Lattice Damage after Implantation

Metallization: in this stage, a thin-film metal layer is produced which interconnects various circuit elements on the chip. Metallization also produces metalized area around the edge of the chip, also referred as bonding pads. Metal film can be deposited by physical vapor deposition (PVD) and chemical vapor deposition (CVD).

Wafer Testing: Number of ICs is produced in a single silicon wafer, which are subjected to test for errors or faulty ones. Testers or wafer probes are equipments used to test the correctness of the IC's by inspecting the output response for the streams of input.

Chip Cutout/Packaging: Individual IC from wafer is cut out using a diamond scribe. Verified ones are mounted in an IC package which encapsulates the IC. Packaging prevents physical damage and corrosion, also supports electrical contact. Through hole package and surface mount package are examples of IC packaging. Single In-line Packaging and Dual In-line Packaging are types of through hole packaging.

PHOTOLITHOGRAPHY

Photolithography is the process which transfers a pattern from a mask to a light-sensitive chemical photoresist on the substrate. The word photolithography is from the Greek origin: photo means light, litho means stone and graphy means writing. It uses optical radiation to create patterns of complex circuit on a wafer. The various steps involved in photolithographic process are deposit barrier layer, photoresist coating, soft bake, mask alignment and exposure, develop photoresist, hard bake, etch window in barrier layer and remove photoresist.

The various steps of photolithography are explained below:

A. Deposit Barrier Layer

Barrier layers are the materials which are required to be laid on the substrate. It may be silicon dioxide, silicon nitride, poly-silicon, metals, etc. Different methods can be used for barrier formation: thermal oxidation, chemical vapor deposition, sputtering and vacuum evaporation. Silicon dioxide as a barrier layer is used to isolate one layer from another. For instance, it is used in electrical isolation of multilevel metallization. Silicon Dioxide can be grown using dry oxidation which uses O_2 gas in a chamber or wet oxidation in which the wafer is submerged in water. When heat is applied to the oxidation process, it increases the rate of SiO_2 growth.

B. Photoresist Coating

Photoresist is a substance which changes its characteristics when exposed to UV light. Before photoresist is coated, hexamethyldisilazane (HMDS) is used on the surface to improve adhesion. After that, liquid photoresist is coated over barrier layer using spin coating method. In this method, the wafer is held on vacuum chuck which is spun at about 3000-6000 rpm for about 15-30 seconds. Appropriate spinner rotational speed and viscosity of resist are essential factors to define photoresist's thickness which is about few micro-meters.

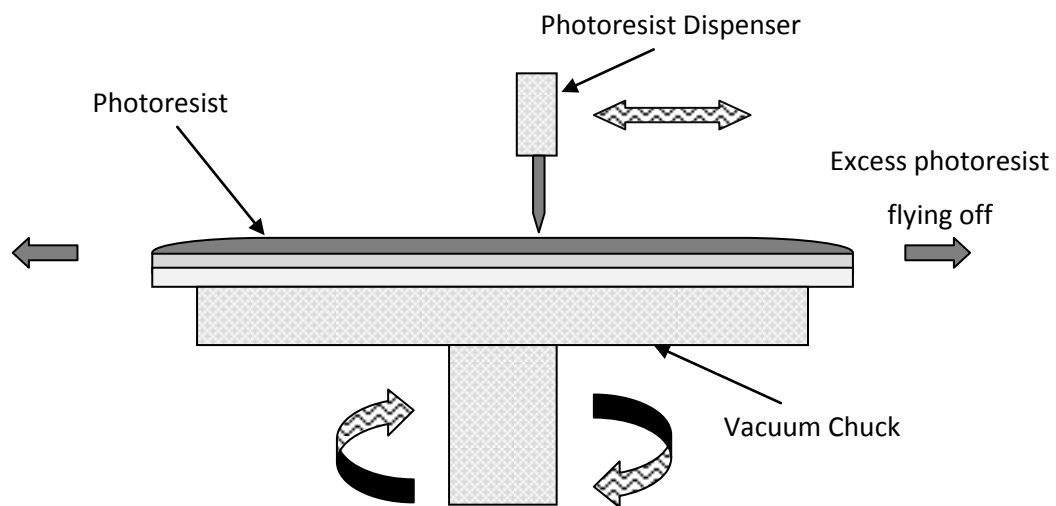


Figure 8.8: Photoresist Coating – Spin Coating Method

Types of Photoresist

Positive photoresist is insoluble in normal state but becomes soluble when exposed to UV light. Negative photoresist is soluble in normal state but becomes insoluble when exposed to UV light.

C. Soft Bake or Pre Bake

Soft bake is simply the process of heating the wafer which removes the solvent from the photoresist. Baking time and temperature depend on the type of photoresist used and baking method. Different baking methods include hotplate, oven baking and microwave baking.

D. Mask Alignment and Exposure

Mask is simply an opaque plate with holes to pass UV rays. It contains pattern to be formed on wafer. Mask is aligned with the wafer accurately with the help of special device: steppers use automatic pattern recognition and alignment systems. Alignment marks are available on the mask and on wafer so as to make alignment more precise.

Once the mask has been precisely aligned, the photoresist is exposed through the pattern on the mask with a controlled amount of UV light. Exposure will cause exposed positive photoresist to become soluble whereas if negative photoresist is used then exposed part of it becomes insoluble. There are three primary exposure methods: contact, proximity, and projection.

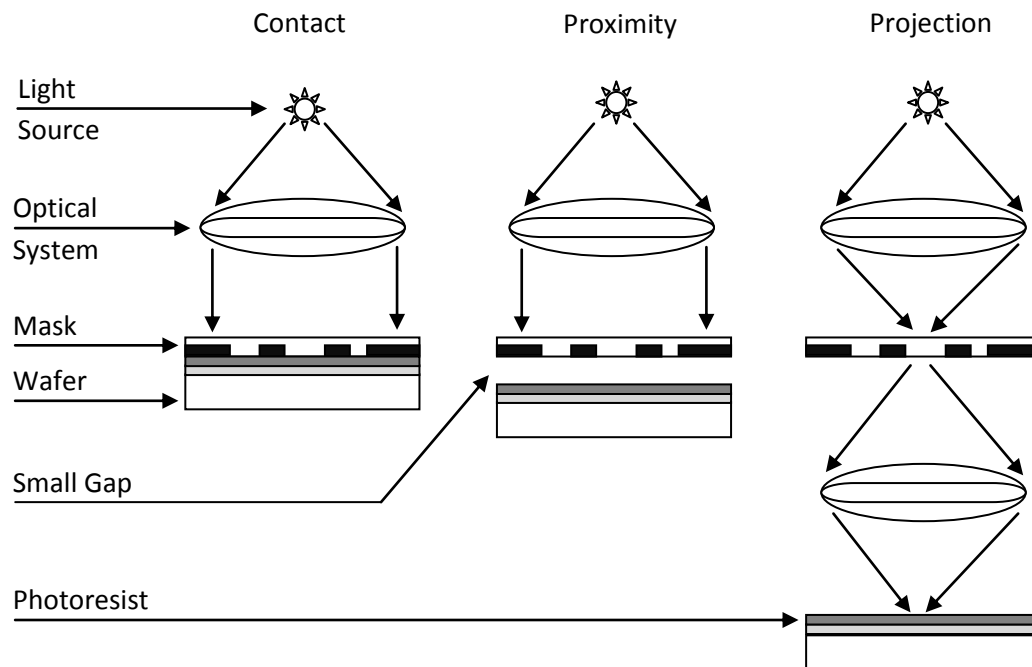


Figure 8.9: Different Exposure Methods

Contact Printing: In contact printing, the resist-coated silicon wafer and mask are brought into physical contact when exposed to UV light. This method results in very high resolution but the

debris, trapped between the resist and the mask, can damage the mask and cause defects in the pattern.

Proximity Printing: In this method, small gap is maintained between wafer and the mask during exposure. The gap minimizes the risk of mask damage at the expense of resolution.

Projection Printing: in this printing method, an image of the patterns on the mask is projected onto the resist-coated wafer. High gap eliminates the risk of mask damage and high resolution is possible. For high resolution, only a small portion of the mask is imaged and stepped over the surface of the wafer.

E. Develop Photoresist

Barrier layer is exposed when the soluble photoresist is chemically washed away using a developer solution. In immersion develop method the photoresist-coated wafer is immersed in a developer solution. Then, it is rinsed with DI water and dried using spin dry method.

F. Hard Bake or Post Bake

Hard bake is used to stabilize and harden developed photoresist. It not only improves adhesion of the photoresist but also removes traces of solvent or developer solution. But, however, improper post bake can cause resist removal more difficult. Baking time and temperature can vary based on type of photoresist and baking method.

G. Etch Window in Barrier Layer

As hardened photoresist does not shield all part of barrier layer, etching method is implemented to remove the barrier layer which was left uncovered. Two methods of etching can be implemented: wet etch, also known as chemical etching, and dry etch, also known as plasma etching. In Wet etching method, wafer is submerged in HF acid and unprotected barrier layer is removed. Dry etch method uses plasma which collides with the surface and removes the layers of target material.

H. Remove Photoresist

Finally, the remaining photoresist is stripped from the surface exposing the required barrier layer. Photoresist can be removed by using solvent strippers, which cause the resist to swell and lose adhesion from the substrate. Another method of photoresist removal is by burning the resist in an oxygen plasma system and this process is called Resist Ashing.

The photolithography process can be summarized diagrammatically as:

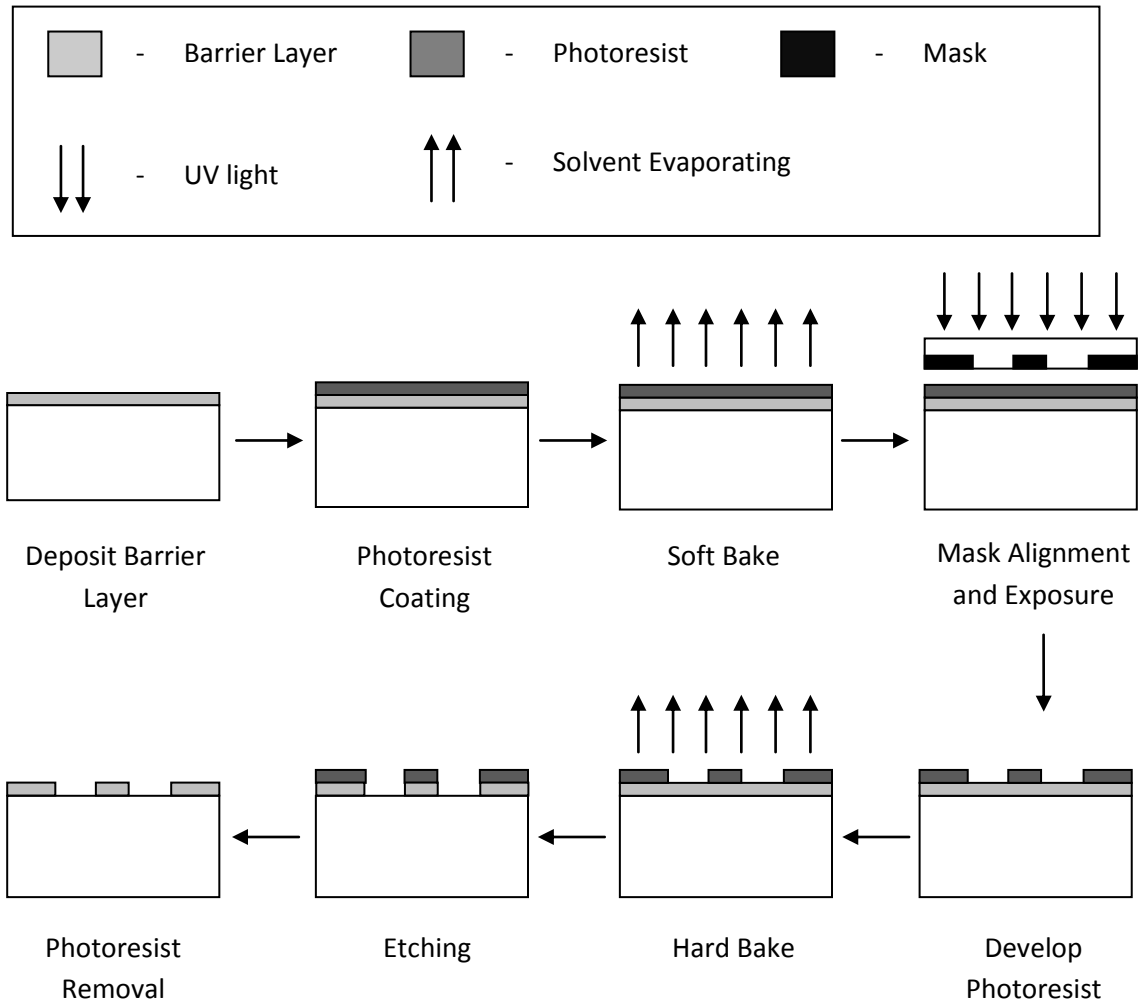


Figure 8.10: Various steps of photolithography process

8.2 Full-Custom (VLSI) IC Technology

Full-Custom IC technology includes VLSI (Very Large Scale Integrated Circuit) design in which the designer designs the complete transistor-level circuit for every processor, memory and other components used in the design. In this technology, first the designer creates layouts for basic components. And then, components are placed and connected, which are later translated to masks. Finally, the masks are given to the manufacturer for fabrication of IC of final design. The design steps are shown in the figure 8.11.

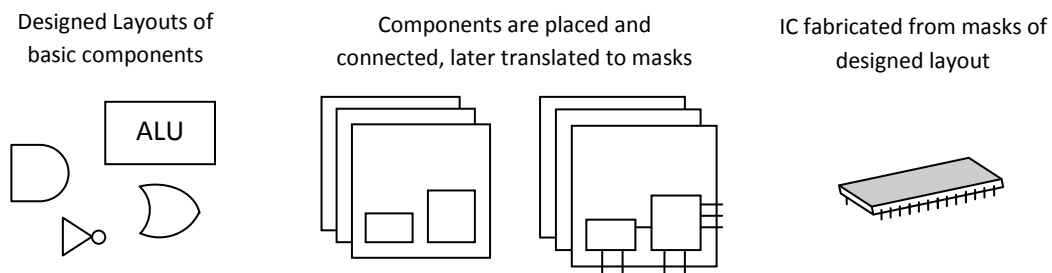


Figure 8.11: Full-Custom IC Technology

Placement, routing and sizing are few important physical design tasks that should be done carefully for an efficient layout design. Placement represents the task of placing and orienting the transistors on the IC. Routing is the task of connecting wires between the transistors. In sizing, width of each wire along with size of transistor is taken into considerations. Placement and routing should be done so as to avoid overlapping of transistors and wires. Placement also defines the length of wire required to connect transistors. Large size of wires and transistors provide better performances, but it increases power consumption and demands more silicon area in the IC. Compact layout can lead to an efficient design. For instance, transistor placed at closer distance requires shorter connecting wires, which further decreases the silicon size in the IC. In early days, compact designs were implemented using hand layout technique which is generally used for small and critical components. Today, however, physical design tools are used for automatic layout of the design which runs for hours or days to generate the optimized layout for better performance.

Advantages

Excellent efficiency: With respect to power consumption, performance and size, full-custom IC technology can be highly efficient. Since layout design is done by the designer, the components can be placed closer to each other which can be connected using short wires. Such layout yields optimum performance, size and power.

No wasted area and no unused transistors: In full-custom design, the required transistors for the circuit are placed on the IC. But there are no unused transistors which prevents wasted area.

Disadvantages

High NRE Cost and Long Time-to-market: Designing a complete layout, even with the help of CAD tools, can be time-consuming and prone to error. In addition to that, creating masks for

every layer of IC adds more time in design process. Also, manufactured IC may contain errors leading to requirement of several re-spins. All these factors cause full-custom IC technology to have high NRE cost and long time-to-market.

8.3 Semi-Custom (ASIC) IC Technology

In semi-custom IC technology, designer does not require to create full-custom layout rather connects the pre-positioned building blocks. The use of chip with pre-existing gates will lessen the design work of layout and mask creation. So, the NRE cost is reduced while the time-to-market is relatively fast as compared to full custom IC technology. But, however, there will be a reduction in performance in terms of power, size and speed. Two types of semi-custom IC technologies are described in the following paragraphs.

Gate Array Semi-Custom IC Technology

In a gate array IC technology, a chip with arrays of pre-designed logic gates is utilized to implement the desired circuit. Here, the masks for transistor and gate levels are already designed, so the designer has the task of connecting pre-designed gates to achieve the desired implementation. In this technology, a set of masks of predefined gates are provided to the designer who then provides the connections among gates to implement required circuit. Masks of connections are generated and all masks are used to fabricate the IC.

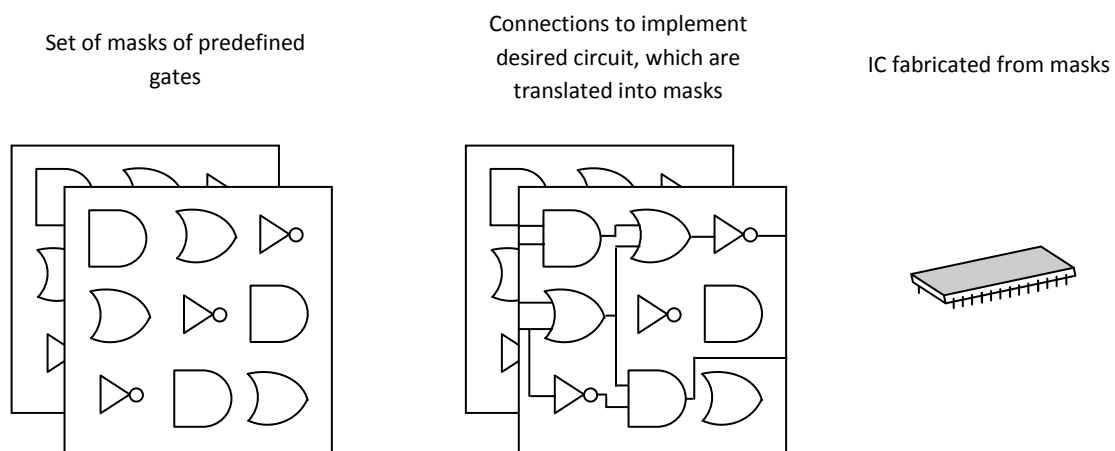


Figure 8.12: Gate Array Semi Custom IC Technology

This technology results in fast and relatively inexpensive design cycles. But, gates are placed in advance which may result in many unused gates, since all instances of each type of gate may not

be required in our desired circuit. Also, the fixed placement of gates can result in long routing wires between gates as the connection is not known while gates are already placed.

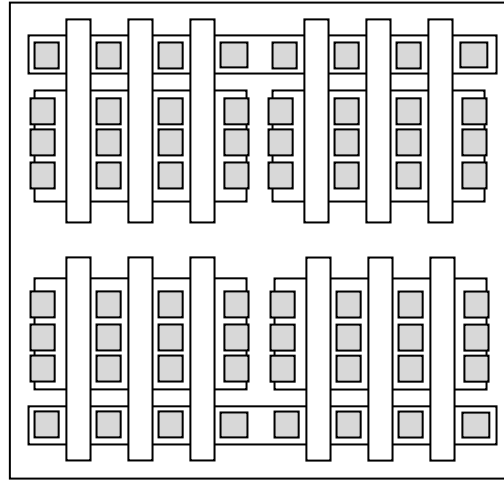


Figure 8.13: A simplified gate array layout

Standard Cell Semi-Custom IC Technology

In standard cell semi-custom IC technology, functional blocks, which are also called cells, with known electrical characteristics are utilized in the design to achieve very high gate density and good electrical performance. Cells may include logic gates such as NAND, NOR etc and other function blocks like multiplexor, flip-flop etc. In this technology, designers are facilitated with a library of predesigned cells from which the designer selects the required cells that are needed in the desired circuit. Masks of cells are created after the cells are placed and connected. Also the masks of connections among cells are generated. Using those masks, the IC is fabricated.

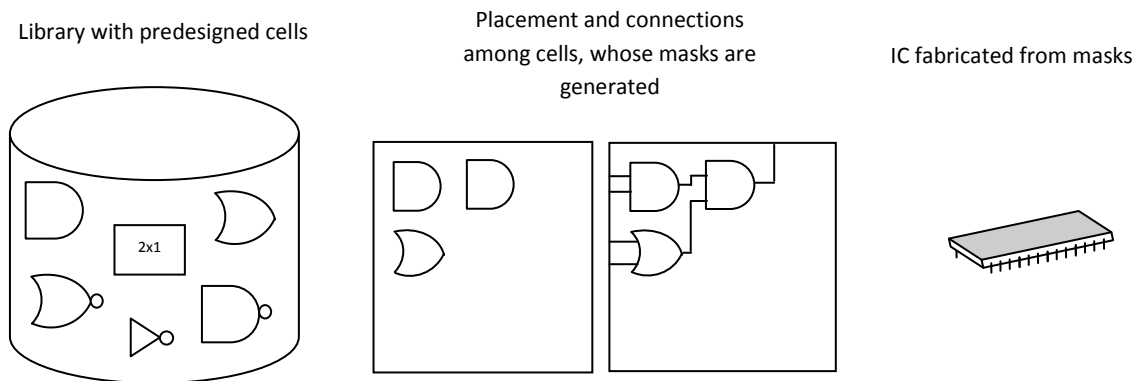


Figure 8.14: Standard Cell Semi-Custom IC Technology

The designer selects the cell, its position and its routing mechanism. So, it requires more NRE cost and longer time-to-market as compared to gate-array technology but still requires less than that of full-custom. However, the efficiency is better compared to gate array but less efficient than full-custom design. Hence, standard cell design lies between gate array and full custom design in terms of NRE cost, time-to-market and performance.

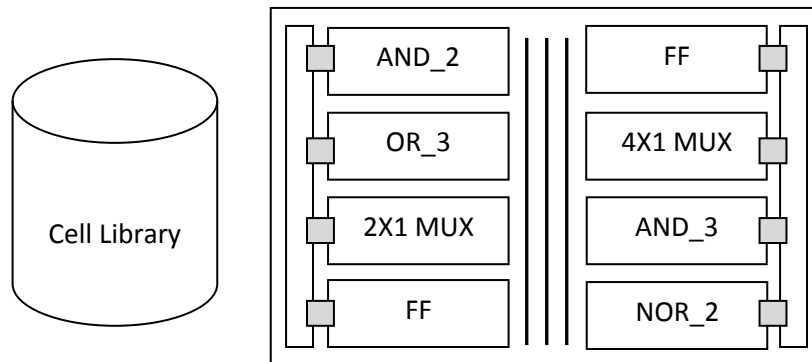


Figure 8.15: A simple standard cell layout

8.4 Programmable Logic Device (PLD) IC Technology

In programmable logic device IC technology, there exist programmable circuits which are programmed by the designer to implement the required design. Programming, in this case, means creating or breaking connections between wires that connect gates, either by blowing a fuse with high current, or setting a bit in a programmable switch. In this technology, a pre-fabricated chip with no logic function programmed is made available to the designer who then programs the required portions of the chip to implement the desired functionality.

It offers the designer the facility of changing design functions even after it has been programmed. PLD can be programmed, erased, and reprogrammed number of times, allowing easier prototyping and design modification.

There is a wide variety of PLD types, including Simple PLD, Complex PLD, GAL (Generic Array Logic), FPGA (Field-Programmable Gate Array) as well as many others left unmentioned. Programmable Logic Array (PLA) and Programmable Array Logic (PAL) are two examples of Simple Programmable Logic Devices (SPLD). Programmable Logic Array (PLA) consists of two planes of logic arrays: a programmable array of AND gates and a programmable array of OR gates. The AND plane and the OR plane give the possibility to computer any function expressed

as a sum of products. Every AND gate in AND plane is associated with inputs and complement of inputs to generate any product term. And, OR gate generates the sum of AND gate outputs. The example of PLA is shown in the figure below.

Example 2: Implement the following truth table using PLA

A	B	C	F1	F2	F3	F4
0	0	0	0	0	1	1
0	0	1	0	1	0	1
0	1	0	0	1	0	1
0	1	1	0	1	0	1
1	0	0	0	1	0	1
1	0	1	0	1	0	1
1	1	0	0	1	0	1
1	1	1	1	1	0	0

$$F1 = ABC$$

$$F2 = A + B + C$$

$$F3 = A'B'C'$$

$$F4 = A' + B' + C'$$

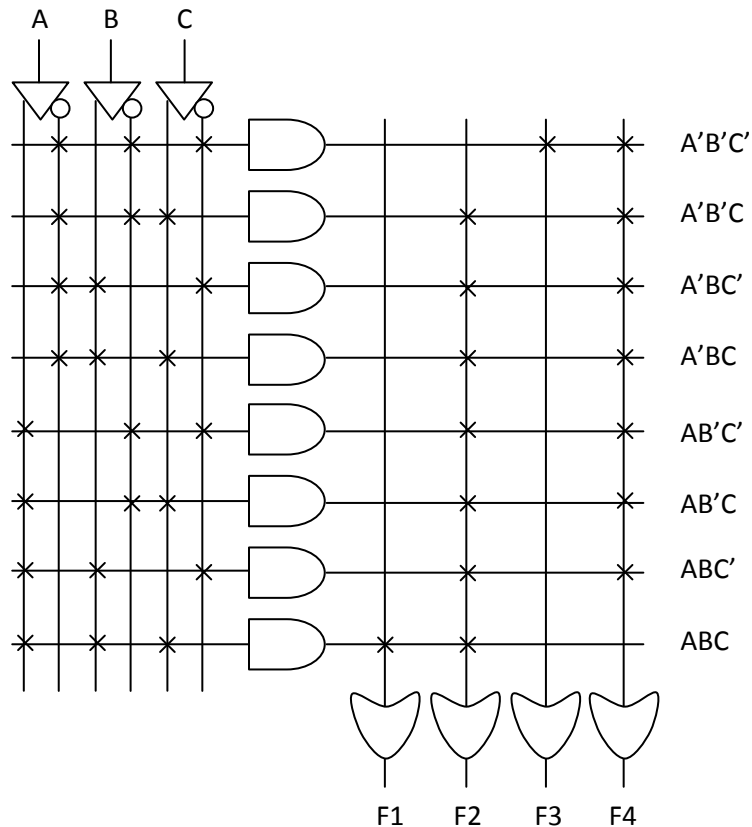


Figure 8.16: PLA implementation for given truth table

Programmable Array Logic uses just one programmable array: fixed OR matrix and programmable AND matrix. It decreases number of expensive programmable components which further reduces size and delay. PLA and PAL are generally used for low-complexity problems which require fairly high speed. As the complexity grows, Complex Programmable Logic Devices (CPLD) must be used. CPLDs are the integration of numerous SPLDs with added programmable interconnect between them. CPLD is a combination of fully programmable AND/OR array which perform a multitude of logic functions and microcells which perform combination or sequential logic. CPLDs may use analog sense amplifiers to boost the performance but at the cost of very high current requirements.

- **Intel 8051 Micro-controller family, its architecture and instruction sets**
- **Assembly language programming**
- **Interfacing with seven segment display**

9.1 Intel 8051 Micro-controller family, its architecture and instruction sets

A. Introduction

Microcontroller is a small computer on a single IC which contains processor core along with memory, I/O ports and other features. Microcontrollers are used in embedded applications in which systems are controlled automatically to carry out certain application. Almost every system using microcontroller performs control-oriented tasks. Several peripheral devices are inbuilt within the microcontroller to carry out the specified function. Timers, ADC and serial communication devices are few examples of peripheral devices.

B. Block Diagram

The general block diagram of the microcontroller is shown in the figure below:

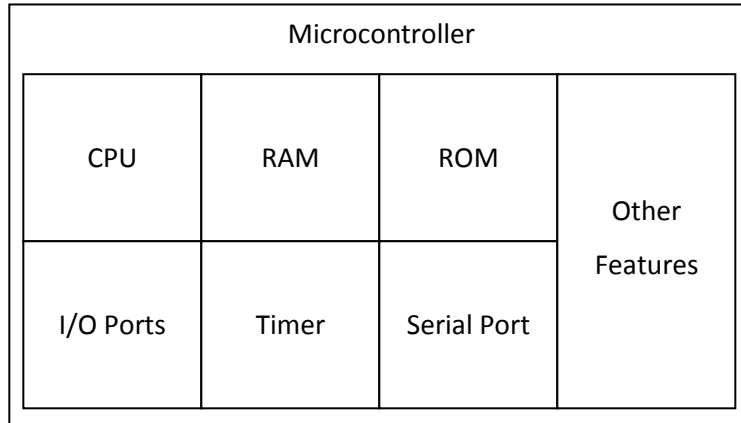


Figure 9.1: General block diagram of microcontroller

C. Comparison with Microprocessor

Many would easily presume that microcontroller and microprocessor to be similar. However, the following table will make a clear distinction between microcontroller and microprocessor.

SN	Microprocessor	Microcontroller
1.	General purpose processors	Special purpose processors
2.	It contains complete functional CPU only	In addition to functional CPU, it has timers, I/O ports, internal RAM and ROM, and other features
3.	Designer can select the size of memory, number of I/O ports, timers etc to be used	Size of memory, number of I/O ports, timers etc are fixed for a particular microcontroller

4.	Clock speed in very high in GHz range	Clock speed is low in MHz range
5.	Powerful addressing modes and many instructions are available to move data between memory and CPU	It focuses on bit handling instructions along with byte processing instructions.
6.	Access time for external memory and I/O devices is more	Access time for on-chip memory and I/O devices is less
7.	Microprocessor based systems are expensive and consumes more power	Microcontroller based systems are cheap and consumes less power

D. Criteria For Choosing a Microcontroller

- It must meet the computational needs of the task efficiently and cost effectively. Other considerations includes
 - Speed, packaging (DIP(dual line package), QFP(quad flat package)), power consumption, amount of RAM and ROM, number of I/O pins and the timer on the chip
 - Ease to amendments, cost per units
- It must provide flexibility to develop products around it. Some of the considerations include availability of an assembler, debugger, C compiler, emulator, technical support.
- It along with other reliable resources must be readily available in required quantities at any instant of time.

E. Comparison of 8051 Family Members

Each member of 8051, somehow, differs from each other. Though the instruction sets are almost common, the features provided can vary. The 8031 microcontroller is also referred as ROMLESS 8051 as all features are common except ROM space. The following table shows comparison of 8051, 8052 and 8031 microcontrollers.

Table 9.1: Comparison of three microcontrollers

Feature	8051	8052	8031
ROM	4K	8K	0K
RAM (bytes)	128	256	128
Timers	2	3	2
I/O Pins	32	32	32

Serial Port	1	1	1
Interrupt Sources	6	8	6

F. 8051 ARCHITECTURE

Internal Block Diagram of 8051

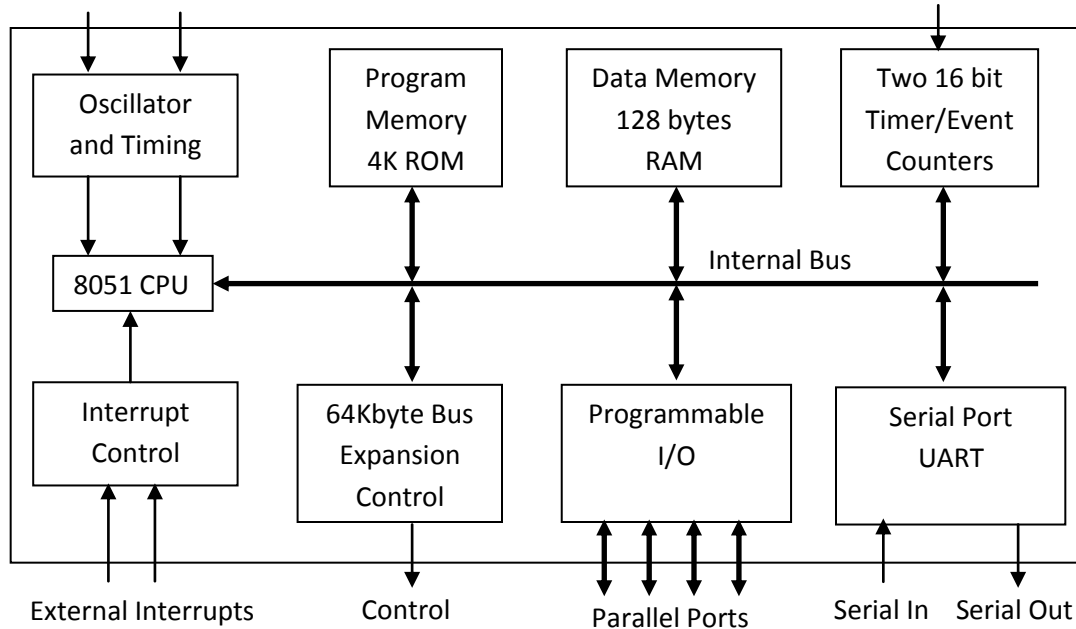


Figure 9.2: Internal block diagram of 8051 Architecture

Features of 8051 Architecture

- **Eight bit CPU with registers A and B:** Register A or Accumulator is used for mathematical and data transfer operations. Register B is used for multiplication and division purpose.
- **Sixteen bit program counter (PC) and data pointer (DPTR):** PC points to the address of next instruction to be executed from ROM while DPTR is used to point to the memory addresses for internal and external code access and external data access. DPTR is made up of two 8 bit registers, DPH and DPL.
- **Eight bit program status word (PSW):** Four flags in PSW are used to represent the outcomes of mathematical operations; Carry (CY), Auxiliary Carry (AC), Overflow (OV), and Parity (P) flags. It also consists two register select bits which select the particular register bank; RS1 and RS0 determines which register bank is being used out of four register banks.
- **Eight bit stack pointer (SP):** SP points to the stack which is the area to store and retrieve data quickly for some operations. It follows last in first out technique.

- **Internal ROM:** It consists of 4 Kbytes or memory space as program memory. Look up tables can also be stored which can be accessed using appropriate instruction.
- **Internal RAM:** It consists of 128 bytes of memory space as data memory.
 - Four Register Banks, each with eight registers (R0 – R7): Bank 0 occupies address from 00H to 07H and consecutive addresses are used by bank 1, bank 2 and bank 3. Total 32 registers are available from address 00H to 1FH. Bank 0 is selected as default. RS1 = 0 and RS0 = 1 in PSW register will select the register bank 1.
 - Sixteen bytes of bit addressable memory: The address from 20H to 2FH of RAM is bit addressable. It is useful in bit manipulating operations. Each bit can be addressed using direct address from 00H to 7FH.
 - Eighty bytes of general purpose data memory: The memory space from address 30H to 7FH can be used for various operations when required.

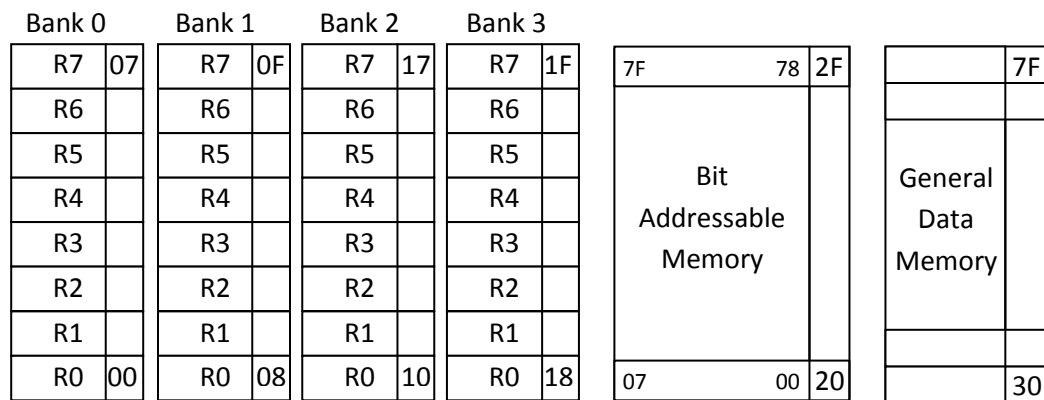


Figure 9.3: Internal RAM Organization

- **Thirty two I/O pins arranged as four 8 bit ports:** Four ports are bidirectional and can be used for input and output. Some of the pins are multifunctional which provide other functions along with input and output.
- **Two 16 bit timer/counters (T0 and T1):** Each counter can be programmed to count internal clock pulses, acting as a timer, or programmed to count external pulses as a counter. This selection as well as mode of operation of counter can be set by using Timer Mode Control (TMOD) register.
- **Full Duplex serial data receiver/transmitter:** Register Serial Control (SCON) controls serial data communication, and pins RXD and TXD are used to connect to other devices supporting serial

communication. The Serial Buffer (SBUF) register is used to hold data in serial communication process.

- **Two external interrupts and three internal interrupt sources:** Interrupt Enable (IE) register selects which interrupt is to be selected and enabled. INT0 and INT1 pins are used by external circuitry to interrupt processor. Timer overflow (TF), receive interrupt (RI), and transmit interrupt (TI) are internal interrupts.
- **The 8051 Oscillator and Clock:** It is the circuitry that generates the clock pulses by which all internal operations are synchronized. Time for particular instruction execution can be calculated based on number of machine cycles required by the instruction. For AT89S51, operating frequency is 11.0592 and its machine cycle consists of 12 clocks. So, time period for one machine cycle is 12 times the time period of single pulse which is equivalent to 1.085 μ s.

8051 Special Function Registers (SFRs)

They are a group of specific internal registers that use internal RAM and their address lie between 80H and FFH. Some SFRs are bit addressable which allows programmer to access each bit of the register. The list of SFRs is given in the table below:

Table 9.2: List of Special Function Registers

Name	Description	RAM Address	Access Level
A	Accumulator	0E0H	Bit Addressable
B	Register B, for multiplication and division	0F0H	Bit Addressable
PSW	Program Status Word	0D0H	Bit Addressable
SP	Stack Pointer	81H	
DPH	Data Pointer Higher Byte	83H	
DPL	Data Pointer Lower Byte	82H	
IE	Interrupt Enable	0A8H	Bit Addressable
IP	Interrupt Priority	0B8H	Bit Addressable
P0	Port 0	80H	Bit Addressable
P1	Port 1	90H	Bit Addressable
P2	Port 2	0A0H	Bit Addressable
P3	Port 3	0B0H	Bit Addressable
PCON	Power Control	87H	

SCON	Serial Port Control	98H	Bit Addressable
SBUF	Serial Port Data Buffer	99H	
TMOD	Timer/Counter Mode Control	89H	
TCON	Timer/Counter Control	88H	Bit Addressable
TL0	Timer 0 Low Byte	8AH	
TH0	Timer 0 High Byte	8CH	
TL1	Timer 1 Low Byte	8BH	
TH1	Timer 1 High Byte	8DH	

G. Pin Descriptions

Pins 1 – 8 (PORT 1)

These eight pins represent **PORT 1** which can be used as an input/output port. Since it is internally pulled up, it can be used without any external pull up registers configuration.

Pin – 9 (RESET)

RESET pin is used to set different registers to its initial values. The RESET pin must be set high for 2 machine cycles.

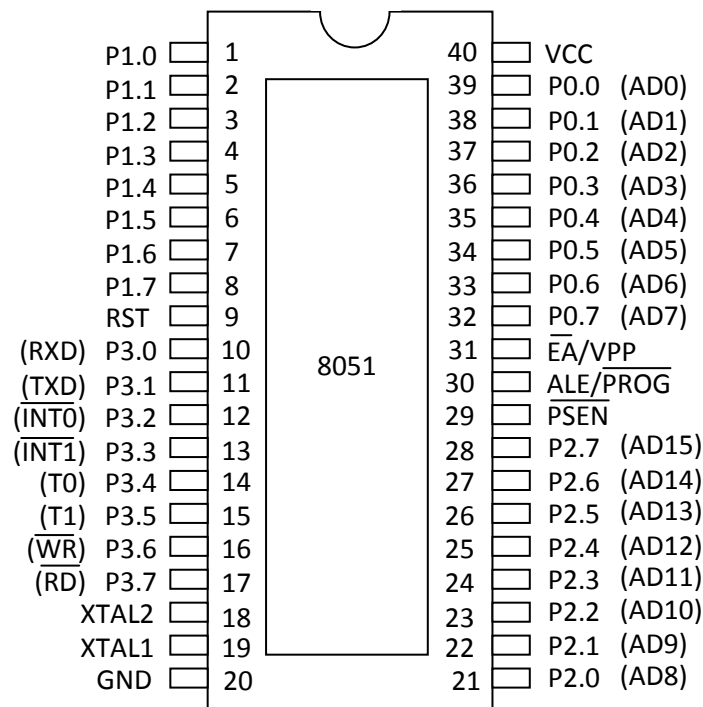


Figure 9.4: Pin descriptions of 8051 microcontroller

Pins 10 – 17 (PORT 3)

These pins together called **PORT 3** are bi directional and multifunctional in nature. Similar to port 1, it can be used as input or output without any external pull up registers configuration. Besides I/O, it supports serial communication (RXD and TXD), external interrupts (INT0 and INT1), timers (T0 and T1), and control signals (WR and RD) for external memory.

Pins 18 – 19 (XTAL)

These pins are used to connect an external crystal to provide system clock.

Pin – 20 (GND, 0V)**Pins 21 – 28 (PORT 2)**

These pins are bidirectional and multifunctional in nature. **PORT 2** may be used as an input or output port similar to port 1. The alternate use of port 2 is to provide a high-order address in conjunction with the port 0 to address external memory.

Pin 29 (PSEN)

Program Store Enable (PSEN) is connected to Output Enable (OE) pin of external memory being interfaced. It is an active low output signal. When this pin is reset, microcontroller can read content of external memory location.

Pin 30 (ALE)

Address Latch Enable (ALE) is used to select address or data signal that are required while interfacing external memory. It is active high output signal and when it goes high, the lower address provided by port 0 is latched into the external address latch. This pin is also the program pulse input during flash programming

Pin 31 (EA)

External Access enables or disables access of program from external memory. It must be connected to GND to fetch code from external program memory locations. It should be strapped to VCC for program executions of internal memory.

Pins 32 - 39 (PORT 0)

PORT 0 is a collection of open drain bidirectional I/O pins. It can be configured as low-order address or data bus while accessing external memory.

Pin 40 (VCC, +5V)

H. Minimum Hardware Configuration

Power Supply: Pin 40 is connected to +5VDC, Pin 20 is grounded. Pin 31 is connected to VCC, representing the code is accessed from internal memory.

Reset Circuit: Charging of capacitor makes RST high, which ensures two machine cycles on RST pin. After completion of charge, capacitor blocks DC causing RST low.

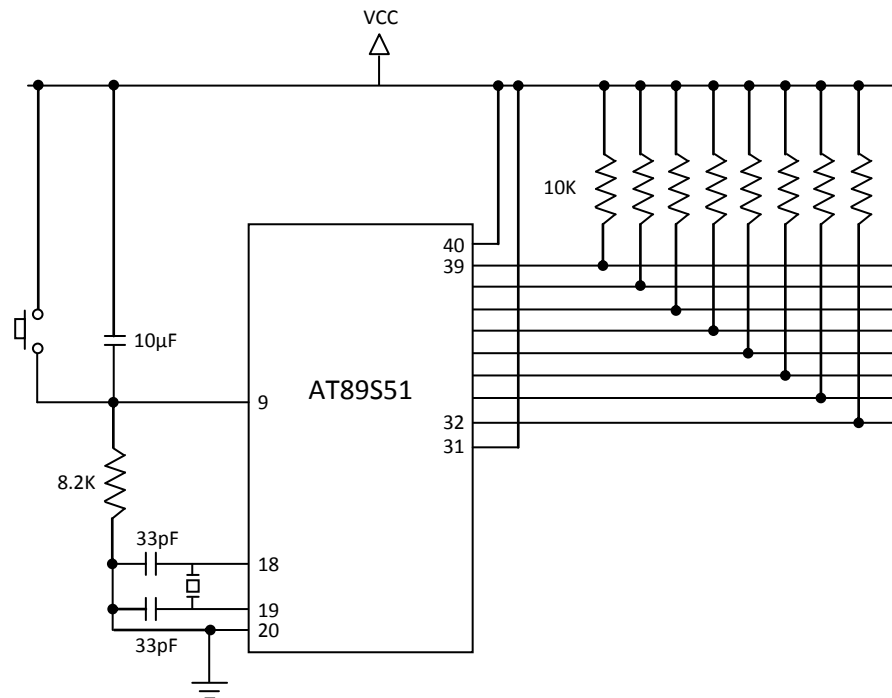


Figure 9.5: Minimum configuration for microcontroller to operate

Oscillator circuit: Ceramic capacitors of value between $20\mu\text{F}$ – $40\mu\text{F}$ are used as stabilizing capacitors. They act as loading capacitor and adjust the crystal frequency by shifting the frequency to a lower value.

Pull up Circuit: Pins of PORT 0 are open drain, so require pull up circuit. Each pin must be connected externally to a 10K ohm pull-up resistor.

I. 8051 Instruction Sets

Different symbols are used in the instruction whose meaning is clarified below:

#data – represents 8 bit data

Rn – represents one of eight registers (R0, R1, R2, R3, R4, R5, R6, R7)

@Ri – represents address pointed by value of Ri. Ri can be either R0 or R1.

direct – represents direct byte addressable memory

bit – direct bit addressable memory

C – Carry, A – Accumulator, B – Register B

addr11 – 11 bit address, addr16 – 16 bit address and rel – 8 bit relative address

1. Data Transfer Instructions

SN	Mnemonics	Operation	Description
1	MOV A, Rn	$A \leftarrow Rn$	MOV instruction is used to transfer data involving registers, memory and immediate data
2	MOV A, direct	$A \leftarrow [\text{direct}]$	
3	MOV A, @Ri	$A \leftarrow [Ri]$	
4	MOV A, #data	$A \leftarrow \text{data}$	
5	MOV Rn, A	$Rn \leftarrow A$	
6	MOV Rn, direct	$Rn \leftarrow [\text{direct}]$	
7	MOV Rn, #data	$Rn \leftarrow \text{data}$	
8	MOV direct, A	$[\text{direct}] \leftarrow A$	
9	MOV direct, Rn	$[\text{direct}] \leftarrow Rn$	
10	MOV direct, direct	$[\text{direct}] \leftarrow [\text{direct}]$	
11	MOV direct, @Ri	$[\text{direct}] \leftarrow [Ri]$	
12	MOV direct, #data	$[\text{direct}] \leftarrow \text{data}$	
13	MOV @Ri, A	$[Ri] \leftarrow A$	
14	MOV @Ri, direct	$[Ri] \leftarrow [\text{direct}]$	
15	MOV @Ri, #data	$[Ri] \leftarrow \text{data}$	
16	MOV DPTR, #data16	$DPTR \leftarrow \text{data16}$	MOVC is used to read data from code memory (ROM)
17	MOVC A, @A + DPTR	$A \leftarrow [A + DPTR]$	
18	MOVC A, @A + PC	$A \leftarrow [A + PC]$	
19	MOVX A, @Ri	$A \leftarrow [Ri]$	MOVX is used to move data to and from external RAM. R0, R1 and DPTR are used to hold address of RAM
20	MOVX A, @DPTR	$A \leftarrow [DPTR]$	
21	MOVX @Ri, A	$[Ri] \leftarrow A$	
22	MOVX @DPTR, A	$[DPTR] \leftarrow A$	
23	PUSH direct	$\text{Stack} \leftarrow [\text{direct}]$	
24	POP direct	$[\text{direct}] \leftarrow \text{Stack}$	
25	XCH A, Rn	$A \leftarrow Rn, Rn \leftarrow A$	

26	XCH A, direct	$A \leftarrow [\text{direct}], [\text{direct}] \leftarrow A$	
27	XCH A, @Ri	$A \leftarrow [Ri], [Ri] \leftarrow A$	
28	XCHD A, @Ri		

2. Arithmetic Instructions

SN	Mnemonics	Operation	Description
1	ADD A, Rn	$A \leftarrow A + Rn$	Accumulator is one of the sources as well as destination for every ADD and SUB instructions
2	ADD A, direct	$A \leftarrow A + [\text{direct}]$	
3	ADD A, @Ri	$A \leftarrow A + [Ri]$	
4	ADD A, #data	$A \leftarrow A + \text{data}$	
5	ADDC A, Rn	$A \leftarrow A + Rn + C$	
6	ADDC A, direct	$A \leftarrow A + [\text{direct}] + C$	
7	ADDC A, @Ri	$A \leftarrow A + [Ri] + C$	
8	ADDC A, #data	$A \leftarrow A + \text{data} + C$	
9	SUBB A, Rn	$A \leftarrow A - Rn - C$	
10	SUBB A, direct	$A \leftarrow A - [\text{direct}] - C$	
11	SUBB A, @Ri	$A \leftarrow A - [Ri] - C$	
12	SUBB A, #data	$A \leftarrow A - \text{data} - C$	
13	INC A	$A \leftarrow A + 1$	
14	INC Rn	$Rn \leftarrow Rn + 1$	
15	INC direct	$[\text{direct}] \leftarrow [\text{direct}] + 1$	
16	INC @Ri	$[Ri] \leftarrow [Ri] + 1$	
17	DEC A	$A \leftarrow A - 1$	
18	DEC Rn	$Rn \leftarrow Rn - 1$	
19	DEC direct	$[\text{direct}] \leftarrow [\text{direct}] - 1$	
20	DEC @Ri	$[Ri] \leftarrow [Ri] - 1$	
21	INC DPTR	$DPTR \leftarrow DPTR + 1$	
22	MUL AB	$A \leftarrow \text{Lower Byte}$ $B \leftarrow \text{Higher byte}$	
23	DIV AB	$A \leftarrow \text{Quotient}$ $B \leftarrow \text{Remainder}$	
24	DA A	Decimal Adjust accumulator	

3. Logical Instructions

SN	Mnemonics	Operation	Description
1	ANL A, Rn	$A \leftarrow A \text{ AND } Rn$	Logical AND, OR and XOR allows direct operation on memory address as well.
2	ANL A, direct	$A \leftarrow A \text{ AND } [\text{direct}]$	
3	ANL A, @Ri	$A \leftarrow A \text{ AND } [Ri]$	
4	ANL A, #data	$A \leftarrow A \text{ AND } \text{data}$	
5	ANL direct, A	$[\text{direct}] \leftarrow [\text{direct}] \text{ AND } A$	
6	ANL direct, #data	$[\text{direct}] \leftarrow [\text{direct}] \text{ AND } \text{data}$	
7	ORL A, Rn	$A \leftarrow A \text{ OR } Rn$	
8	ORL A, direct	$A \leftarrow A \text{ OR } [\text{direct}]$	
9	ORL A, @Ri	$A \leftarrow A \text{ OR } [Ri]$	
10	ORL A, #data	$A \leftarrow A \text{ OR } \text{data}$	
11	ORL direct, A	$[\text{direct}] \leftarrow [\text{direct}] \text{ OR } A$	
12	ORL direct, #data	$[\text{direct}] \leftarrow [\text{direct}] \text{ OR } \text{data}$	
13	XRL A, Rn	$A \leftarrow A \text{ XOR } Rn$	
14	XRL A, direct	$A \leftarrow A \text{ XOR } [\text{direct}]$	
15	XRL A, @Ri	$A \leftarrow A \text{ XOR } [Ri]$	
16	XRL A, #data	$A \leftarrow A \text{ XOR } \text{data}$	
17	XRL direct, A	$[\text{direct}] \leftarrow [\text{direct}] \text{ XOR } A$	
18	XRL direct, #data	$[\text{direct}] \leftarrow [\text{direct}] \text{ XOR } \text{data}$	
19	CLR A	$A \leftarrow 0$	
20	CPL A	$A \leftarrow A'$	
21	RL A	Rotate A left	
22	RLC A	Rotate A left through C	
23	RR A	Rotate A right	
24	RRC A	Rotate A right through C	
25	SWAP A	Swap nibbles of A	

4. Bit Manipulation and Program Branching Instructions

SN	Mnemonics	Operation	Description
----	-----------	-----------	-------------

1	CLR C	$C \leftarrow 0$	
2	CLR bit	$\text{bit} \leftarrow 0$	
3	SETB C	$C \leftarrow 1$	
4	SETB bit	$\text{bit} \leftarrow 1$	
5	CPL C	$C \leftarrow C'$	
6	CPL bit	$\text{bit} \leftarrow \text{bit}'$	
7	ANL C, bit	$C \leftarrow C \text{ AND bit}$	
8	ANL C, /bit	$C \leftarrow C \text{ AND bit}'$	
9	ORL C, bit	$C \leftarrow C \text{ OR bit}$	
10	ORL C, /bit	$C \leftarrow C \text{ OR bit}'$	
11	MOV C, bit	$C \leftarrow \text{bit}$	
12	MOV bit, C	$\text{bit} \leftarrow C$	Short jumps must be within
13	JC rel	Jump if $C \leftarrow 1$	128 to +127 bytes of the
14	JNC rel	Jump if $C \leftarrow 0$	contents of PC
15	JB bit, rel	Jump if $\text{bit} \leftarrow 1$	
16	JNB bit, rel	Jump if $\text{bit} \leftarrow 0$	Long Jumps and calls can be
17	JBC bit, rel	Jump if $\text{bit} \leftarrow 1$, and $b \leftarrow 0$	used for any location within
18	ACALL addr11	Absolute jump to routine	64 Kbyte address space
19	LCALL addr16	Long jump to routine	
20	RET	Return from subroutine	Absolute jumps and calls can
21	RETI	Return from interrupt	be used for address within
22	AJUMP addr11	Absolute jump	2Kbyte range
23	LJUMP addr16	Long jump	
24	SJMP rel	Short jump	
25	JMP @ A + DPTR	Jump relative to DPTR	
26	JZ rel	Jump if A is zero	
27	JNZ rel	Jump if A is not zero	
28	CJNE A, direct, rel		
29	CJNE A, #data, rel	Compare and jump if not	
30	CJNE Rn, #data, rel	equal	
31	CJNE @Ri, #data, rel		

32	DJNZ Rn, rel	Decrease and jump if not	
33	DJNZ direct, rel	zero	

J. Addressing Modes in 8051

Immediate Addressing Mode

The source operand is a constant value which must be preceded by # sign. It is used to load direct values into registers. For example, MOV A, #25H will assign 25 to register A.

Register Direct Addressing Mode

The operand is a register which holds the data to be manipulated. For example, ADD A, R5 will add content of A and R5, and store back in A.

Register Indirect Addressing Mode

Register is used to point the effective address of the operand. Registers R0, R1 and DPTR are used as pointer registers which must be preceded by @ sign. For example, MOV A, @R0 represents copying the contents of the address in R0 to the accumulator.

Direct Addressing Mode

The operand represents the actual address of RAM in the instruction. For instance, MOV A, 80H moves the data of 80H into accumulator.

Relative Addressing

A relative address or offset is added to the PC to form the actual address. Generally used in jump instructions.

Absolute Addressing Mode

In instructions, 11-bit or 16-bit absolute address is specified as the operand. ACALL and AJMP instructions use 11-bit address while LCALL and LJMP use 16-bit address.

Indexed Addressing Mode

Index value or displacement is added to the base address to generate the effective address of the operand. For instance, MOVC A, @A + DPTR uses indexed addressing mode. The content pointed by (A + DPTR) address of ROM is copied to accumulator.

9.2 Assembly language programming

An assembly language program consists of series of statements which include assembly language instructions and directives. Assembly language instruction represents the operation to be carried out by the processor. Every instruction is composed of mnemonic followed by one, two or no

operand. Mnemonic represents the actual operation to be done which operands are data items being manipulated. Directives are used to give directions to the assembler. Generally used directives are DB, ORG, END, and EQU. The DB directive is used to define 8-bit data. The ORG represents the beginning of the program address while END represents the end of program. The EQU directive is used to define constant within a program. The numbers used must be followed by H to represent hex value otherwise the value will be taken as decimal.

The assembly language program, in general, is written using following format:

```
[label:]      Mnemonic    [operands]    [; comments]
```

Example 1: Read the content of port1 and port2, OR those contents and store the result in external RAM location 0310.

Problem Analysis: Port1 (with address 90H) and Port2 (with address A0H) provide eight bit data, so after OR operation the final result will also be of eight bit. Hence, single byte memory location is enough to store the result. However, to store the result in external RAM, the address of external RAM must be loaded into DPTR register and MOVX instruction should be used for data transfer.

Source Code:

```
ORG 00H
MOV A, 90H           ; copy the data of port1 to A
ORL A, 0A0H          ; OR the contents of A with port2, and stored in A
MOV DPTR, #0310H     ; DPTR used to point the external RAM address
MOVBX @DPTR, A       ; move the content of A to RAM location pointed by DPTR
END
```

Example 2: Read the content of internal RAM locations 27H and 28H, add them and store the result in RAM locations 30H and 31H.

Problem Analysis: The largest possible value at memory locations can be FFH, so the maximum value of final result will be $FFH + FFH = 01FEH$. Hence, two bytes of memory is required to store the result. To solve this, ADD instruction must be used to add two data while ADDC or JNC/JC can be used to add the carry.

Source Code:

```
ORG 00H
MOV 30H, #00H        ; [30] ← 00, assigns zero to memory location 30H
MOV A, 27H           ; A ← [27], assigns content of location 27H to accumulator
```

```

ADD A, 28H          ;  $A \leftarrow A + [28]$ , adds content of A and memory location 28H
MOV 31H, A          ;  $[31] \leftarrow A$ , moves the lower byte or result to 31H address of RAM
MOV A, #00H         ;  $A \leftarrow 00$ , reset accumulator
ADDC A, 30H         ;  $A \leftarrow A + [30] + C$ , to extract the value of carry
MOV 30H, A          ;  $[30] \leftarrow A$ , move upper byte to location 30.
END

```

Example 3: Add the content of internal RAM location 29H and port1, and store the result in RAM locations 30H and 31H in BCD form.

Problem Analysis: In BCD, the largest possible value can be 99, so the maximum value of final result will be $99 + 99 = 198$. Hence, two bytes of memory is required to store the result. To solve this, ADD instruction must be used to add two BCD data and DA instruction after addition. However, DA instruction is not required for upper byte of result as it is less than 10. But, had there been more numbers causing upper byte to exceed more than 9, DA would have been required for upper byte as well.

Source Code:

```

ORG 00H
MOV 30H, #00H       ;  $[30] \leftarrow 00$ , assigns zero to memory location 30H
MOV A, 90H          ;  $A \leftarrow [90]$ , assigns content of Port1 to accumulator
ADD A, 29H          ;  $A \leftarrow A + [29]$ , adds content of accumulator and location 29H
DA A                ; Adjust the content of accumulator to BCD form
MOV 31H, A          ;  $[31] \leftarrow A$ , moves the lower byte or result to 31H address of RAM
MOV A, #00H         ;  $A \leftarrow 00$ , reset accumulator
ADDC A, 30H         ;  $A \leftarrow A + [30] + C$ , to extract the value of carry
MOV 30H, A          ;  $[30] \leftarrow A$ , move upper byte to location 30.
END

```

Example 4: Add 10 bytes of data of RAM location starting from address 20H. Store lower byte at 30H and upper byte at 31H.

Problem Analysis: Use of direct RAM address can make program complex, so R0 or R1 can be used to point the one byte address of RAM location. The register (R0 or R1) then can be incremented to access data of consecutive locations. The final sum after addition of 10 bytes of data can result in

two bytes of data. So, carry must be checked after addition and value in register must be incremented accordingly when carry is resulted after addition.

Source Code:

```

ORG 00H
MOV R0, #20H      ; R0 ← 20H, assigning starting address of RAM to R0
MOV R5, #0AH      ; R5 ← 0AH, counter for 10 bytes of data
MOV R6, #00H      ; R6 ← 00H, used for lower byte result
MOV R7, #00H      ; R7 ← 00H, used for upper byte result

```

HOME:

```

MOV A, @R0        ; A ← [R0], move content of RAM location pointed by R0 to A
ADD A, R6          ; A ← A + R6, adds data of RAM location in each iteration
JNC NEXT          ; checking if carry is generated after addition
    INR R7         ; R7 ← R7 + 1, carry after each addition is added to form upper byte

```

NEXT:

```

MOV R6, A          ; R6 ← A, move partial sum to R6
INR R0             ; R0 ← R0 + 1, increment the address of RAM to access next byte
DJNZ R5, HOME      ; decrease R5 by 1 and jump to HOME if R5 is not equal to zero
MOV 30H, R6        ; [30] ← R6, move the lower byte to memory address 30H
MOV 31H, R7        ; [31] ← R7, move the upper byte to memory address 31H
END

```

Example 5: ASCII character string is stored in the program memory starting at 300H. Send each character to port 2.

Problem Analysis: Since the data is stored in the program memory (ROM), DPTR must be used to point the address of ROM. Using MOVC instruction, the data can be retrieved and further manipulated as required.

Source Code:

```

ORG 00H
MOV DPTR, #300H    ; load address of data into DPTR

```

HOME:

```

CLR A              ; A ← 0, clear the content of A
MOVC A, @A + DPTR  ; A ← [A + DPTR], load ROM content of address (A + DPTR) to A

```

```

JZ EXIT          ; jump out of loop if last character is detected which is 0
MOV P2, A        ; P2 ← A, move content of accumulator to port 2
INC DPTR         ; DPTR ← DPTR + 1, to point to next character of ROM
SJMP HOME        ; Continue the loop

ORG 300H

DB    "ASSEMBLY PROGRAM", 0

EXIT:
NOP
END

```

9.3 Interfacing with seven segment display

A. Seven Segment Configurations

Pin Configurations

The figure shows the pin configuration of seven segment display which consists of 10 pins; eight pins to control the leds and two common pins which are grounded or connected to VCC based on common cathode or common anode configuration.

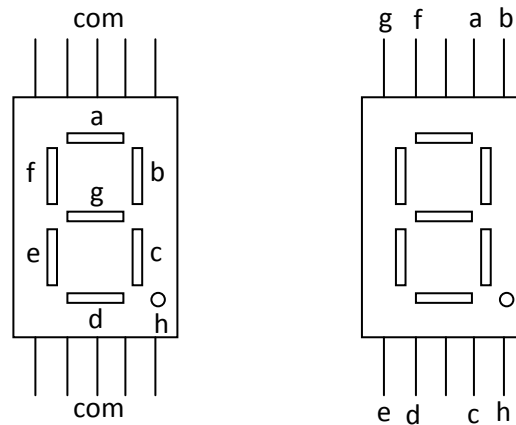


Figure 9.6: Pin Configurations of Seven Segment Display

Modes of Configurations

There are two modes of configurations: Common Anode Configuration and Common Cathode Configuration. In common anode configuration, anodes of all leds are connected together to form a common pin which must be connected to high logic voltage. In common cathode configuration, cathodes of all leds are connected together to form a common pin which must be connected to low logic voltage.

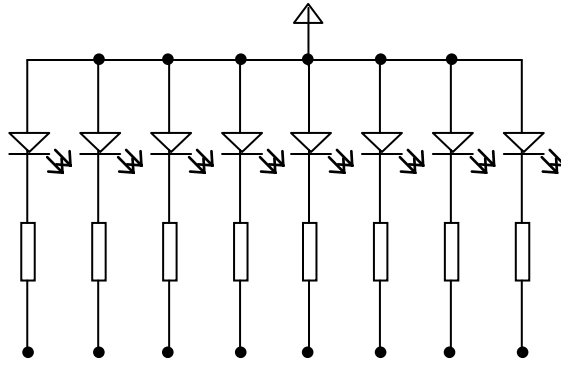


Figure 9.7: Common Anode Configurations

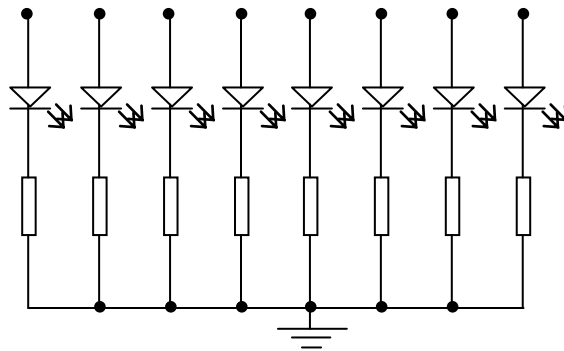


Figure 9.8: Common Cathode Configurations

Lookup table of HEX equivalent

For common anode configurations, low logic must be provided to the pins of the seven segment display to glow the particular led. The equivalent hex values are sent through the port of microcontroller. However, designer must be aware of the driving circuit which should provide the equivalent hex to the pins of seven segment display.

Common Anode Configurations									For Common Cathode mode
Digits	Individual LEDs Illuminated								HEX
	h	g	f	e	d	C	b	a	
0	1	1	0	0	0	0	0	0	0xC0
1	1	1	1	1	1	0	0	1	0xF9
2	1	0	1	0	0	1	0	0	0xA4
3	1	0	1	1	0	0	0	0	0xB0
4	1	0	0	1	1	0	0	1	0x99
5	1	0	0	1	0	0	1	0	0x92

6	1	0	0	0	0	0	1	0	0x82	0x7D
7	1	1	1	1	1	0	0	0	0xF8	0x07
8	1	0	0	0	0	0	0	0	0x80	0x7F
9	1	0	0	1	0	0	0	0	0x90	0x6F

B. Interfacing Seven Segment

Before connecting the seven segment display to the port of microcontroller, the current requirement of the seven segment display along with the source current and sink current capacity of the microcontroller must be examined. However, it is always better to use the driving circuit rather to connect seven segment display directly to the port of microcontroller.

Hardware Connections

Any port of the microcontroller can be used to connect to the seven segment display through the driving circuit. The driving circuit may be in the form of an IC which can sink or source high current. The circuit configurations can vary depending on the designer. IC ULN2003 is an example of a current sinker while IC L293D can be a good source of current for driving circuits.

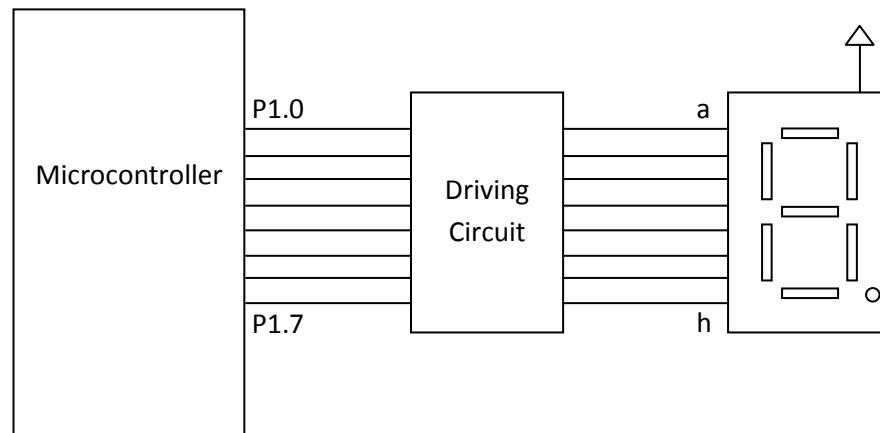


Figure 9.9: Connection of common anode seven segment with microcontroller

Coding Implementations

In Assembly language, simple MOV instructions can be used to transfer HEX value to the seven segment display. For example, `MOV P2, #92` will display the digit 5 for common anode configurations. Using C programming language, simple assignment can be done. For example, `P2 =`

0x92; will display digit 5. However, appropriate delay or repetition mechanism must be used to ensure that the digit displays for certain duration and becomes observable to the designer or user.

Delay Calculation in assembly program

Actual time of the execution of instructions can be determined by making use of the operating frequency of the microcontroller and machine cycles required by the instructions. The total machine cycles required by the instruction is multiplied by time duration of one machine cycle to calculate the total time.

DELAY:

```

        MOV R4,#64H          ; MC = 1, executes only once
AGAN:   MOV R3,#0AH          ; MC = 1, executes 100 times
AGA:    DJNZ R3,AGA          ; MC = 2, executes 100 x 10 = 1000 times
        DJNZ R4,AGAN         ; MC = 2, executes 100 times
        RET                  ; MC = 2, executes 1 time

```

In 8051, crystal frequency is 11.0592MHz and one machine cycle (MC) is equal to 12 clock cycles.

So, 1 MC = 1.085µs.

Total machine cycles in DELAY subroutine = 1 + 1x100 + 2x1000 + 2x100 + 2 = 2303

Total time duration = 2303 x 1.085µs = 2.5ms

FEW SOLVED EXAMPLES

1. Write an assembly and C language program to generate a pulse of 50% duty cycle at pin P2.3 of 8051 microcontroller.

➔ **Problem Analysis:** Duty cycle of 50% represents equal ON and OFF time at pin P2.3, so an arbitrary delay is required after setting P2.3 and after resetting P2.3.

Source Code:

```

ORG 00H
CLR P2.3
Back:
    CLR P2.3
    LCALL DELAY
    SETB P2.3
    LCALL DELAY
    SJMP Back

    ORG 300H
DELAY:
    MOV R5,#64H
AGAIN: MOV R4,#0FFH
AGAN:  MOV R3,#08H
AGA:   DJNZ R3,AGA
        DJNZ R4,AGAN
        DJNZ R5,AGAIN
    RET

END

```

```

In C programming language

#include<at89x52.h>
void delay(unsigned char x)
{
    int i, j;
    for(i=0;i<x;i++)
        for(j=0;j<1275;j++);
}

void main()
{
    P2_3 = 0;
    while(1)
    {
        P2_3 = 1;
        delay(50);
        P2_3 = 0;
        delay(50);
    }
}

```

2. **Control the LED connected at 2.1 by a SWITCH which is connected to P1.3. ON/OFF status of LED is defined by ON/OFF status of SWITCH.**

➔ **Problem Analysis:** This is a simple data movement problem in which a bit from one pin P1.3 of microcontroller is assigned to another pin P2.1. The code can vary based of LED configuration and SWITCH configuration used. Two cases are given below:

CASE I: The SWITCH will generate high logic when pressed and LED will glow when high logic is assigned to pin 2.1.

Source Code:

```
ORG 00H
CLR P2.1
SETB P1.3
Back:
    MOV C, P1.3
    MOV P2.1, C
    SJMP Back
END
```

```
#include<at89x52.h>
void main()
{
    P2_3 = 0;
    P1_3 = 1;
    while(1)
    {
        P2_1 = P1_3;
    }
}
```

CASE II: The SWITCH will generate low logic when pressed and LED will glow when high logic is assigned to pin 2.1.

Source Code:

```
ORG 00H
CLR P2.1
SETB P1.3
Back:
    MOV C, P1.3
    CPL C
    MOV P2.1, C
    SJMP Back
END
```

```
#include<at89x52.h>
void main()
{
    P2_3 = 0;
    P1_3 = 1;
    while(1)
    {
        P2_3 = !P1_3;
    }
}
```

3. Using an assembly and C language program, generate a pulse of 75% duty cycle at pin P1.7 when the switch connected to P1.1 is ON.

➔ **Problem Analysis:** Duty cycle of 75% represents ON time three times more than OFF time at pin P1.7. The status of P1.1 must be checked continuously. Based on the logic level at P1.1 after button is pressed, delay in case of ON time must be three times more than that of OFF time. In the code below, we assume the switch connects P1.1 to ground when pressed and P1.1 connected to VCC when not pressed. So, a low logic at P1.1 will generate the required pulse using appropriate delay.

Source Code:

```

                ORG 00H
                CLR P1.7
                SETB P1.1
BACK:
                MOV C, P1.1
                JC BACK
                SETB P1.7
                LCALL DELAY
                LCALL DELAY
                LCALL DELAY
                CLR P1.7
                LCALL DELAY
                SJMP BACK

                ORG 300H
DELAY:
                MOV R5, #64H
AGAIN:
                MOV R4, #0FFH
AGAN:
                MOV R3, #08H
AGA:
                DJNZ R3, AGA
                DJNZ R4, AGAN
                DJNZ R5, AGAIN
                RET

END

```

```

#include<at89x52.h>
#define OUT P1_7
#define SW P1_1
void delay(unsigned int x)
{
    int i,j;
    for(i=0;i<x;i++)
        for(j=0;j<1275;j++);
}
void main()
{
    OUT = 0;
    SW = 1;
    while(1)
    {
        if(SW == 0)
        {
            OUT = 1;
            delay(300);
            OUT = 0;
            delay(100);
        }
    }
}

```

4. Write an assembly and C language program to generate a count from 0 to 9 using a seven segment display. Use Common Cathode configurations.

➔ **Problem Analysis:** The equivalent HEX values are directly assigned one by one to the required port. Certain delay after each digit display can be placed to control the speed of count.

Source Code:

```

ORG 00H
    MOV P2, #00H
    MOV P2, #7FH
BACK:    LCALL DELAY
    MOV P2, #3FH
    MOV P2, #6FH
    LCALL DELAY
    LCALL DELAY
    MOV P2, #06H
    SJMP BACK
    LCALL DELAY
    MOV P2, #5BH
    ORG 300H
    LCALL DELAY
    MOV P2, #4FH
    DELAY:
    MOV R5, #64H
    LCALL DELAY
    AGAIN:    MOV R4, #0FFH
    MOV P2, #66H
    AGAIN:    MOV R3, #08H
    LCALL DELAY
    AGA:      DJNZ R3, AGA
    MOV P2, #6DH
    DJNZ R4, AGAN
    LCALL DELAY
    DJNZ R5, AGAIN
    MOV P2, #7DH
    RET
    LCALL DELAY
    MOV P2, #07H
    END
    LCALL DELAY

```

ALTERNATIVE CODE:

Problem Analysis: The HEX values are not directly assigned rather stored in memory and accessed using data pointer (DPTR). DPTR is used to point to the HEX values and MOVC instruction must be used to access the HEX values. MOVC uses accumulator to represent offset as well as data. The value of accumulator is added to DPTR to represent the address of memory and finally the data is stored into accumulator.

Source Code:

```

ORG 00H

    MOV P2, #00H

    MOV R6, #00H

    MOV DPTR, #DIGITS

MAIN:

    MOV A, R6
    MOVC A, @A+DPTR
    MOV P2, A
    LCALL DELAY
    INC R6
    CJNE R6, #0AH, MAIN

    MOV R6, #00H
    SJMP MAIN

DELAY:

    MOV R3, #0F0H

DEL1: MOV R2, #0FAH
DEL2: DJNZ R2, DEL2
      DJNZ R3, DEL1

RET

DIGITS:

    DB 3FH, 06H, 5BH, 4FH, 66H
    DB 6DH, 7DH, 07H, 7FH, 6FH

END

```

```

#include<at89x52.h>

#define DISPLAY P2

void delay(unsigned int x)
{
    unsigned int i,j;
    for(i=0;i<x;i++)
        for(j=0;j<1275;j++);
}

char digits[] = {0x3F, 0x06, 0x5B, 0x4F,
0x66, 0x6D, 0x7D, 0x07, 0x7F, 0x6F};

void main()
{
    unsigned char i;
    DISPLAY = 0x00;
    while(1)
    {
        for(i=0;i<10;i++)
        {
            DISPLAY = digits[i];
            delay(100);
        }
    }
}

```


5. A PUSH BUTTON is connected to P1.1, increase the count in SEVEN SEGMENT when the button is pressed.

➔ Source Code:

```

ORG 00H
    SETB P1.1
    MOV P2, #00H
    MOV R6, #00H
    MOV DPTR, #DIGITS
MAIN:
    MOV A, R6
    MOVC A, @A+DPTR
    MOV P2, A
    MOV C, P1.1
    JC MAIN
    LCALL DELAY
    INC R6
    CJNE R6, #0AH, MAIN

    MOV R6, #00H
    SJMP MAIN
DELAY:
    MOV R3, #0F0H
DEL1:  MOV R2, #0FAH
DEL2:  DJNZ R2, DEL2
       DJNZ R3, DEL1
RET
DIGITS:
    DB 3FH, 06H, 5BH, 4FH, 66H
    DB 6DH, 7DH, 07H, 7FH, 6FH
END

```

```

#include<at89x52.h>
#define DISPLAY P2
#define SW P1_1
void delay(unsigned int x)
{
    unsigned int i,j;
    for(i=0;i<x;i++)
        for(j=0;j<1275;j++);
}
char digits[] = {0x3F, 0x06, 0x5B, 0x4F,
0x66, 0x6D, 0x7D, 0x07, 0x7F, 0x6F};
void main()
{
    unsigned char i = 0;
    DISPLAY = 0x00;
    while(1)
    {
        DISPLAY = digits[i];
        if(SW == 0)
        {
            i++;
            if (i>9)
                i = 0;
            delay(20);
        }
    }
}

```

6. Using two seven segment displays, build a down counter which counts from 99 to 00. Make appropriate assumptions wherever necessary.

➔ **Description:** Two ports can be used to display two digits. However, it is better to use a single port for both seven segment displays and control their display through control lines. Two digits separately are sent at different instant of time. Corresponding seven segment display must be enabled while the digits are sent at different instant. However, the time duration between sending of lower digit and upper digit must not be high. High delay can lead to flickering effect which causes both digits to be displayed one by one rather than simultaneously. Also the speed of count can be controlled by using appropriate repetition of each display of both digits.

Source Code:

```

ORG 00H

    MOV P2, #00H
    MOV R6, #09H           ; counter for lower byte
    MOV R7, #09H           ; counter for upper byte
    MOV R5, #07h           ; to control speed of counter
    MOV P3, #00H
    MOV DPTR, #LABEL1      ; loads the starting address of hex code list to DPTR

MAIN:
    MOV A, R6
    SETB P3.0              ; activates 2nd display to display lower byte
    CLR P3.1               ; deactivates the 1st display
    LCALL DISPLAY
    LCALL DELAY

    MOV A, R7
    SETB P3.1              ; activates 1st display
    CLR P3.0               ; deactivates the 2nd display
    LCALL DISPLAY
    LCALL DELAY

    DJNZ R5, MAIN          ; repetition of same display to control speed
    MOV R5, #07H

```

```

    DEC R6                                ; decrease value of R6
    CJNE R6, #-1, MAIN                    ; compare unless R6 becomes less than zero

    MOV R6, #09H
    DEC R7
    CJNE R7, #-1, MAIN

    MOV R7, #09H
    SJMP MAIN

DISPLAY:
    MOVC A, @A+DPTR                      ; load HEX value to accumulator from memory
    MOV P2, A                            ; sending HEX value to SEVEN SEGMENT through P2
    RET

DELAY:
    MOV R3, #F0H

DEL1:  MOV R2, #0FAH
DEL2:  DJNZ R2, DEL2
       DJNZ R3, DEL1
    RET

LABEL1:                                ; represents starting address of HEX value list
    DB 3FH, 06H, 5BH, 4FH, 66H, 6DH, 7DH, 07H, 7FH, 6FH

END

```

Equivalent code in C language

Description: Within infinite loop, the first loop is used to select the upper byte, the second loop is used to select the lower byte and the third loop is used to control the speed of count.

Source Code:

```

#include<at89x52.h>
#define DISPLAY P2
#define SEL0 P3.1
#define SEL1 P3.0

```

```

void delay(unsigned int x)
{
    unsigned int i,j;
    for(i=0;i<x;i++)
        for(j=0;j<1275;j++);
}

char digits[] = {0x3F, 0x06, 0x5B, 0x4F, 0x66, 0x6D, 0x7D, 0x07, 0x7F, 0x6F};

void main()
{
    char i, j, k;
    DISPLAY = 0x00;
    while(1)
    {
        for(i = 9; i >= 0; i++)                // loop for upper byte
        {
            for(j = 9; j >= 0; j--)            // loop for lower byte
            {
                for(k = 0; k<7; k++)           // loop to control the speed of count
                {
                    DISPLAY = digits[j];
                    SEL0 = 0;
                    SEL1 =1;
                    delay(5);
                    DISPLAY = digits[i];
                    SEL0 = 1;
                    SEL1 = 0;
                    delay(5);
                }
            }
        }
    }
}

```

- **Introduction**
- **VHDL Code Structure**
- **Data types, Data Objects and Operators**
- **Statements in VHDL**
- **Standard Architectures**
- **FSM Design**

10.1 Introduction

VHDL is a hardware description language. It is used to describe the behavior of an electronic system, which further enables designer to implement the physical system. VHDL stands for VHSIC Hardware Description Language, where VHSIC is an acronym for Very High Speed Integrated Circuits.

The main purpose of VHDL is to model and synthesize digital circuits. Simulation and testing of the design for the optimum operation can be done using VHDL model of the system. Also, digital integrated circuits for particular operations can be created using VHDL or other hardware description languages. Finally, VHDL code can be used to create actual functional system. Hence, VHDL code can be used either to implement the circuit in a programmable device or can be forwarded for fabrication.

VHDL Invariants

- It is not case sensitive.
- It is not sensitive to white space.
- Comments begin with two consecutive dashes (“--”).
- Parenthesis usage is optional in many cases.
- Every statement in VHDL is terminated with a semicolon.
- Statements are inherently concurrent. Only statements placed inside a PROCESS, FUNCTION, or PROCEDURES are executed sequentially.

10.2 VHDL Code Structure

Fundamental VHDL Units

VHDL code comprises of at least the three fundamental sections: LIBRARY Declaration, ENTITY and ARCHITECTURE. LIBRARY is a collection of pre-defined set of codes that can be re-used or shared by various designs. ENTITY specifies the I/O connections of the system. ARCHITECTURE contains the code that describes how the circuit should function.

LIBRARY DECLARATION

The general form is:

```
LIBRARY LIBRARY_NAME;
```

```
USE LIBRARY_NAME.PACKAGE_NAME.PACKAGE_PARTS;
```

Example:

```
LIBRARY IEEE;
```

```
USE IEEE.STD_LOGIC_1164.ALL;
```

The libraries *STD* and *WORK* are made visible by default, so they are not required to declare. However, *STD_LOGIC_1164* package of *IEEE* library must be declared when *STD_LOGIC* data type is used in the design. Similarly, for *SIGNED* and *UNSIGNED* data types and its related arithmetic and comparison operations, package *STD_LOGIC_ARITH* of LIBRARY *IEEE* must be declared.

ENTITY

The VHDL ENTITY declaration describes the interface or the external representation of the circuit. An ENTITY is a list of all input and output pins with its specification such as data type and data direction mode.

Its syntax is:

```
ENTITY ENTITY_NAME IS
  PORT (
    PORT_NAME: SIGNAL_MODE SIGNAL_TYPE;
    PORT_NAME: SIGNAL_MODE SIGNAL_TYPE;
    ...
  );
END ENTITY_NAME;
```

Here, ENTITY_NAME and PORT_NAME are identifiers. The SIGNAL_MODE which defines the direction of signal can be IN, OUT, INOUT, or BUFFER. IN and OUT are unidirectional pins, while INOUT is bidirectional. The data type or SIGNAL_TYPE can be BIT, STD_LOGIC, INTEGER, etc.

Example 1: Entity of AND gate with two inputs each of one bit.

```
ENTITY AND_GATE IS
  PORT (
    IN_A : IN STD_LOGIC;
    IN_B : IN STD_LOGIC;
    OUT_Z : OUT STD_LOGIC
  );
END AND_GATE;
```

Example 2: Entity of 4x1 MUX with each input of three bits

```
ENTITY MUX IS
  PORT (
    A, B, C, D : IN STD_LOGIC_VECTOR(2 DOWNTO 0);
    SEL : IN STD_LOGIC_VECTOR(1 DOWNTO 0);
```

```

        Z : OUT STD_LOGIC_VECTOR(2 DOWNTO 0)
    );
END MUX;

```

ARCHITECTURE

The ARCHITECTURE describes how the circuit should function. It describes the internal implementation of the associated entity. There are several models that are followed by architecture to describe the operation of the circuit.

The general form of ARCHITECTURE is:

```

ARCHITECTURE architecture_name OF entity_name IS
    [declarations]
BEGIN
    [code]
END architecture_name;

```

Here, declarative part is optional and includes signal and constant declarations. Code part includes different VHDL statements describing the system to be designed.

Example 1: ARCHITECTURE of AND gate with two inputs each of one bit.

```

ARCHITECTURE AND_ARCH OF AND_GATE IS
BEGIN
    OUT_Z <= IN_A AND IN_B;
END AND_ARCH;

```

10.3 Data types, Data Objects and Operators

A. DATA TYPES

Pre-Defined Data Types

VHDL contains a series of pre-defined data types. Such data type definitions can be found in various packages or libraries.

- Package STANDARD of library STD includes BIT, BOOLEAN, INTEGER, and REAL.
- Package STD_LOGIC_1164 of library IEEE includes STD_LOGIC and STD_ULOGIC.

Various pre-defined data types are listed in the table below:

SN	TYPE	LEVEL/RANGE	DESCRIPTION
1.	BIT	2 LOGIC LEVEL	0, 1
2.	BIT_VECTOR	2 LOGIC LEVEL	0, 1

3.	STD_LOGIC	8 VALUED LOGIC	X, 0, 1, Z, W, L, H, -
4.	STD_LOGIC_VECTOR	8 VALUED LOGIC	X, 0, 1, Z, W, L, H, -
5.	STD_ULOGIC	9 VALUED LOGIC	U, X, 0, 1, Z, W, L, H, -
6.	STD_ULOGIC_VECTOR	9 VALUED LOGIC	U, X, 0, 1, Z, W, L, H, -
7.	BOOLEAN	2 VALUES	TRUE, FALSE
8.	INTEGER	-2147483647 TO +2147483647	32 BIT NUMBER
9.	NATURAL	0 TO +2147483647	
10.	REAL	-1.0E38 TO +1.0E38	
11.	SIGNED	POSITIVE AND NEGATIVE	USED IN ARITHMETIC OPERATIONS
12.	UNSIGNED	POSITIVE	
13.	PHYSICAL	TIME, VOLTAGE	USED IN SIMULATION

Here, various logic levels represent: 'U' – Unresolved, 'X' – Forcing Unknown, '0' – Forcing Low, '1' – Forcing High, 'Z' – High Impedance, 'W' – Weak unknown, 'L' – Weak Low, 'H' – Weak High, '-' – Don't care. STD_LOGIC levels are intended for simulation only. When a node has two STD_LOGIC signals connected, then conflicting logic levels are resolved automatically in case of STD_LOGIC whereas such conflict is not resolved in STD_ULOGIC. For arithmetic operations using STD_LOGIC, packages STD_LOGIC_SIGNED and STD_LOGIC_UNSIGNED must be used.

User Defined Data Types

VHDL allows users to define their own data types. There are two categories of user-defined data types.

- User-Defined Integer Type

General form:

TYPE TYPE_NAME IS RANGE LOW_VALUE TO HIGH_VALUE;

Example:

TYPE TEMPERATURE IS RANGE -125 TO 125;

TYPE MARKS IS RANGE 0 TO 100;

- User-Defined Enumerated Type

General Form:

TYPE TYPE_NAME **IS** (VALUE1, VALUE2... VALUEN);

Example:

TYPE COLOR **IS** (RED, GREEN, BLUE, WHITE);

Based on bits requirement encoding of enumerated type is done sequentially and automatically, unless specified.

B. DATA OBJECTS

An object is an item in VHDL that has both name and a specific type. Commonly used data objects are signals, variables and constants.

CONSTANTS are used to assign default values in the code. It can be declared in PACKAGE, ENTITY or ARCHITECTURE. Declaring **CONSTANTS** in PACKAGE makes it global, since PACKAGE can be used by several entities. If it is declared in an ENTITY, it can be shared by all ARCHITECTURE that follows that ENTITY. When defined within ARCHITECTURE the scope of **CONSTANTS** are limited to that ARCHITECTURE only.

Declaration:

CONSTANT name:TYPE:= value;

Examples:

CONSTANT high:STD_LOGIC:=’1’;

CONSTANT count:INTEGER:=10;

SIGNAL is used to pass value in and out of the circuit and within internal units. It simply represents interconnection of circuit. All ports of ENTITY are signals by default. The change in the SIGNAL may not be updated immediately, since the value is more likely to get updated after the completion of its corresponding PROCESS, FUNCTION or PROCEDURE. Similar to **CONSTANT**, it can be declared in PACKAGE, ENTITY or ARCHITECTURE.

Declaration:

SIGNAL name: **TYPE** [range] [:= initial_value];

The part inside the square bracket may or may not be present depending upon data types used and requirement of initialization.

Examples:

SIGNAL start: STD_LOGIC:=’0’;

SIGNAL count: INTEGER RANGE 0 TO 100;

VARIABLE represents the local information. Its value cannot be passed out directly. The change in value is immediately updated; new value can be promptly used in next line of code. It can be declared and used inside a PROCESS, FUNCTION or PROCEDURE.

Declaration:

VARIABLE name: type [range] [:= initial value];

Examples:

VARIABLE count: INTEGER:=0;

VARIABLE a : STD_LOGIC_VECTOR(7 DOWNTO 0);

C. OPERATORS

The various operators supported by VHDL are tabulated below:

ASSIGNMENT OPERATORS

SN	Operator	Assign Value To	Examples
1.	<=	Signal	X <= '1'; Y<="101";
2.	:=	Variable, constant, generic, and for initialization	Z := "1001"; Z is a variable
3.	=>	Individual Elements or with OTHERS	W <= (0=> '1', OTHERS => '0') LSB assigned 1 and others as 0

LOGICAL OPERATORS

SN	Operators	Description/Example	Supported Data Type
1.	NOT	Inverts the signal, High Precedence	BIT STD_LOGIC STD_ULOGIC BIT_VECTOR STD_LOGIC_VECTOR STD_ULOGIC_VECTOR
2.	AND	Result high when both inputs is high	
3.	OR	Result high when one of the inputs is high	
4.	NAND	X <= a NAND b	
5.	NOR	Z <= NOT a NOR B	
6.	XOR	Complements the bit when XORed with 1	
7.	XNOR	Complements the bit when XNORed with 0	

RELATIONAL OPERATORS

SN	Operators	Description
1.	=	Equal to

2.	/=	Not equal to
3.	<	Less than
4.	>	Greater than
5.	<=	Less than or equal to
6.	>=	Greater than or equal to

ARITHMETIC OPERATORS

SN	Operator	Meaning	Description
1.	+	Addition	
2.	-	Subtraction	
3.	*	Multiplication	
4.	/	Division	Limited to powers of two
5.	**	Exponentiation	Limited to powers of two
6.	MOD	Modulus	X MOD Y results value with sign of Y
7.	REM	Remainder	X REM Y results value with sign of X
8.	ABS	Absolute Value	
Avoid Using MOD operator when dealing with negative numbers			

SHIFT OPERATORS

SN	Operator	Meaning	Description
1.	SLL	Shift Left Logic	Zeros are fed from one end and bits are lost from other end. Sign bit never changes in arithmetic shift. "10101" SLL 3 results in "01000"
2.	SRL	Shift Right Logic	
3.	SLA	Shift Left Arithmetic	
4.	SRA	Shift Right Arithmetic	
5.	ROL	Rotate Left	"1001" ROL 2 results in "0110"
6.	ROR	Rotate Right	

CONCATENATION OPERATOR

The concatenation operator (&) is used to combine values of similar data type. The following example will illustrate the use of concatenation operator.

Example:

```

signal A, B : std_logic_vector (3 downto 0);    -- Signal A and B of 4 bits
signal C : std_logic_vector (5 downto 0);      -- Signal C of 6 bits

```

```

signal D : std_logic_vector (7 downto 0);      -- Signal D of 8 bits
C <= A & "00" ;                               -- 4 bits of A and two bits "00" assigned to C
D <= B & A ;                                   -- 4 bit of A and B combined and assigned to D

```

10.4 STATEMENTS IN VHDL

A. CONCURRENT STATEMENTS

Concurrent Signal Assignment

Syntax:

Target <= expression;

Examples:

```

A <= B NAND C;
X <= (D OR E) AND (F AND G);

```

Conditional Signal Assignment

Syntax:

*Target <= expression when condition else
 expression when condition else
 expression;*

Example:

```

Z <=  '1' when (L='0' AND M='0') else
      '1' when (L='1' AND M='1') else
      '0';

```

Selective Signal Assignment

Syntax:

*with choose_expression select
 target <= expression when choices,
 expression when choices;*

Example:

```

with SEL select
      M_OUT <=    A3 when "11",
                  A2 when "10",
                  A1 when "01",

```

A0 when "00",
'0' when others;

Process Statement

Syntax:

```
label: process(sensitivity list)
begin
    sequential statements
end process label;
```

B. SEQUENTIAL STATEMENTS

Signal Assignment

Syntax:

target <= *expression*;

Example:

```
A <= B NAND C;
X <= (D OR E) AND (F AND G);
```

IF statements

Syntax:

```
if (condition) then
    { sequence of statements }
elsif (condition) then
    { sequence of statements }
else
    { sequence of statements }
end if;
```

Example:

```
if (SEL = "111") then F_OUT <= D(7);
elsif (SEL = "110") then F_OUT <= D(6);
elsif (SEL = "101") then
    F_OUT <= D(1);
elsif (SEL = "000") then
    F_OUT <= D(0);
```

```

else F_OUT <= '0';
end if;

```

CASE statements

Syntax:

```

case (expression) is
  when choices =>
    sequential statements
  when choices =>
    sequential statements
  when others =>      -- (optional)
    sequential statements
end case;

```

Example:

```

case (ABC) is
  when "100" =>
    F_OUT <= '1';
  when "011" =>
    F_OUT <= '1';
  when "111" =>
    F_OUT <= '1';
  when others =>
    F_OUT <= '0';
end case;

```

10.5 Standard Architectures

A. Dataflow Style Architecture

Dataflow style architecture specifies a circuit as a concurrent representation of the flow of data through the circuit. In this modeling, the internal working of a system is implemented using concurrent statements. It can be used for small and primitive circuits but not for complex designs. In this style of architecture, whenever there is a change in signal of right hand side, the expression is evaluated and assigned to left hand side.

Example:

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY HALF_ADDER IS
    PORT(
        A, B: IN STD_LOGIC;
        S, C: OUT STD_LOGIC
    );
END HALF_ADDER;

ARCHITECTURE HALF_ADDER_ARCH OF HALF_ADDER IS
BEGIN
    S <= A XOR B;
    C <= A AND B;
END HALF_ADDER_ARCH;

```

B. Behavior Style Architecture

The behavioral style architecture models how the circuit outputs will behave to the circuit inputs. This model may not reflect how the circuit is implemented when it is synthesized. Process statement is the core part of behavioral style architecture. In this style, the internal working is implemented using sequential statements within process statements.

Example:

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY HALF_ADDER IS
    PORT(
        A, B: IN STD_LOGIC;
        S, C: OUT STD_LOGIC
    );
END HALF_ADDER;

```



```

ARCHITECTURE HALF_ADDER_ARCH OF HALF_ADDER IS
BEGIN
    PROCESS_ADDER: PROCESS (A, B)
    BEGIN
        S <= A XOR B;
        C <= A AND B;
    END PROCESS PROCESS_ADDER;
END HALF_ADDER_ARCH;

```

C. Structural Style Architecture

The structural style architecture is a modular approach to coding which supports hierarchical design which is essential to understand complex digital designs. Modular designs enhance understandability by combining low-level functionality into modules. These modules can be reused in different designs resulting in save of design time. VHDL structural model may not be efficient for simple designs. However, the following are the general steps for writing structural model code.

- Initially the entity and architecture implementations for the individual gates or modules which are within our system must be defined.
- The entity declaration of our system is done, similar to other models.
- Different components used in our design are declared within the declarative part of architecture. Component declaration is similar to entity declaration, only keyword entity must be replaced by keyword component.
- Internal signals, which are the intermediate output signals of one module fed into another module as input signals, are declared.
- Finally, instances of all modules are created and mapped in the architecture body. Mapping can be done using direct mapping or implied mapping. In direct mapping, each of the internal signals and signals of entity of the system are directly associated with the signals of corresponding components. Whereas in implied mapping, only internal signals and signals of entity of the system are listed. Though it uses less space, but it requires the signals be placed in the proper order.

Example: To implement $Z = (A \text{ AND } B) \text{ OR } (C \text{ AND } D)$ using structural model

<pre> LIBRARY IEEE; USE IEEE.STD_LOGIC_1164.ALL; ENTITY AND_GATE IS PORT (X, Y: IN STD_LOGIC; W: OUT STD_LOGIC); END AND_GATE; ARCHITECTURE AND_AH OF AND_GATE IS BEGIN PROCESS(X, Y) BEGIN W <= X AND Y; END PROCESS; END AND_AH; </pre>	<pre> LIBRARY IEEE; USE IEEE.STD_LOGIC_1164.ALL; ENTITY OR_GATE IS PORT (X, Y: IN STD_LOGIC; W: OUT STD_LOGIC); END OR_GATE; ARCHITECTURE OR_ARCH OF OR_GATE IS BEGIN PROCESS(X, Y) BEGIN W <= X OR Y; END PROCESS; END OR_ARCH; </pre>
---	---

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY TEST IS
  PORT (
    A, B, C, D: IN STD_LOGIC;
    Z: OUT STD_LOGIC
  );
END TEST;

ARCHITECTURE TEST_ARCH OF TEST IS
  COMPONENT AND_GATE
    PORT (

```

```

        X, Y: IN STD_LOGIC;
        W: OUT STD_LOGIC
    );
END COMPONENT;
COMPONENT OR_GATE
    PORT (
        X, Y: IN STD_LOGIC;
        W: OUT STD_LOGIC
    );
END COMPONENT;
SIGNAL E, F: STD_LOGIC;
BEGIN
    U1: AND_GATE PORT MAP (X => A, Y => B, W => E);
    U2: AND_GATE PORT MAP (X => C, Y => D, W => F);
    U3: OR_GATE PORT MAP {X => E, Y => F, W => Z};
END HALF_ADDER_ARCH;

```

10.6 FSM Design

Finite State Machines (FSM) constitute a special modeling technique for sequential logic circuits. The digital systems, in general, can be expressed as a sequence of actions which can be realized using FSM.

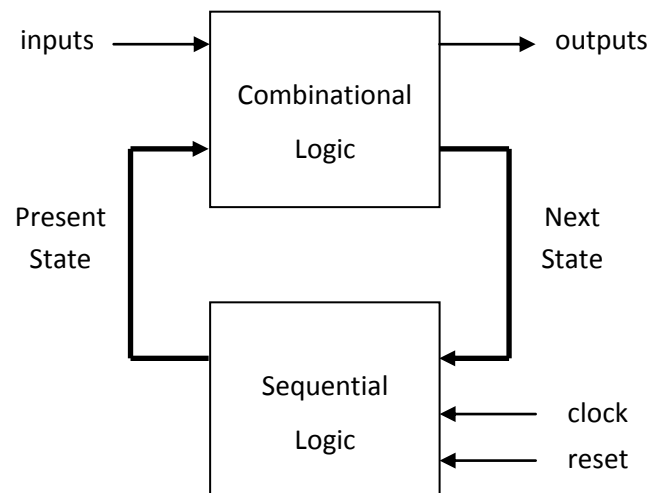


Figure 10.1: General block diagram of Finite State Machine

A FSM is specified by five entities: symbolic states, input signals, output signals, next-state function and output function. A state specifies a unique internal condition of a system and the FSM transits from one state to another with time. The next-state function is used to determine the next state of the system. The output function specifies the value of the output signals. The general block diagram of FSM is shown in the figure 10.1.

FSM consists of two sections; combinational and sequential logic. The combinational part has two inputs – external input and present state – and two outputs; next state and external output. Whereas, the sequential section has three inputs – clock, reset, and next state – and one output in a form of present state. Since the flip flops are implemented in sequential logic, clock and reset are part of this section.

If the output of the machine depends not only on the present state but also on the current input, then it is called a Mealy machine. Otherwise, if it depends only on the current state, it is called a Moore machine.

A. Design of Sequential Section

PROCESS statement is required for sequential section. The clock and reset signals appear in the sensitivity list of **PROCESS** statement. When reset is asserted, present state will be set to initial state of the system. In other cases, present state will change to next state at the proper clock edge. A typical design template for the sequential section is given as:

```
PROCESS (RESET, CLOCK)
BEGIN
    IF(RESET = '1') THEN
        PRESENT_STATE <= INITIAL_STATE;
    ELSIF (CLOCK'EVENT AND CLOCK = '1') THEN
        PRESENT_STATE <= NEXT_STATE;
    END IF;
END PROCESS;
```

B. Design of Combinational Section

In this section, the code does not need to be sequential; concurrent code can be used. If sequential is implemented then the input and present state will be the part of sensitivity list of PROCESS statement. Within the PROCESS statement, CASE statement is used to implement the actions and conditions of each state. A basic template for the combinational section is shown as:

```

PROCESS (INPUT, PRESENT_STATE)
BEGIN
    CASE PRESENT_STATE IS
        WHEN STATE0 =>      -- WITHIN WHEN STRUCTURE OF CASE,
            ...              -- MAY CONTAIN ACTIONS AND CONDITIONS
        WHEN STATE1 =>
            ...              -- NUMBER OF WHEN STRUCTURE IS
        WHEN STATE2 =>      -- DEFINED BY NUMER OF STATES IN FSM
            ...
        WHEN OTHERS =>
            ...
    END CASE;
END PROCESS;

```

FEW SOLVED EXAMPLES

PROBLEM 1: Simple NAND Gate with two inputs, each input of single bit

```

LIBRARY IEEE
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY NAND_GATE IS

    PORT (
        IN_A, IN_B: IN STD_LOGIC;
        X_OUT: OUT STD_LOGIC
    );
END NAND_GATE;

ARCHITECTURE NAND_ARCH OF NAND_GATE IS
BEGIN
    PROC: PROCESS (IN_A, IN_B)
    BEGIN
        X_OUT <= IN_A NAND IN_B;
    END PROCESS PROC;
END NAND_ARCH;

```

PROBLEM 2: Write a VHDL code to implement 4 X 1 MUX with each input of 3 bits.

```

LIBRARY IEEE
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY MUX_4X1 IS

    PORT (
        IN_A, IN_B : IN STD_LOGIC_VECTOR (2 DOWNTO 0);
        IN_C, IN_D : IN STD_LOGIC_VECTOR (2 DOWNTO 0);
        SEL : IN STD_LOGIC_VECTOR(1 DOWNTO 0);
        Z_OUT : OUT STD_LOGIC_VECTOR (2 DOWNTO 0)
    );
END MUX_4X1;

ARCHITECTURE MUX_ARCH OF MUX_4X1 IS

```

```

BEGIN

  PROC: PROCESS (IN_A, IN_B, IN_C, IN_D, SEL)
    BEGIN
      IF (SEL = "00") THEN
        Z_OUT <= IN_A;
      ELSIF (SEL = "01") THEN
        Z_OUT <= IN_B;
      ELSIF (SEL = "01") THEN
        Z_OUT <= IN_C;
      ELSE
        Z_OUT <= IN_D;
      END IF;
    END PROCESS PROC;
END MUX_ARCH;

```

PROBLEM 3: Write a VHDL code to implement D flip flop.

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY DFLIPFLOP IS
  PORT (
    D, CLK : IN STD_LOGIC;
    Q : OUT STD_LOGIC
  );
END DFLIPFLOP;

ARCHITECTURE BEHAVIORAL OF DFLIPFLOP IS
BEGIN
  PROCESS (D, CLK)
    BEGIN
      IF (CLK'EVENT AND CLK = '1') THEN
        Q <= D;
      END IF;
    END
  END

```

```

    END PROCESS;
END BEHAVIORAL;

```

PROBLEM 4: Implement a counter that counts from 0 to 9 using VHDL code

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;
USE IEEE.STD_LOGIC_ARITH.ALL;
USE IEEE.STD_LOGIC_UNSIGNED.ALL;

ENTITY COUNTER_CODE IS
    PORT (
        CLR : IN STD_LOGIC;
        CLK : IN STD_LOGIC;
        Q : OUT STD_LOGIC_VECTOR (3 DOWNTO 0)
    );
END COUNTER_CODE;

ARCHITECTURE BEHAVIORAL OF COUNTER_CODE IS
    SIGNAL TMP: STD_LOGIC_VECTOR (3 DOWNTO 0);
BEGIN
    PROCESS (CLK, CLR)
    BEGIN
        IF (CLK'EVENT AND CLK = '0') THEN
            IF (CLR = '1') THEN
                TMP <= "0000";
            ELSE
                TMP <= TMP + 1;
            END IF;
        END IF;
    END PROCESS;
    Q <= TMP;
END BEHAVIORAL;

```


Problem 5: Write a VHDL code to detect a sequence of “1001”

→ The FSM for the detection of the sequence is given by following diagram.

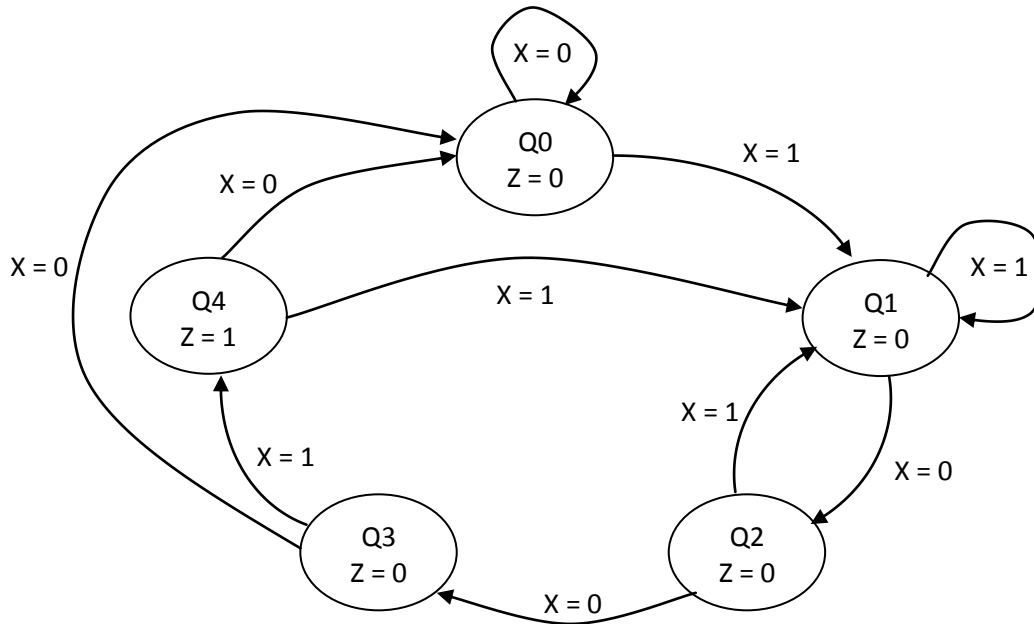


Figure 10.2: FSM for detection of sequence “1001”

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY SEQUENCE_STATE IS
  PORT (
    X : IN STD_LOGIC;
    CLK, RESET : IN STD_LOGIC;
    Z : OUT STD_LOGIC
  );
END SEQUENCE_STATE;

ARCHITECTURE BEHAVIORAL OF SEQUENCE_STATE IS
  TYPE STATE IS (Q0, Q1, Q2, Q3, Q4);
  SIGNAL PS, NS: STATE;

BEGIN
  SYNC_PROC: PROCESS (CLK, RESET)
  BEGIN

```

```

IF (RESET = '1') THEN
    PS <= Q0;
ELSIF (RISING_EDGE(CLK)) THEN
    PS <= NS;
END IF;
END PROCESS SYNC_PROC;

COMB_PROC: PROCESS(PS,X)
BEGIN
    CASE PS IS
        WHEN Q0 => Z <= '0';
            IF(X='1') THEN
                NS <= Q1;
            ELSE
                NS <= Q0;
            END IF;
        WHEN Q1 => Z <= '0';
            IF(X='0') THEN
                NS <= Q2;
            ELSE
                NS <= Q1;
            END IF;
        WHEN Q2 => Z <= '0';
            IF(X='0') THEN
                NS <= Q3;
            ELSE
                NS <= Q1;
            END IF;
        WHEN Q3 => Z <= '0';
            IF(X='1') THEN
                NS <= Q4;
            ELSE
                NS <= Q0;

```

```

        END IF;
    WHEN Q4 => Z <= '1';
        IF(X='1') THEN
            NS <= Q1;
        ELSE
            NS <= Q0;
        END IF;
    WHEN OTHERS =>
        NS <= Q0;
    END CASE;
END PROCESS COMB_PROC;
END BEHAVIORAL;

```

PROBLEM 6: Calculate the GCD of two numbers using VHDL

➔ Functionality code to calculate the GCD of two numbers is given as

```

int X, Y;
while(1)
{
    while(!GO);
    X = NUM1;
    Y = NUM2;
    while(X != Y)
    {
        if(X<Y)
            Y = Y - X;
        else
            X = X - Y;
    }
    GCD = X;
}

```

THE FSM FOR THE ABOVE CODE CAN BE REPRESENTED BY FOLLOWIING DIAGRAM

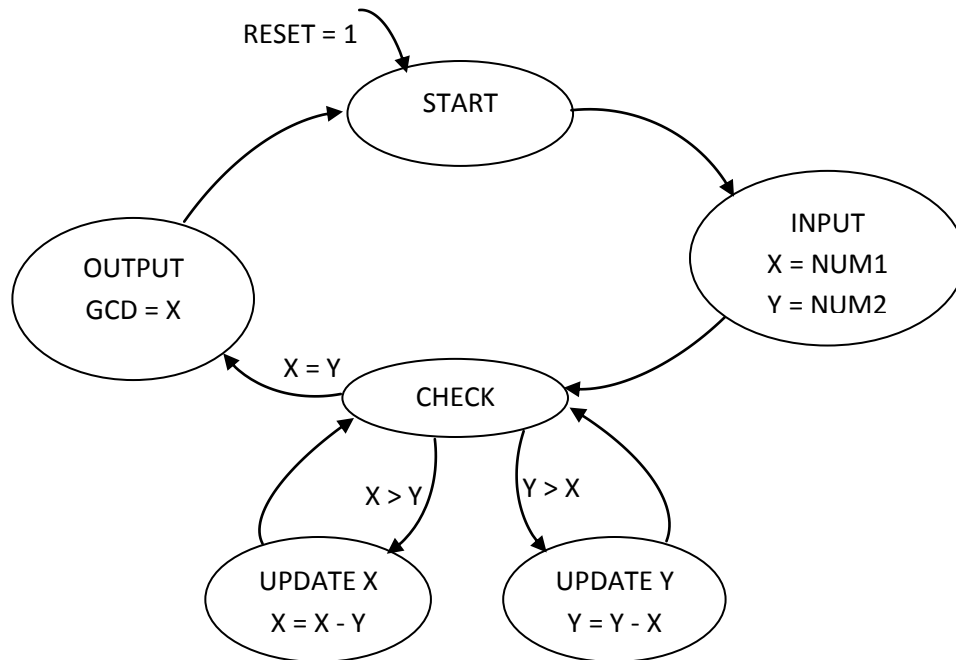


Figure 10.2: FSM for GCD processor

VHDL CODE

```

LIBRARY IEEE;
USE IEEE.STD_LOGIC_1164.ALL;

ENTITY FSM_GCD IS
  PORT (
    RESET, CLK: IN STD_LOGIC;
    GO: IN STD_LOGIC;
    NUM1, NUM2: IN INTEGER;
    GCD: OUT INTEGER
  );
END FSM_GCD;

ARCHITECTURE BEHAVIORAL OF FSM_GCD IS
  TYPE STATE IS (START, INPUT, CHECK, UPDATEX, UPDATE Y, OUTPUT);
  SIGNAL PS, NS: STATE;
BEGIN

```

```

SEQ_PROC: PROCESS (CLK, GO, RESET)
BEGIN
    IF (GO = '1') THEN
        IF (RESET = '1') THEN
            PS <= START;
        ELSIF (RISING_EDGE(CLK)) THEN
            PS <= NS;
        END IF;
    END IF;
END PROCESS SEQ_PROC;

```

```

COMB_PROC: PROCESS (NUM1, NUM2, PS)
    VARIABLE X, Y: INTEGER;
BEGIN
    CASE PS IS
        WHEN START =>
            GCD <= 0;
            NS <= INPUT;
        WHEN INPUT =>
            X := NUM1;
            Y := NUM2;
            NS <= CHECK;
        WHEN CHECK =>
            IF(X>Y) THEN
                NS <= UPDATEX;
            ELSIF(X<Y) THEN
                NS<= UPDATEY;
            ELSE
                NS <= OUTPUT;
            END IF;
        WHEN UPDATEX =>
            X := X - Y

```

```
        NS <= CHECK;  
    WHEN UPDATEY =>  
        Y := Y - X;  
        NS <= CHECK;  
    WHEN OUTPUT =>  
        GCD <= X;  
        NS <= INPUT;  
    WHEN OTHERS =>  
        GCD <= 0;  
        NS <= INPUT;  
    END CASE;  
END PROCESS COMB_PROC;  
END BEHAVIORAL;
```