

Computer Organization & InstructionIntroduction

\* Computer architecture acts as the interface between the hardware and lowest level software.

\* Computer architecture is defined as study of the structure, behavior, and design of computers.

Computer Organization. It refers to the operational units and their interconnections that realize the architectural specification. It describes the function of and design of the various units of digital computer that store and process information. The attributes in computer organization refers to.

- \* Control signals.
- \* Computer/Peripheral interface.
- \* Memory technology.

Computer Hardware: It consists of electronic circuits, displays, magnetic and optical storage media, electromechanical equipment and communication facilities.

Computer architecture: It is concerned with the structure and behavior of the computer. It includes the information formats, the instruction set and techniques for addressing memory.



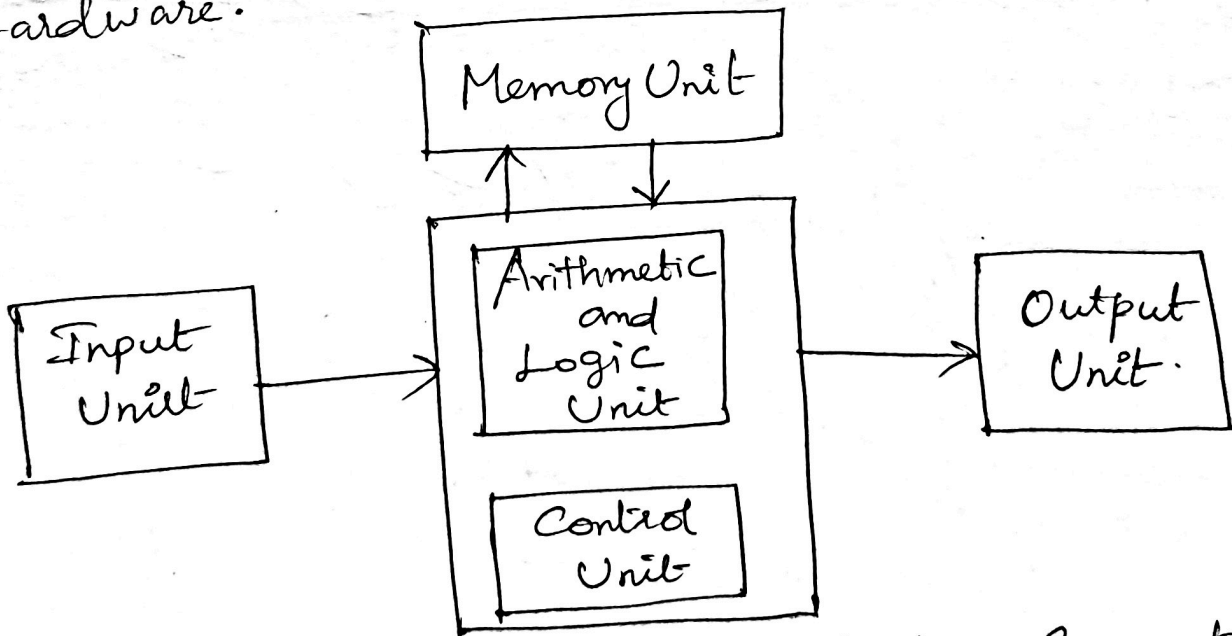
→ The attributes of Computer architecture refers to the

- \* Instruction Set.
- \* Data representation.
- \* I/O Mechanism.
- \* Addressing Techniques.

## Basics of Computer System.

A Computer in its simplest form is a fast electronic calculating machine. that accepts digitized information from the user, Processes it according to a sequence of instructions and provides the processed information to the user.

Hardware: Electronic Components inter connected in the Computer System as a whole is referred to as Hardware.



Components of Computer System Connected as Hardware.



Input Unit: \* It Provides the data to the Computer System from the outside. So it links the external environment with the Computer. It takes the data from the input devices, converts it into machine language and then loads it into the Computer Systems. \* Most Commonly Used input devices are keyboard and mouse.

Output Unit \* It provides the results of Computer process to the users. i.e) It links the computer with the external environment. \* Different output devices are monitors, printers, speakers, headphones etc.

Memory Unit: \* Memory Unit is used to store programs and data.

\* Primary Storage memory and Secondary Storage memory devices form a total memory unit.

\* Primary memory is directly accessible by CPU. Secondary memory is not accessible by CPU.

\* Primary memory (main memory) is fast Semiconductor RAM.

\* Secondary Storage memories such as magnetic tapes, magnetic disk are used for the storage of larger amount of data.



## Arithmetic and Logic Unit:

\* ALU is responsible for performing arithmetic and logical operations.

\* To perform these operations operands from main memory are brought into internal registers of processor.

\* After performing operation the result is either stored in the register or memory location.

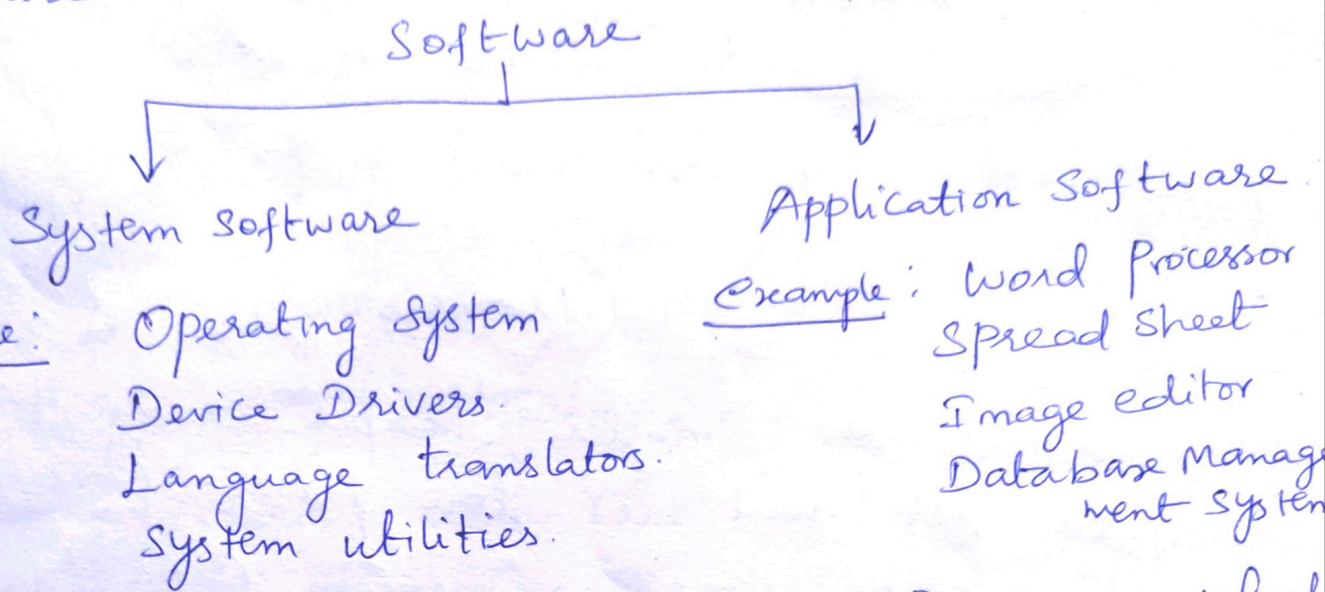
Control Unit \* A Control Unit co-ordinates and controls all the activities amongst the functional units.

\* Control Unit transfers data from Storage Unit to arithmetic logic unit when calculation need to be performed.

\* The combine unit of arithmetic logic unit and control unit is called as Central Processing Unit.



Software. Computer Software is divided into two broad categories: System software and application software.



System Software. \* It is collection of Programs which is needed in the Creation, Preparation and execution of other program.

\* Typical System software includes editors, assemblers, linker, loader, Compiler, debuggers.

Application Software: Allows to perform specific task on a Computer using the Capabilities of a Computer.

## Evolution of Computer Systems

### First Generation (vacuum Tubes)

\* First Electronic Computer ENIAC (Electronic Numerical Integrator and Computer)

\* It was made up of more than 18000 Vacuum tubes and 1500 relays.

\* It was able to perform nearly 5000 additions or subtractions per second.



\* Its data memory consists of 20 accumulators each capable of storing a ten digit decimal number.

Features. weight - 30 tons, area - 15000 Sq. ft.

Power Consumption - 140 kW.

Draw back - It was wired in for specific computation.

Second Generation: (Transistors).

\* Transistors are smaller, cheaper and they dissipate less heat i.e. low power consumption.

\* Greater speed, larger memory capacity and smaller size than first generation.

\* Due to addition of a set of index register and arithmetic circuits, CPU can handle both floating point and fixed point operations.

\* Introduction of more complex arithmetic and logic units and control unit to support higher level programming languages. for example FORTRAN.

\* Magnetic core memories and magnetic drum storage devices were more widely used.

Third Generation: (Integrated Circuits).

\* Uses integrated circuits.

\* Integrated circuit technology enabled lower cost, faster processors and development of memory chips.



\* Integrated circuit technology allowed to increase memory size and number of I/O ports.

\* Magnetic core memories were replaced by integrated circuit memories.

\* Various techniques were introduced to improve the performance of the computer. These are.

→ Microprogramming

→ Parallel processing a) Multiprocessing b) Pipelining.

→ Sharing resources.

\* Introduction of IBM System/360 - first planned family of computer products.

Feature of IBM 360-370.

\* It uses 32 bit or 4 byte word format

\* It has 16 general register and 4 double length floating point registers.

\* Increasing speed, Increasing number of I/O Port.

VLSI (Very Large Scale Integrated Circuits)

\* Uses Large Scale Integration and Very Large Scale Integration (VLSI) technology for computer design.

\* VLSI technology made it possible to fabricate an entire CPU, main memory or similar device with a single IC that can be mass produced at very low cost.



\* Fourth generation Products are portable notebook Computers, desktop Computers. and workstation<sup>e</sup> interconnected by local area network, wide area network. and the internet.

### ULSI (Fifth Generation)

\* Ultra Large Scale Integration is the Process of integrating or embedding millions of transistors on a single silicon semiconductor microchip

\* It is a successor to large scale integration (LSI) and Very large scale integration.

\* Parallel Processing Technique Used.



## Performance:

(5)

\* The performance is an important attribute of a Computer. It is an important Criteria for Selection of a Computer.

## Defining Performance:

\* The Computer User is always interested in reducing the time between the start and completion of the program or event. (e) reducing the execution time.

\* The execution time is also referred to as response time. Reduction in response time increases the throughput. (the total amount of work done in a given time).

\* The Performance of the Computer is directly related to throughput and hence it is reciprocal of execution time.

$$\text{Performance}_A = \frac{1}{\text{Execution time}_A}$$

This means that for two Computers A and B if the performance of A is greater than the performance of B, we have.

$$\text{Performance}_A > \text{Performance}_B$$
$$\frac{1}{\text{Execution time}_A} > \frac{1}{\text{Execution time}_B}$$

$$\text{Execution time}_B > \text{Execution time}_A$$



The execution time on B is longer than that on A.  
If A is faster than B.

→ "A is n times as fast as B" to mean

$$\frac{\text{Performance}_A}{\text{Performance}_B} = n.$$

If A is n times faster than B then the execution time on B is n times longer than it is on A.

$$\frac{\text{Performance}_A}{\text{Performance}_B} = \frac{\text{Execution time}_B}{\text{Execution time}_A} = n$$

Measuring Performance:

CPU Time

\* CPU execution time or simply CPU time is the time the CPU spends computing for particular task and does not include time spent waiting for I/O or running other programs.

\* CPU time can be divided into the CPU time spent in the program called User CPU time and CPU time spent in the Operating System performing tasks on behalf of the program, called System CPU time.

Performance Metrics:

\* Users and designers often examine performance using different metrics.



CPU execution time for a Program =

(6)

CPU clock cycles for a Program  $\times$  clock cycle time

alternatively, because clock rate and clock cycle time are inverse.

CPU execution time for a Program =  $\frac{\text{CPU clock cycle for a Program}}{\text{clock rate}}$ .

This formula makes it clear that the hardware designer can improve performance.

Problem:

If Computer A runs a Program in 10 seconds and Computer B runs the same Program in 25 sec. How much faster is A and B?

$$\frac{\text{Performance}_A}{\text{Performance}_B} = \frac{\text{Execution time}_B}{\text{Execution time}_A} = n$$

Thus the performance ratio is

$$\frac{25}{10} = 2.5$$

A is therefore 2.5 times faster than B.

Computer A runs a Program in 12 seconds with a 3 GHz clock. We have to design a Computer B such that it can run the same Program with in 9 seconds. Determine the clock rate for Computer B. Assume that due to increase in clock rate CPU design of Computer B is effected and it requires 1.2 times as many clock cycles as Computer A for execution this Program.



Soln

$$\text{Clock rate}_A = 3 \times 10^9 \text{ cycle/sec.}$$

$$\text{Cpu time}_A = 12 \text{ Seconds.}$$

$$\text{Cpu time}_B = 9 \text{ Sec.}$$

$$\text{Cpu time}_A = \frac{\text{CPU clock cycles}_A}{\text{clock rate}_A}$$

$$12 \text{ sec} = \frac{\text{CPU clock cycle}_A}{3 \times 10^9 \text{ cycles/sec}}$$

$$\begin{aligned} \text{CPU clock cycle}_A &= 12 \text{ sec} \times 3 \times 10^9 \text{ cycle/sec} \\ &= 36 \times 10^9 \text{ cycles.} \end{aligned}$$

Cpu time for Computer B is

$$\text{Cpu time}_B = \frac{\text{CPU clock cycles}_B}{\text{clock rate}_B} = \frac{1.2 \times \text{CPU clock cycle}_A}{\text{clock rate}_B}$$

$$9 \text{ Seconds} = \frac{1.2 \times 36 \times 10^9 \text{ cycle}}{\text{clock rate}_B}$$

$$\begin{aligned} \text{Clock rate}_B &= \frac{1.2 \times 36 \times 10^9 \text{ cycle}}{9 \text{ sec}} = 4.8 \text{ cycle/sec} \\ &= 4.8 \text{ GHz.} \end{aligned}$$

### Power Wall

\* Power wall refers to the Representational wall signifying the Peak power Constraint of a System.

\* Running a Processor at high clock speeds allows for better performance. However when Processor runs at a higher speed, it generates more heat and consumes more power.



Thus when there is an  $\uparrow$  increase in clock rate, there is an increase in power consumed.

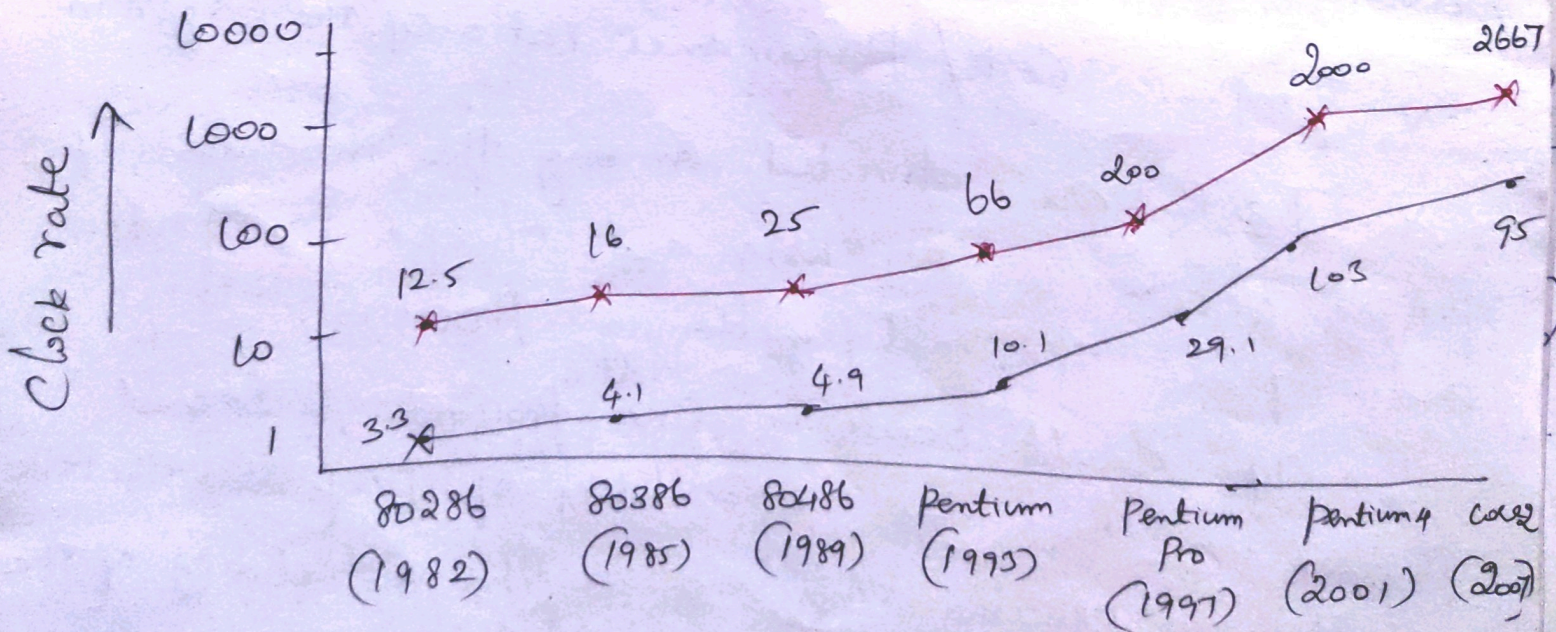
\* Fig shows the increase in clock rate and power of Intel microprocessors. The graph is flattened in the recent years because there is a practical limit for cooling microprocessors and hence increase in power dissipation.

\* Power consumed by a CPU is given by

$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{frequency}$$

$$P = C V^2 f$$

The  $f$  is the function of clock rate. The  $C$  is a function of both the number of transistors connected to an output and technology.





## From Uniprocessor to Multiprocessors

→ The Performance of the Computers has drastically increased when the technology has drifted from Uniprocessor Systems to Multiprocessor System. The performance of the processors also increased significantly.

→ Uniprocessor System is a type of architecture that is based on a single Computing Unit. All the operations were done sequentially on the same unit. Multiprocessor Systems are based on executing instructions on multiple Computing Units.

→ The multiprocessor architecture is based on Flynn Taxonomy.

### Advantages of Multiprocessor

- ✦ Improves Cost/Performance ratio of the System
- ✦ Tasks are divided among the modules. If failure occurs, it is easier and cheaper to find and replace the malfunctioning Processor.
- ✦ If fault occurs in one Processor, a second Processor can take the responsibility of performing the task of failure processor. This improves the reliability of the Systems.



## ① Single Instruction, Single Data (SISD) ⑧

\* This is a Uniprocessor machine. which is Capable of executing a single instruction operating on a single data stream.

\* Most conventional computers have SISD architecture.

\* All the instructions and data to be processed have to be stored in primary memory.

## ② Multiple Instruction, Single Data (MISD)

\* An MISD computing system is a multiprocessor machine capable of executing different instructions on different processing elements but all of them operating on the same dataset.

## ③ Single Instruction Multiple Data (SIMD)

\* This machine is capable of executing the same instruction on all the CPUs but operating on different data streams.

## ④ Multiple Instruction, Multiple Data (MIMD)

\* This is capable of executing multiple instructions on multiple data sets.

\* Each PE in the MIMD model has separate instruction and data streams.

\* Unlike SIMD and MISD machines, PEs in MIMD machines work asynchronously.



## Instruction Types (or) Operation.

⑨

A Computer has a Set of instructions that allows the User to formulate any data Processing task.

It can be categorized into following instruction types:

- \* Data transfer
- \* Arithmetic
- \* Logical
- \* Program Control
- \* Input/output
- \* String Manipulation
- \* Translate Instruction
- \* Processor Control Instruction.

### ① Data Transfer Instruction.

Which copy information from one location to another either in the processor's internal register set or in external Main Memory.

Example:

MOVE → Transfer data from the Source location to destination location.

LOAD → Transfer data from a Memory location to CPU register.

STORE → Transfer data from a CPU register to Memory location.

PUSH → Transfer data from the source to stack

XCHG → Exchange; Swap the contents of the source and destination.

### ② Arithmetic: which perform operations on numerical data.

Example: ADD → Calculate Sum of two operands.

ADC → Add with Carry. Calculate Sum of two operands and the carry bit.

SUB → Subtract; Calculate difference of two numbers.

MUL → Multiply; Calculate Product of Two Operands.



DIV → Divide.

NEG → Negative. Change Sign of Operand.

INC → Increment; add 1 to operand.

DEC → Decrement; Subtract 1 from operand.

③ Logical operation which include Boolean and other non-numerical operations.

Example: NOT → Complement the operand

OR → perform bit wise logical OR of operand.

AND → perform bit wise logical AND of operand.

ROT → Rotate, Shift operand (left or right)

TEST → Test for specified condition and set flags.

④ Control Transfer Instruction:

which change the sequence in which Programs are executed.

(eg.) JUMP → Branch; Enter the specified address into PC; Unconditional transfer.

JUMP IF → Branch on Condition; Enter the specified address into PC only if the specified condition is satisfied; Conditional transfer.

INT → Interrupt; Creates a software interrupt.

⑤ Input-Output Instruction.

which cause information to be transferred between the processor or its main memory and external I/O device.

(eg.) IN → Input; Transfer data from specified I/O port or device to destination.

OUT → Output; Transfer the data from specified source to I/O port or device.



## ⑥ String Manipulation Instruction.

(10)

These manipulate string of byte, word, double etc

- eg. MOVSB → Move byte or word of string  
LODSB → Load byte or word of string  
CMPSB → Compare byte or word of string.

## ⑦ Translate Instruction.

These convert data from one format to the another.

eg. XLAT → Translate; Convert the given code into another by table lookup.

PACK → Convert the unpacked decimal number into packed decimal.

## ⑧ Processor Control Instruction:

These control processor operations.

eg. HLT → Halt; stop the program execution.

NO-OP → No operating nothing

CMC → Complement 'Carry' flag.

CLC → clear 'Carry' flag.

WAIT → Stops program execution, resume execution when condition is satisfied.



# Addressing Modes.

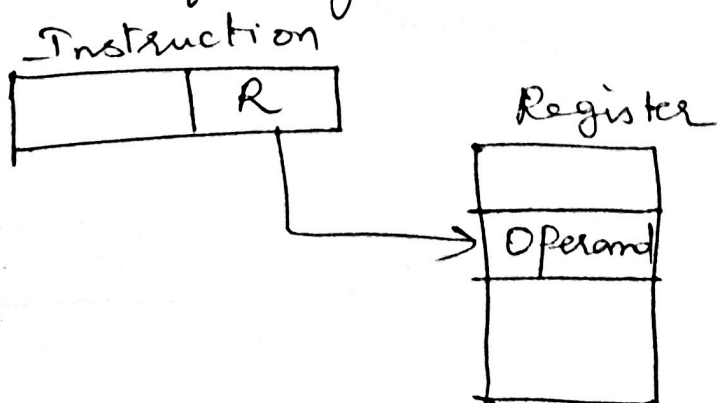
The different ways in which the location of an operand is specified in an instruction are referred to as addressing modes.

Types of addressing modes are.

- ① Register addressing mode.
- ② Absolute (or) Direct addressing mode
- ③ Immediate addressing mode
- ④ Indirect addressing mode.
- ⑤ Register Indexed addressing mode.
- ⑥ Displacement addressing mode.
- ⑦ Relative addressing mode.
- ⑧ Base register addressing mode.
- ⑨ Index addressing mode.

① Register addressing mode: The operand is the contents of processor register.

Example:  $\text{MOV } R_2, R_1$ ; This instruction copies the contents of register  $R_2$  to register  $R_1$ .

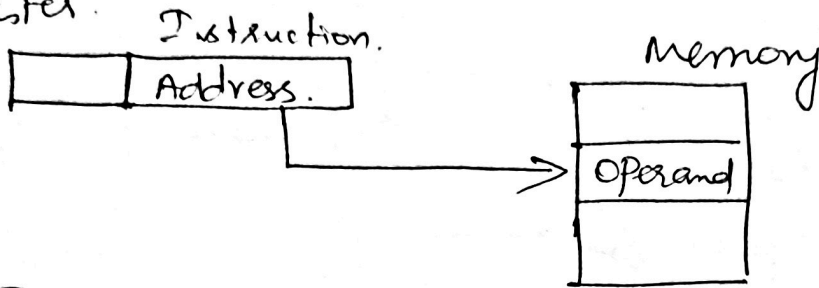




## 2) Absolute or Direct addressing mode:

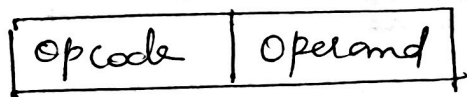
The address of the location of the operand is given explicitly as a part of the instruction.

Example Mov 2000, This instruction copies the contents of memory location 2000 into the A register.



## 3) Immediate addressing mode:

In this mode of addressing the operand is a part of the instruction and is specified in the address field.



Mov #20 A; This instruction is given explicitly in

Copies operand 20 in the register A.

Sign # in front of the value of an operand is used to indicate that this value is an immediate operand.

## 4) Indirect addressing mode:

In this addressing mode, the instruction contains the address of memory which refers the address of the operand. It is indirectly storing the address of the target location in another memory location.



⑤ Register Indirect addressing mode:

The effective address of the operand is the contents of a register or the main memory location whose address is given explicitly in the instruction.

MOV (R0), A; The instruction copies the contents of memory addressed by the contents of register R0 into the register A.

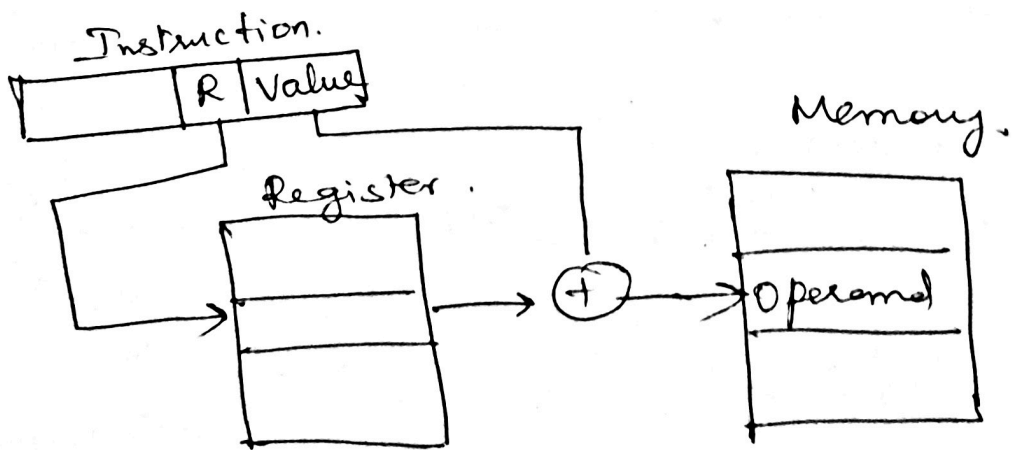
⑥ Displacement addressing mode:

→ This addressing mode combines the capabilities of indirect addressing and register indirect addressing.

→ In this addressing mode, instruction has two address fields: value and referenced register.

→ The effective address is computed by adding contents of referenced register to value.

$$EA = \text{Value} + (R)$$





# Relative addressing

(12)

This addressing mode is commonly used to specify the target address in branch instruction for example.

```
JNZ BACK.
```

This instruction causes program execution to go to the branch target location identified by the name BACK, if the branch condition is satisfied.

# Base Register addressing

\* In this addressing mode, the referenced register contains the main memory address and address field contains the displacement.

\* Displacement is usually integer number.

$$EA = (R) + \text{Displacement}$$

MOV [R+8], A : This instruction copies the contents of memory whose address is determined by adding the contents of register R and displacement 8 to the register A.

# Index addressing mode:

The effective address of the operand is generated by adding a constant value to the contents of a register.

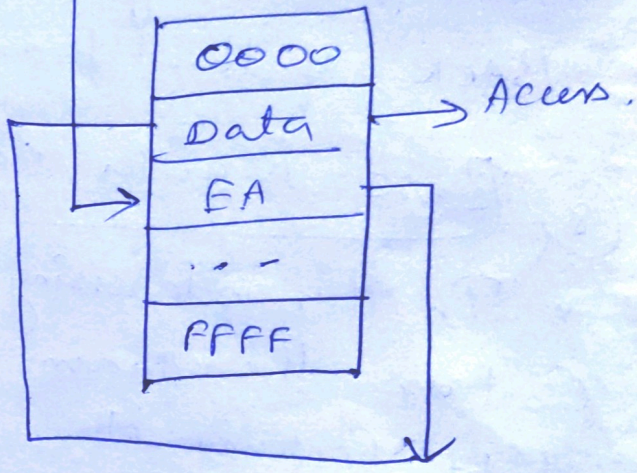
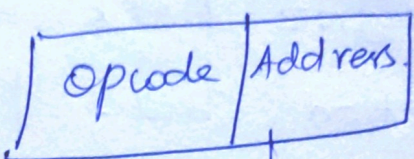
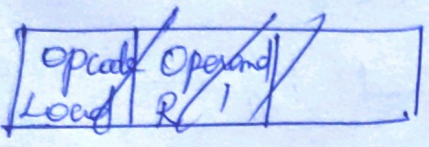
```
MOV 20(R1), R2
```

The above instruction loads the contents of register R2 into the memory location whose address is



Calculated by addition ~~of~~ of the contents of register R<sub>1</sub> and constant value (offset or displacement) 20.

④ Indirect addressing





## Logical data:

→ most of the Processors interpret data as a bit, byte, word or double word. These are referred to as units of data.

→ When data is viewed as  $n$ -bit items of data, each item having the value 0 or 1, it is considered as a logical data.

Logical operation. → logical operators are useful for extracting and inserting groups of bits in a word and it needs to operate on fields of bits within a word or even on individual bits.

Logical operations	C operators	Java operators	MIPS instruction
Shift left	$\ll$	$\ll$	sll
Shift right	$\gg$	$\gg$	srl
Bit by bit AND	$\&$	$\&$	and, andi
Bit by bit OR	$ $	$ $	or, ori
Bit by Bit NOT	$\sim$	$\sim$	nor







MIPS instruction set includes NOR (NOT OR) instruction instead of NOT instruction.

$\text{nor } \$t_0 \ \$t_1, \$t_2.$

$\$t_1$ :	0000	0000	0000	0000	0011	1100	0000	0000
$\$t_2$ :	0000	0000	0000	0000	0000	0000	0000	0000
$\$t_0$ :	1111	1111	1111	1111	1100	0011	1111	1111

$\text{nor } \$t_0 \ \$t_1, \$t_2 \equiv \# \text{ reg } \$t_0 = \text{N}(\text{reg } \$t_1 \mid \text{reg } \$t_2)$

### Shift operation.

→ Another class of operations is called shifts. They move all the bits in a word to the left or right, filling the emptied bits with 0's

→ The actual name of two MIPS shift operations is called shift left logical (sll) and shift right logical (srl).

Example If register  $\$s_0$  contained 9<sub>ten</sub> and the instruction shift left by 4. the new value would be 144.  
left shift by 4.

0000	0000	0000	0000	0000	0000	0000	1001	= 9 <sub>ten</sub>
0000	0000	0000	0000	0000	0000	1001	0000	=
								144 <sub>ten</sub> .



## Control Operations.

→ Decision making instructions alter the control flow i.e.) change the "next" instruction to be executed, like if statement with a go to statement.

→ MIPS assembly language includes two decision making instructions like if else and go to called Conditional branch instruction.

### Decision Making Instruction.

\* The first decision making instruction is

`beq register1, register2, L1;`

branch if equal statement state that if value of register 1 is equal to the value of register 2, then processor can change its sequence to L1. otherwise it execute the next instruction in sequence.

\* Second decision making instruction is

`bne register1, register2, L1;`

branch if not equal statement - state that if value of register 1 is not equal to the value of register 2. then processor can change its sequence to L1. otherwise it execute the next instruction in sequence.



\* Another kind of branch instruction is Unconditional branch instruction. This instruction says that the Processor always follows the branch without testing for any condition.

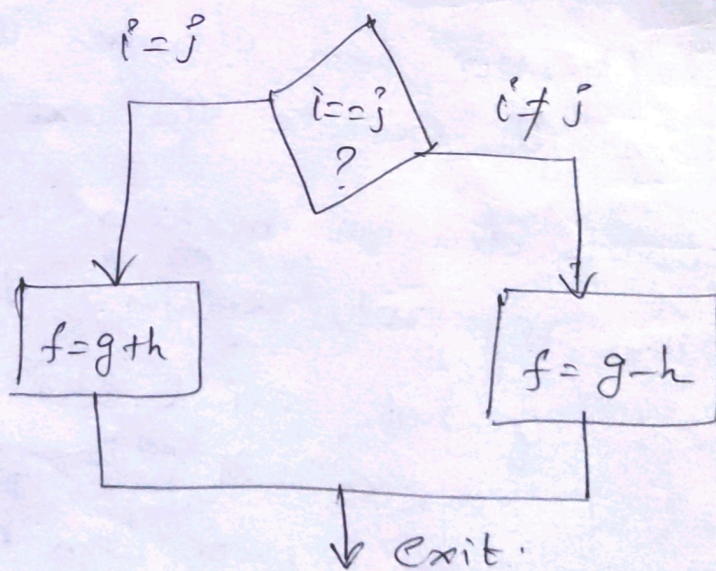
"J" target address  $\rightarrow$  Jump - Instruction allow the Processor to target address.

Compiling if then-else into Conditional Branches.

Example: Compiling 'C' if into MIPS.

In the following code segment  $f, g, h, i$  and  $j$  are variables. If the five variables  $f$  through  $j$  corresponds to the five register  $\$s_0$  through  $\$s_4$

If  $(i == j)$   
 $f = g + h;$   
 else  
 $f = g - h;$



MIPS Code:

first step is Compare for equality

$(i == j)$  if equal to True.

beq  $\$s_3, \$s_4, \text{True};$

sub  $\$s_0, \$s_1, \$s_2;$

j Exit

# branch  $i == j$

#  $f = g - h$  (false)

# go to exit



True: add \$S0, \$S1, \$S2; # f = g+h (true)  
Exit

Write a program to calculate the arithmetic statement:  $Y = (A+B) * (C+D)$  Using three address, two address, one address and zero address.

soln: Using three address instructions.

ADD R1, A, B ;  $R_1 \leftarrow M[A] + M[B]$ .  
ADD R2, C, D ;  $R_2 \leftarrow M[C] + M[D]$ .  
MUL Y, R1, R2 ;  $M[Y] \leftarrow R_1 * R_2$ .

Using two address instruction.

MOV R1, A ;  $R_1 \leftarrow M[A]$ .  
ADD R1, B ;  $R_1 \leftarrow R_1 + M[B]$ .  
MOV R2, C ;  $R_2 \leftarrow M[C]$ .  
ADD R2, D ;  $R_2 \leftarrow R_2 + M[D]$ .  
MUL R1, R2 ;  $R_1 \leftarrow R_1 * R_2$ .  
MOV Y, R1 ;  $M[Y] \leftarrow R_1$ .

Using one address instruction.

LOAD A ;  $AC \rightarrow M[A]$ .  
ADD B ;  $AC \rightarrow AC + M[B]$ .  
STORE T ;  $M[T] \leftarrow AC$ .  
LOAD C ;  $AC \leftarrow M[C]$ .  
ADD D ;  $AC \leftarrow AC + M[D]$ .  
MUL T ;  $AC \leftarrow AC * M[T]$ .



STORE Y ;  $M[Y] \leftarrow AC$ .

Using Zero address Instructions.

PUSH A ;  $TOS \leftarrow A$

PUSH B ;  $TOS \leftarrow B$

ADD ;  $TOS \leftarrow (A+B)$  - Zero address instruction

PUSH C ;  $TOS \leftarrow C$ .

PUSH D ;  $TOS \leftarrow D$ .

ADD ;  $TOS \leftarrow (C+D)$  - Zero address instruction

MUL ;  $TOS \leftarrow (C+D) * (A+B)$

POP Y ;  $M[Y] \leftarrow TOS$ .

In the following code segment a, b, c, d and e are variable. If the five variable f through j correspond to the five registers  $\$S_0$  through  $\$S_4$ , what is the compiled MIPS code for C code.

Soln     $bne \$S_0, \$S_1, Else$     // go to else if  $a \neq b$ .  
           $add \$S_2, \$S_3, \$S_4$     //  $c = d + e$  (skipped if  $a \neq b$ )  
           $j exit$                 // go to exit  
           $Sub \$S_2, \$S_3, \$S_4$     //  $c = d - e$  (skipped if  $a = b$ )



# UNIT-II

## ARITHMETIC

①

### Introduction.

- \* The Arithmetic and Logic Unit (ALU) is responsible for performing arithmetic operations such as add, subtract, division and multiplication and logical operations such as AND, OR Inverting etc
- \* Based on number system two basic data types are implemented in the Computer System.
  1. fixed point numbers.
  2. floating point numbers.

### Binary Addition.

When adding two numbers if the sum of the digits in a given position equals or exceeds, then a carry is propagated.

eg.  $5_{10} + 6_{10}$

0000	0000	0000	0000	0000	0000	0000	0101
0000	0000	0000	0000	0000	0000	0000	0110
<hr/>							
0000	0000	0000	0000	0000	0000	0000	1011
<hr/>							

\*(11)<sub>10</sub>

### Binary Subtraction.

Subtraction operation cannot be performed directly.

- ① One can formulate a subtraction algorithm.
- ② One can negate the subtrahend (ie in  $a-b$  the subtrahend is  $b$ ).  
then perform addition (2's Complement).



$$\begin{array}{r} \overset{0}{1} \overset{1}{0} \overset{1}{0} \overset{1}{1} \quad (-) \\ \underline{101} \\ 01110 \end{array}$$

Ans =  $(1110)_2$

$$\begin{array}{l} 0-0 \\ 1-0 \\ 1-1 \\ 10-1 \end{array} = \begin{array}{l} 0 \\ 1 \\ 0 \\ 1 \end{array}$$

$12_{10} - 5_{10} = 7_{10}$  (Direct method)

$$\begin{array}{r} 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 1100 \\ 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0101 \\ \hline 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0111 \end{array} = (7_{10})$$

$$\begin{array}{r} 10 \\ -1 \\ \hline 1 \end{array}$$

Unsigned binary representation

Addition with 2's Complement

$$\begin{array}{r} 5 = 0101 \\ -5 = 1010 \\ \hline 1011 \end{array}$$

$$\begin{array}{r} 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 1100 \quad (12) \\ 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1111 \quad 1011 \quad (-5) \\ \hline 10000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0000 \quad 0111 \end{array}$$

1

Discard the Carry

Fixed Point addition and Subtraction.

Two's Complement addition and Subtraction.

Subtraction can be done in following ways.

- ① Direct method.
- ② Using 1's Complement
- ③ 2's Complement



Problem Add the numbers  $(0.5)_{10}$  and  $(0.4375)_{10}$  using floating point-addition.

$$(0.5)_{10} = \frac{1}{2} = 0.5 \times 2 = 1.0 \Rightarrow (0.1)_2$$

Normalizing the above value.

$$= 1.0 \times 2^{-1}$$

$$\begin{aligned} 0.4375 &\Rightarrow 0.4375 \times 2 = 0.8750 \\ &0.8750 \times 2 = 1.7500 \\ &0.7500 \times 2 = 1.5000 \\ &1.5000 \times 2 = 1.0000 \end{aligned}$$

$$\Rightarrow (0.0111)_2$$

Normalizing the above Value =  $1.110 \times 10^{-2}$

Step 1: The significant of the number with lesser exponent ( $1.11 \times 10^{-2}$ ) is shifted right until its exponent matches the larger number

$$1.11 \times 10^{-2} = 0.111 \times 2^{-1}$$

Step 2: Add the significant

$$\begin{array}{r} 1.0 \times 2^{-1} \\ 0.111 \times 2^{-1} \\ \hline \end{array}$$

$$1.111 \times 2^{-1}$$

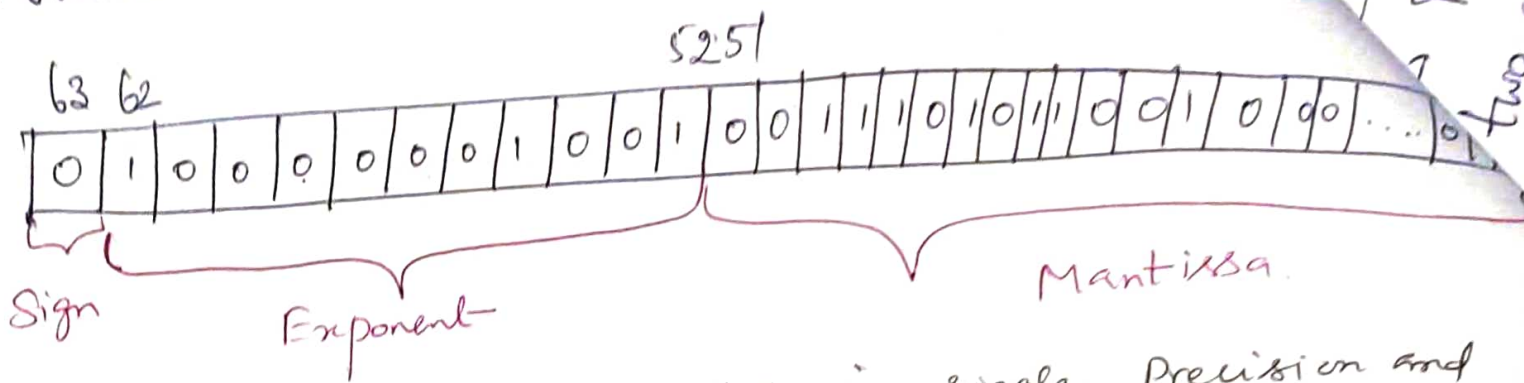
$$\text{Ans } 1.111 \times 2^{-1}$$

Eg. (2) Add the numbers  $0.5_{10}$  and  $-0.4375_{10}$

$$\text{Ans} = 1.000 \times 2^{-4} = 0.0625$$



# Number in double precision format - 8



Eg. 2 represent  $-307.1875_{10}$  in single precision and double precision formats.

## Floating Point Arithmetic

### Addition and Subtraction.

Consider two floating point numbers.

$$A = m_1 r^{e_1} \text{ and}$$

$$B = m_2 r^{e_2} \text{ Assume } e_1 \geq e_2$$

### Procedure.

Step 1: select the number with a smaller exponent and shift its mantissa right a number of steps equal to the difference in exponents  $|e_2 - e_1|$ .  
for example if numbers are  $1.75 \times 10^2$  and  $6.8 \times 10^4$ .  
then the number  $1.75 \times 10^2$  is selected and converted to  $0.0175 \times 10^4$ .

Step 2: Set the exponent of the result equal to larger exponent.

Step 3: Perform addition/subtraction on mantissa and determine the sign of the result.

Step 4: Normalize the result if necessary.



## Fixed Point Number Representation (2)

→ Fixed point/integer numbers are represented in two forms: signed integer and unsigned integer.

→ To represent negative numbers various techniques are used because computer does not have provision to represent negative sign. Techniques to represent signed integer numbers are.

- ① Sign magnitude representation
- ② 1's complement
- ③ 2's complement

## Signed Magnitude Representation

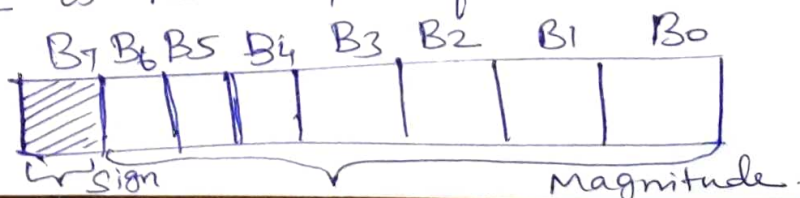
→ left most bit (sign bit) in the binary number represents sign of the number.

→ Here the most significant bit (MSB) represents sign of the number. If MSB is 1 number is negative and if MSB is 0 number is positive. The remaining bits represent magnitude of the number.

→ examples.

+6	=	0000	0110
-14	=	1000	1110
-64	=	1100	0000

→ In case of unsigned 8 bit binary numbers the decimal range is 0 to 255. For signed magnitude 8-bit binary numbers the largest magnitude is reduced from 255 to 127.





## 1's Complement Representation

The 1's Complement of a binary number is the binary number that results when we change all 1's to 0's and 0's to 1's.

## 2's Complement Representation

→ The 2's Complement is the binary number that results when we add 1 to the 1's Complement. It is given as

$$2's \text{ Complement} = 1's \text{ Complement} + 1$$

2's Complement form is used to represent negative numbers.

## Fixed point Addition and Subtraction

→ We can relate addition and subtraction operations of numbers by the following relationship.

$$(\pm A) - (+B) = (\pm A) + (-B)$$

$$(\pm A) - (-B) = (\pm A) + (+B)$$

Therefore we can change subtraction operation to an addition operation by changing the sign of the subtrahend.

## Addition

### Rules for addition

$$0+0=0$$

$$0+1=1$$

$$1+0=1$$

$$1+1=10_2$$

Example Add  $(6)_{10}$  to  $(7)_{10}$   
in binary

$$\begin{array}{r} \phantom{0}1 \phantom{0}1 \phantom{0}1 \phantom{0}0 \\ 0110 \\ 0111 \\ \hline 1101 \\ \hline \end{array} \quad \begin{array}{l} 6_{10} \\ 7_{10} \\ 13_{10} \end{array}$$



all so

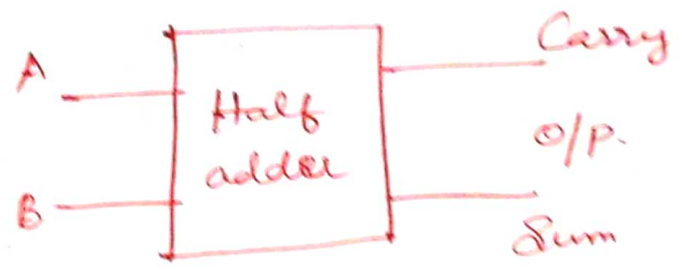
The first three operations produce a sum whose length is one digit but when the last operation is performed sum is two digits. The higher significant bit of this result is called a carry and lower significant bit is called sum.

→ The logic circuit which performs this operation is called a half-adder. The circuit which performs addition of three bits

Half adder

→ The half adder operation needs two binary input: augend and addend bits and two binary output: sum and carry.

Input		Output	
A	B	Carry	Sum
0	0	0	0
0	1	0	1
1	0	0	1
1	1	1	0



Block Schematic of half adder.

Limitations of half adder.

→ In multidigit addition we have to add two bits along with the carry of previous digit addition.

→ Effectively such addition requires addition of three bits. This is not possible with half adder hence half-adders are not used in practice.



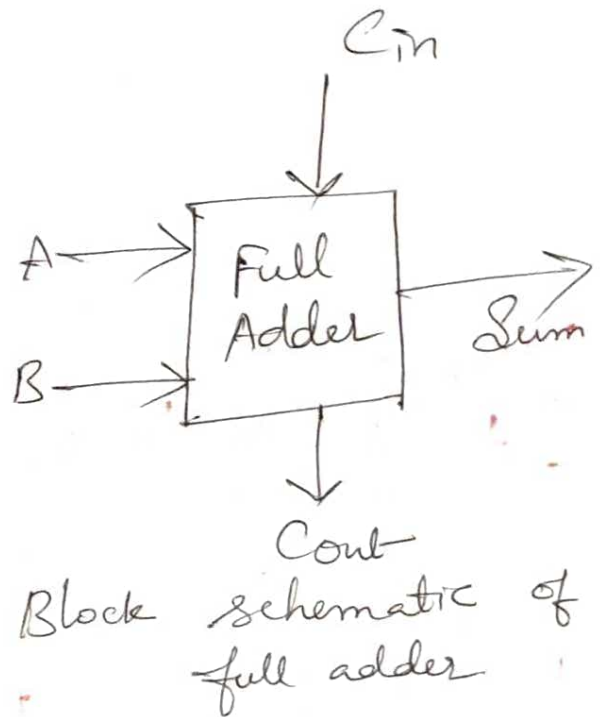
## Full adder

→ A full-adder is a Combinational Circuit that forms the arithmetic sum of three Input bits. It consists of three inputs and two outputs.

→ Two of the input variables, denoted by A and B, represent the two significant bits to be added. The third input  $C_{in}$  represents the carry from the previous lower significant position.

Inputs			Outputs	
A	B	$C_{in}$	Carry	Sum
0	0	0	0	0
0	0	1	0	1
0	1	0	0	1
0	1	1	1	0
1	0	0	0	1
1	0	1	1	0
1	1	0	1	0
1	1	1	1	1

Truth table for full adder.



## Binary Subtraction

→ The subtraction consists of four possible elementary operations.

→ In all operations, each subtrahend bit is subtracted from the minuend bit. In case of second operation the minuend bit is smaller than the subtrahend bit hence 1 is borrowed.



Outputs

Result

### Subtraction Method

(4)

- ① Subtract bits column-wise starting from LSB with borrow if any
- ② Put the difference at the bottom of same column
- ③ Take the borrow, if required from the next column.

### Binary Subtraction Using 1's Complement Method.

→ In a 1's Complement Subtraction, negative number is represented in the 1's Complement form and actual addition is performed to get the desired result. For example operation  $A - B$  is performed using following steps.

- ① Take 1's Complement of B.
- ② Result  $\leftarrow A + 1's \text{ Complement of } B$ .
- ③ If carry is generated then the result is positive and in the true form. Add carry to the result to get the final result.
- ④ If carry is not generated then the result is negative and in the 1's Complement form.

eg. Perform  $(15)_{10} - (28)_{10}$ . Using 6 bit 1's Complement representation.

$$(15)_{10} = (001111)_2 \quad (28)_{10} = 011100$$

$$1's \text{ Complement of } (15)_{10} = 110000$$



$$\begin{array}{r}
 1 \quad 1 \quad \xrightarrow{\text{Carry}} \\
 011100 \quad (28) \\
 110000 \quad (15) \\
 \hline
 \end{array}$$

$$\begin{array}{r}
 \hline
 1001100 \\
 \hline
 \end{array}
 \quad \text{Result}$$

$$\begin{array}{r}
 \xrightarrow{1} \\
 \hline
 001101 \\
 \hline
 \end{array}
 \quad \text{final result } (13)_{10}$$

$$28 = 011100$$

$$1's \text{ Complement of } 28 = 100011$$

$$\text{binary eq. of } (15)_{10} = 001111$$

$$1's \text{ Comp. of } (28)_{10} = 100011$$

$$\begin{array}{r}
 100011 \\
 001111 \\
 \hline
 110010 \quad \text{Result} =
 \end{array}$$

$$\begin{array}{r}
 (15)_{10} \\
 (-28)_{10} \\
 \hline
 (-13)_{10}
 \end{array}$$

$$\text{Verification: } \boxed{001101} = (13)_{10}$$

Binary Equivalent of  $(-13)_{10}$

eg. Perform  $(-4)_{10} - (8)_{10}$ . Using 1's Complement form

$$(-4)_{10} - (-8)_{10} = (-4)_{10} + (8)_{10}$$

$$(4)_{10} = 0100$$

$$(-4)_{10} = 1011$$

1's Complement of 4.

$$\begin{array}{r}
 \boxed{1}1011 \\
 01000 \\
 \hline
 \boxed{1}00011 \\
 \hline
 \end{array}$$

Sign Extension

$$\Rightarrow \boxed{00100} \text{ Result}$$



## by Subtraction <sup>(5)</sup> Using 2's Complement Method

→ In a 2's Complement subtraction negative number is represented in the 2's Complement form and actual addition is performed to get the desired result.

for example, operation.  $A - B$  is performed

Using following steps.

(1) Take 2's Complement of B.

(2) Result  $\leftarrow A + 2$ 's Complement of B.

(3) If carry is generated then result is positive and in the true form. In this case carry is ignored.

(4) If carry is not generated then the result is negative.

eg. Perform  $(28)_{10} - (15)_{10}$  Using 6-bit 2's Complement representation.

$$(15)_{10} = 001111$$

$$\text{2's Complement of } (15)_{10} = 110000$$

$$\begin{array}{r} 110000 \\ \hline 110001 \end{array}$$

$$011100 = (28)_{10}$$

$$110001 = 2$$
's Comp of 15

$$\begin{array}{r} 011100 \\ \hline \times 110001 \\ \hline 001101 \end{array} \Rightarrow \text{Result: Binary equivalent of } (13)_{10}$$



Pbm

Subtract  $(11010)_2 - (10000)_2$   
Complement and 2's Complement

Using 1's  
method.

$$\begin{array}{r} 10000 \\ 01111 \rightarrow \text{1's Complement} \\ \hline 10000 \rightarrow \text{2's Complement} \end{array}$$

Using 1's Complement method

$$\begin{array}{r} 11010 \\ 01111 \\ \hline 101001 \\ \rightarrow 1 \\ \hline 101010 \end{array}$$

Result is +ve.

Using 2's Complement Method

$$\begin{array}{r} 11010 \\ 10000 \\ \hline \times 01010 \\ \hline \downarrow \\ \text{ignore Carry.} \end{array}$$

Result is +ve.



## Point Multiplication (6)

→ Multiplication of fixed point numbers can be accomplished with addition, subtraction and shift operations.

### Unsigned Multiplication

→ Multiplication of Unsigned binary integers is handled similarly to the way it is carried out by hand for decimal numbers.

→ In the binary system the partial products are easily defined when the multiplier bit is 0, the partial product is 0 and when the multiplier is 1 the partial product is the multiplicand.

→ The final product is produced by summing the partial products before carrying operation. Each successive partial product is shifted one position to the left relative to the preceding partial product.

Multiplicand = 1101 (13)

Multiplier = 1001 (9)

$$\begin{array}{r} \phantom{11}1101 \\ 0000 \\ 0000 \\ 1101 \\ \hline 1110101 = 117 \end{array}$$

} partial product

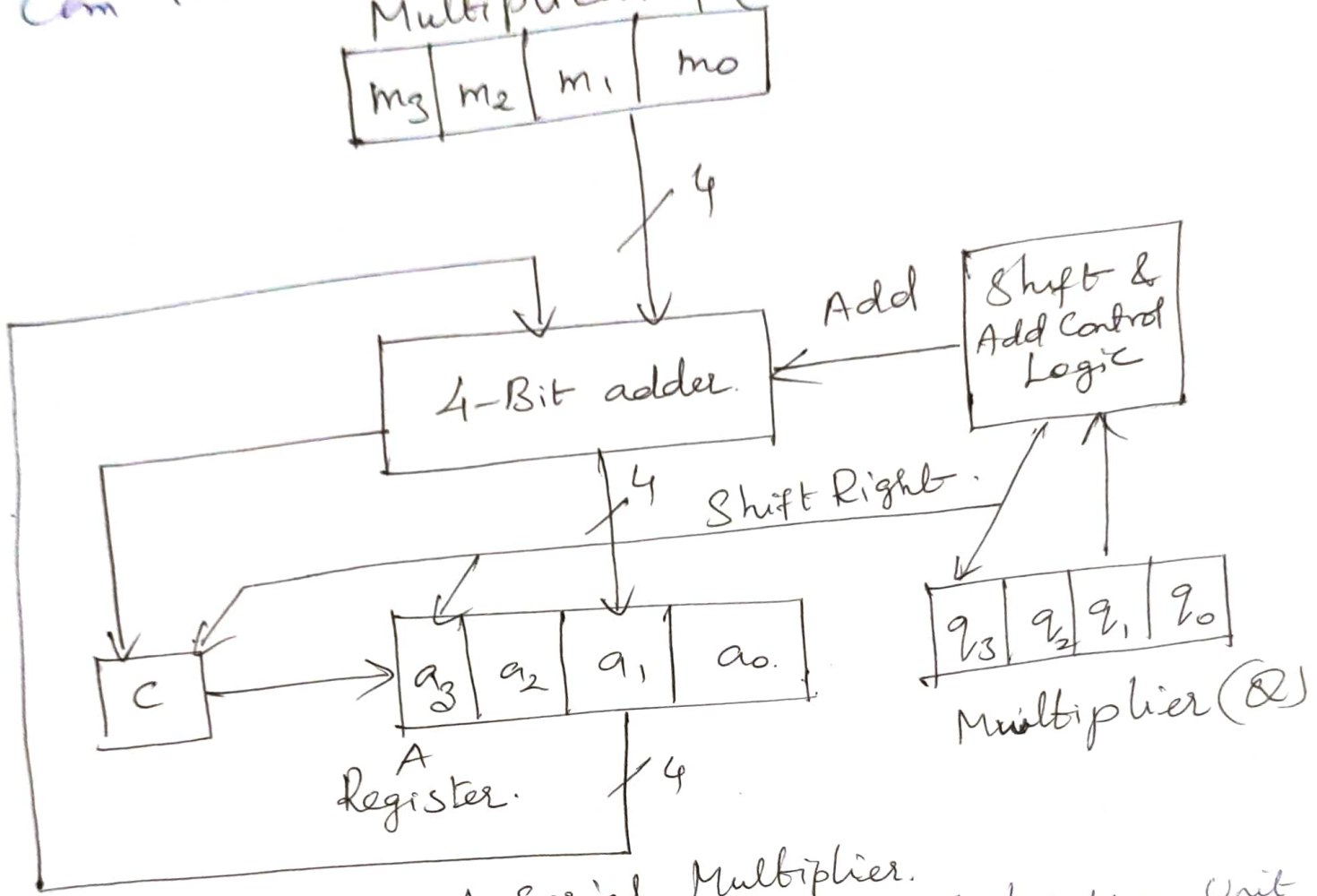
→ When two unsigned  $n$ -bit numbers are multiplied the result can be large as  $2n$  bits.

→ for example two four bit operands result in an eight bit product.



# Hardware Implementation of Integer Multiplication

The hardware implementation of integer multiplication can take a similar form to the manual method.



A Serial Multiplier.

→ It shows a layout of a multiplication unit for four bit numbers in which there is a four bit adder, a Control Unit, three four bit registers and one bit Carry register.

Steps: ① In order to multiply two numbers, the multiplicand is placed in the M-register, the multiplier is placed in the Q-register and the A and C registers are cleared to zero.



al method

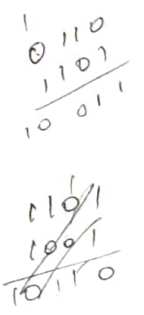
Using multiplication the right most bit of the multiplier determines whether the multiplicand is added into the product at each step.

③ After the multiplicand is added into the product the multiplier and A registers are simultaneously shifted to the right. This has the effect of shifting the multiplicand to the left and exposing the next bit of the multiplier in position  $q_0$ .

Example Multiplicand (M)

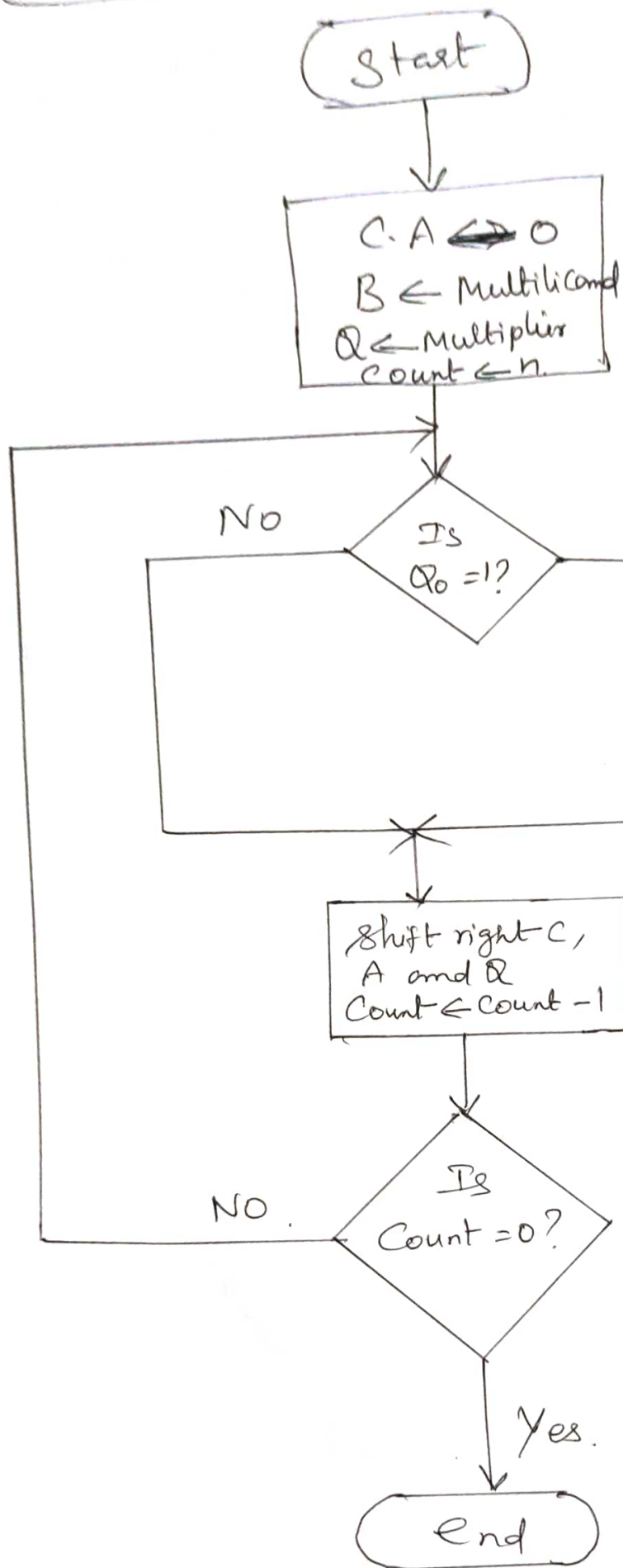
Consider 4-bit multiplier and Multiplicand  
 Multiplicand = 1101 (B) and Multiplier = 1011 (B)

C	A	Q	
0	0000	1011	Add M to A Shift
0	1101	1011	
0	0110	1101	Add M to A Shift
1	0011	1101	
0	1001	1110	No add only shift
0	0100	1111	
1	0001	1111	Add m to A Shift
0	1000	1111	





# Multiplication Operation Steps



→ Bit 0 of multiplier operand (Q<sub>0</sub> of Q register) checked

→ If bit 0 (Q<sub>0</sub>) is ~~the~~ One then multiplicand and partial product are added and all bits of C, A and Q register are shifted to right one bit. so that C bit goes into A<sub>n-1</sub>.

→ A<sub>0</sub> goes into Q<sub>n-1</sub> and Q<sub>0</sub> is lost. If bit 0 (Q<sub>0</sub>) is 0. then no addition is performed. Only Shift Operation is Carried out.

→ Step 1 and 2 are repeated n times to get the desired result in the A and Q reg.

Flow chart for multiplication operation.



Initial Version of the Multiplication Algorithm (8)  
Hardware

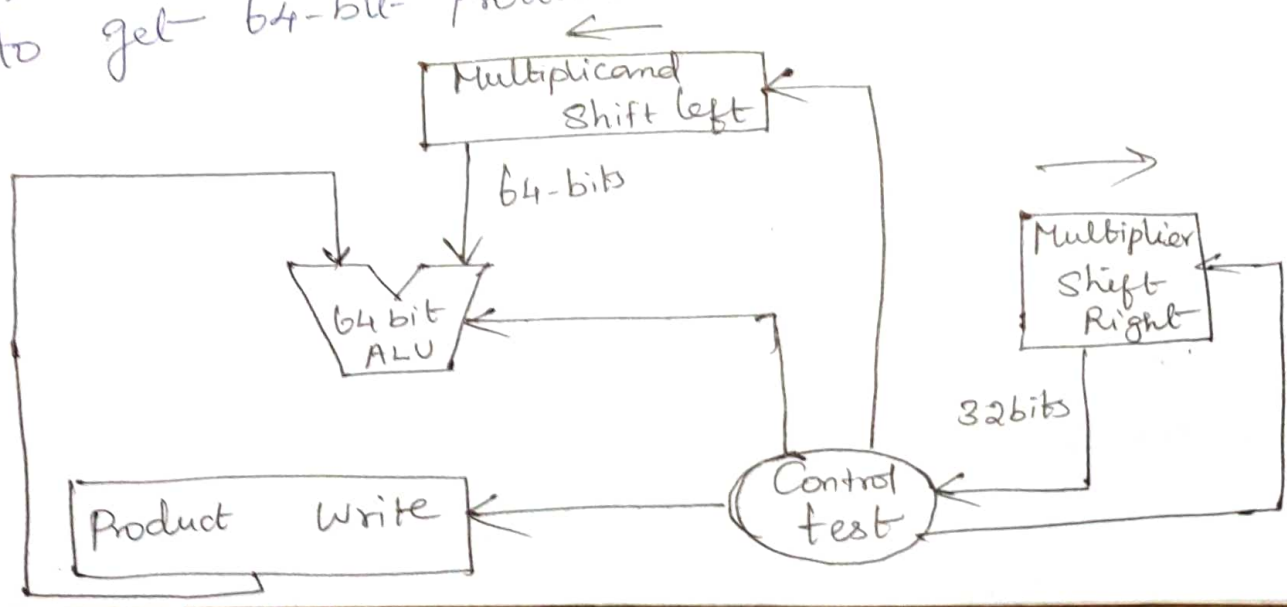
Step 1: The LSB bit of the Multiplier 0 determines whether the multiplicand is added to the Product register.

Step 2: Shift Multiplicand by 1 bit left logically that will have a effect of moving the intermediate operands to the left.

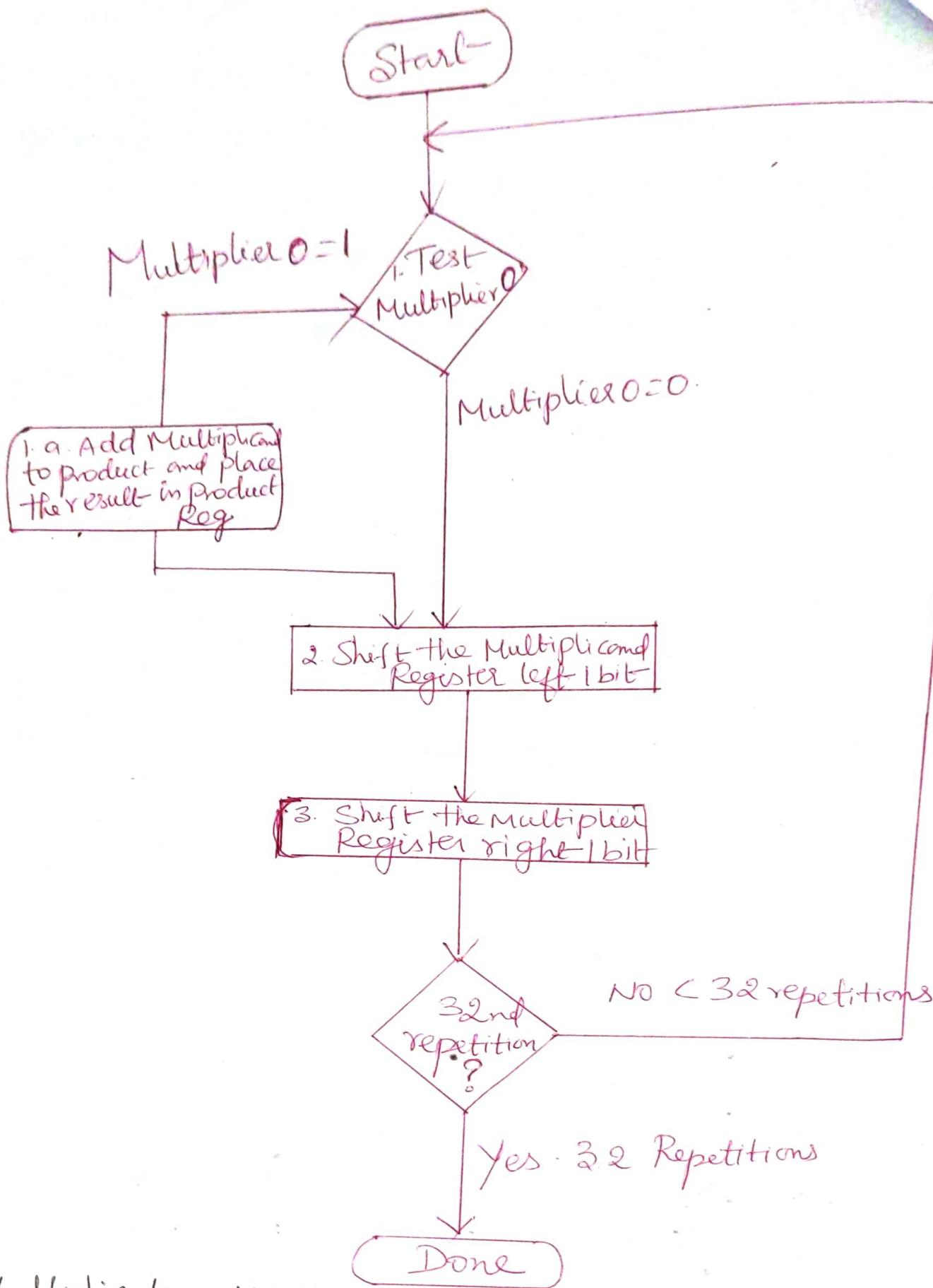
Step 3: Shift the Multiplier right logically by 1 bit that will ~~have~~ give the next bit of the multiplier to be examined.

Multiplication Hardware

- \* It is Constructed with 64 bit multiplicand, ALU and Product register and one 32 bit multiplier.
- \* The 32 bit multiplier is loaded in 32 bit multiplier register, 32 bit multiplicand is loaded in right half of 64 bit multiplicand register and zero in the left half and the 64 bit Product register is initialized to 0.
- \* The above Said Procedure is executed for 32 times to get 64-bit Product result.







Multiplication Hardware flow chart.



For Simplicity assume Use 4-bit Multiply  
 $4_{10} \times 3_{10} = 12_{10}$  or  $0100_2 \times 0011_2$  Same Concept Can be  
 Extended to 32 bit size number.

Iteration.	Steps.	Multiplier	Multiplicand	Product
0	Initial values	00 <u>11</u>	0000 0100	0000 0000
1	1. a. 1 $\rightarrow$ Prod = Prod + M <sub>can</sub> d	0011	0000 0100	0000 0100
	2. Shift left M <sub>can</sub> d.	0011	0000 1000	0000 0100
	3. Shift right Multiplier	00 <u>01</u>	0000 1000	0000 0100
2	1. a. 1 $\rightarrow$ Prod = Prod + M <sub>can</sub> d	0001	0000 1000	0000 1100
	2. Shift left M <sub>can</sub> d	0001	0001 0000	0000 1100
	3. Shift Right multiplier	00 <u>00</u>	0001 0000	0000 1100
3	1. 0 $\rightarrow$ No-operation.	0000	0001 0000	0000 1100
	2. Shift left M <sub>can</sub> d	0000	0010 0000	0000 1100
	3. Shift Right Multiplier	00 <u>00</u>	0010 0000	0000 1100.
4	1. 0 $\rightarrow$ NO operation	0000	0010 0000	0000 1100
	2. $\rightarrow$ Shift left multiplier	0000	0100 0000	0000 1100
	3. Shift right Multiplier	0000	0100 0000	0000 1100



## Binary Division.

(9a)

→ Division is a Reciprocal Operation of multiplication. It is similar to the decimal numbers.

→ In Division Process first the bits of dividend are examined from left to right, until the set of bits examined represent a number greater than or equal to divisor.

→ until this Condition occurs '0's placed in the quotient from left to right.

→ When the Condition is satisfied a '1' is placed in the quotient and divisor is subtracted from the partial dividend. The result is referred to as a partial remainder.

→ from this point onwards the division process is repeated. In each repetition cycle additional bits from the dividend are brought down to the partial remainder until the result is greater than or equal to the divisor and the divisor is subtracted from the result to produce a new partial remainder.

→ The Process continues until all the bits of the dividend are brought down and result is still less than the divisor.

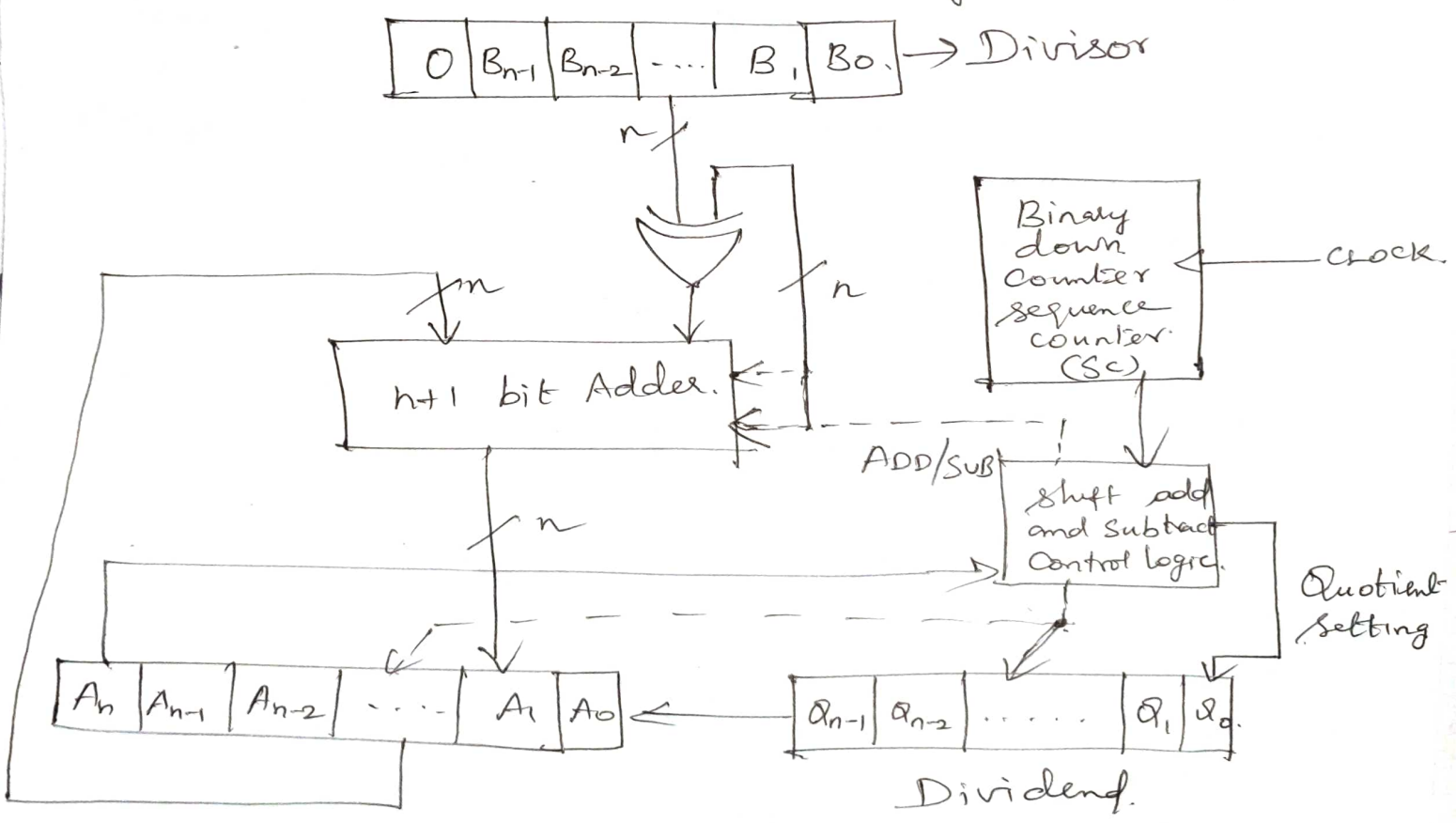


$$\begin{array}{r}
 110 \\
 \overline{) 11011011} \\
 \underline{100100} \\
 000110 \\
 \underline{-110} \\
 00011
 \end{array}$$

← Quotient  
 ← Dividend  
 ← Remainder

Division Example.

Division Algorithm. (Hardware to implement binary Division)





## Division Operation Steps

① b

- ① Shift A and Q left one binary position
- ② Subtract divisor (i.e. add 2's Complement of divisor(B)) from A and place answer back in A ( $A \leftarrow A - B$ )
- ③ If the sign bit of A is 1 Set  $Q_0$  to 0 and add divisor back to A. (that is restore A) otherwise Set  $Q_0$  to 1
- ④ Repeat steps 1, 2, and 3 n times

## Example of Division Process Using the Serial ~~adder~~

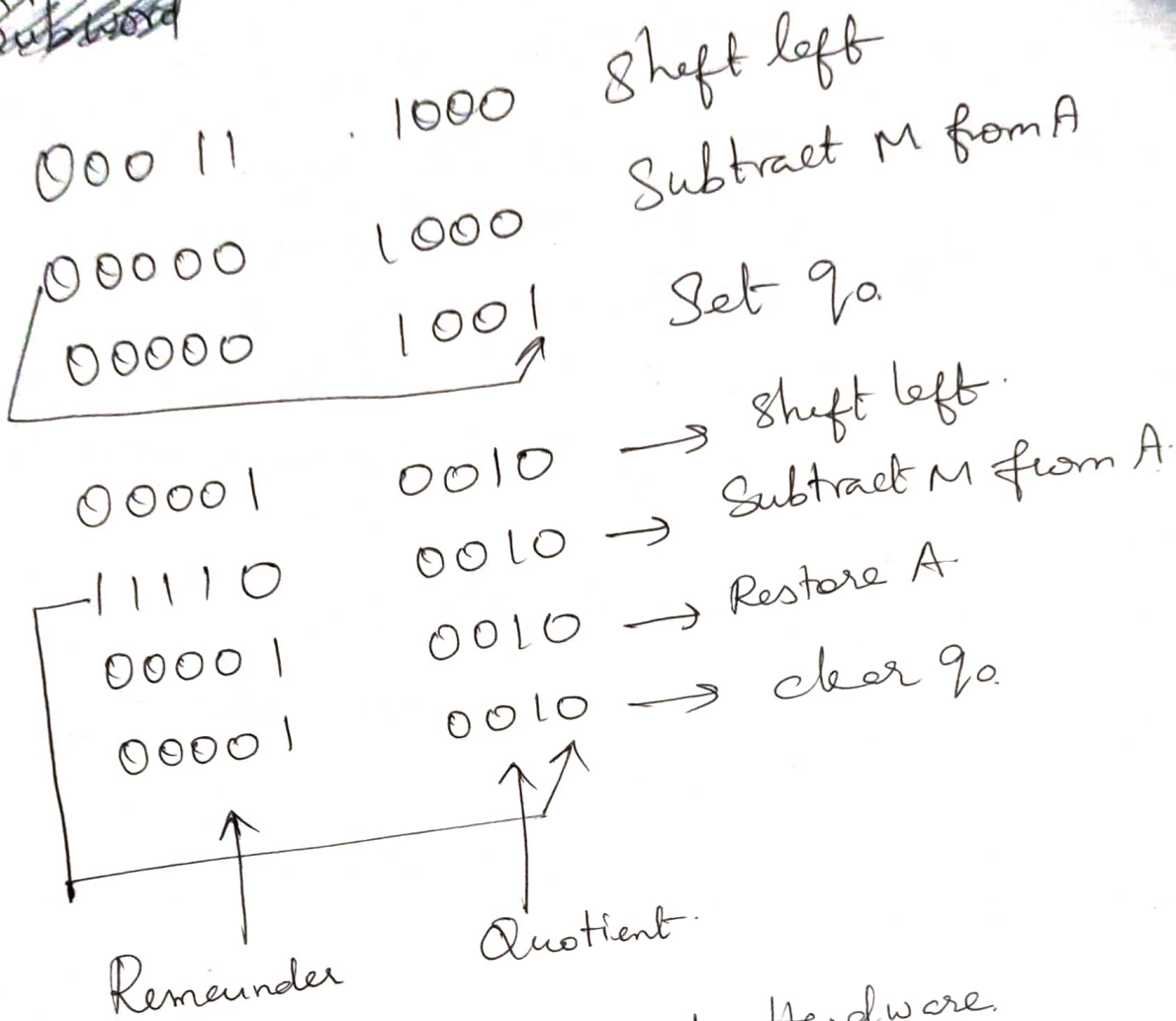
Dividend: 00011  $\rightarrow$  Initial value.

A	Q	
0000	0111	
00000	1110	$\rightarrow$ Shift left
11101	1110	$\rightarrow$ Subtract M from A
00000	1110	$\rightarrow$ Restore A
00000	1110	$\rightarrow$ clear $q_0$ .
00001	1100	$\rightarrow$ Shift left
11110	1100	$\rightarrow$ Subtract M from A.
00001	1100	$\rightarrow$ Restore A.
00001	1100	$\rightarrow$ clear $q_0$ .

$$\begin{array}{r}
 00000 - A \\
 11101 - M \\
 \hline
 11101
 \end{array}$$

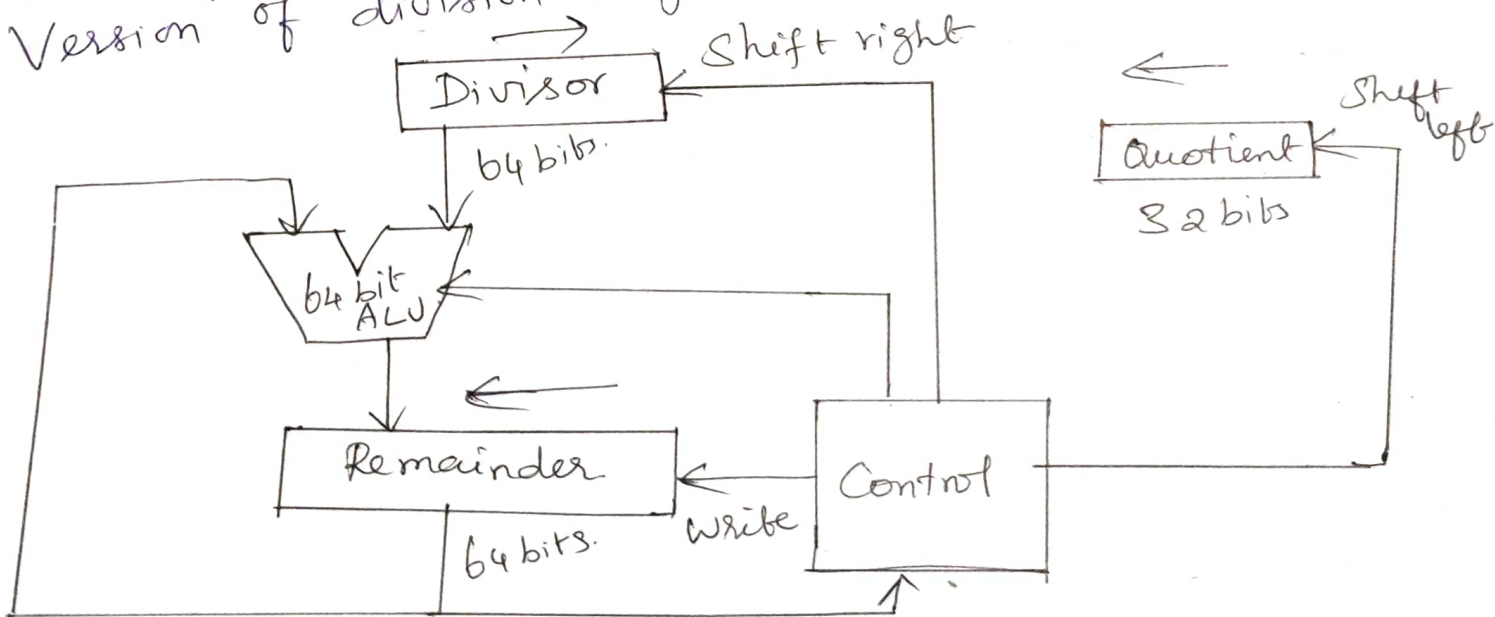


~~Subword~~



A division Algorithm and Hardware.

The following assumptions are made in the first version of division algorithm.



Division Hardware.



(9.6)

Divisor register, ALU and Remainder register are all 64 bits wide. With only the quotient register being 32 bits.

- \* The 32 bits divisor starts in the <sup>left</sup> half of the divisor register and is shifted right  $\wedge$  1 bit each iteration.
- \* 32 bit quotient register is set to 0.
- \* The remainder is initialized with the dividend.
- \* Control decides when to shift the divisor and quotient and when to write the new value into the remainder register.



## Floating Point addition.

Most computers use dedicated hardware for floating point operations as fast as possible.

① Align the decimal number of the smallest exponent number that matches the larger exponent. (Shift the smaller number right until exponent match). (i.e) exponents of the operands are made equal.

② Add/Subtract the significands depending on sign.

③ Normalise the sum by adjusting exponent.

④ Check for overflow

⑤ Rounding off to appropriate numbers of bits

⑥ Result may need further normalization; then go to step 3.

⑦ Set the sign, sign is determined by addition of the significand.

Example: Perform  $9.999_{\text{ten}} \times 10^1 + 1.610_{\text{ten}} \times 10^0$

Assume that we can store only four decimal digits of the significand and two decimal digits of the exponent.

Step 1: Compare the exponent of both the operands

If it equal add the two operand (significand).  
If it is not equal then increase the smaller exponent to the larger exponent.



most  
 10  
 10

$$10 \times 10^{-1} = 0.1610 \times 10^0 = 0.0161 \times 10^1$$

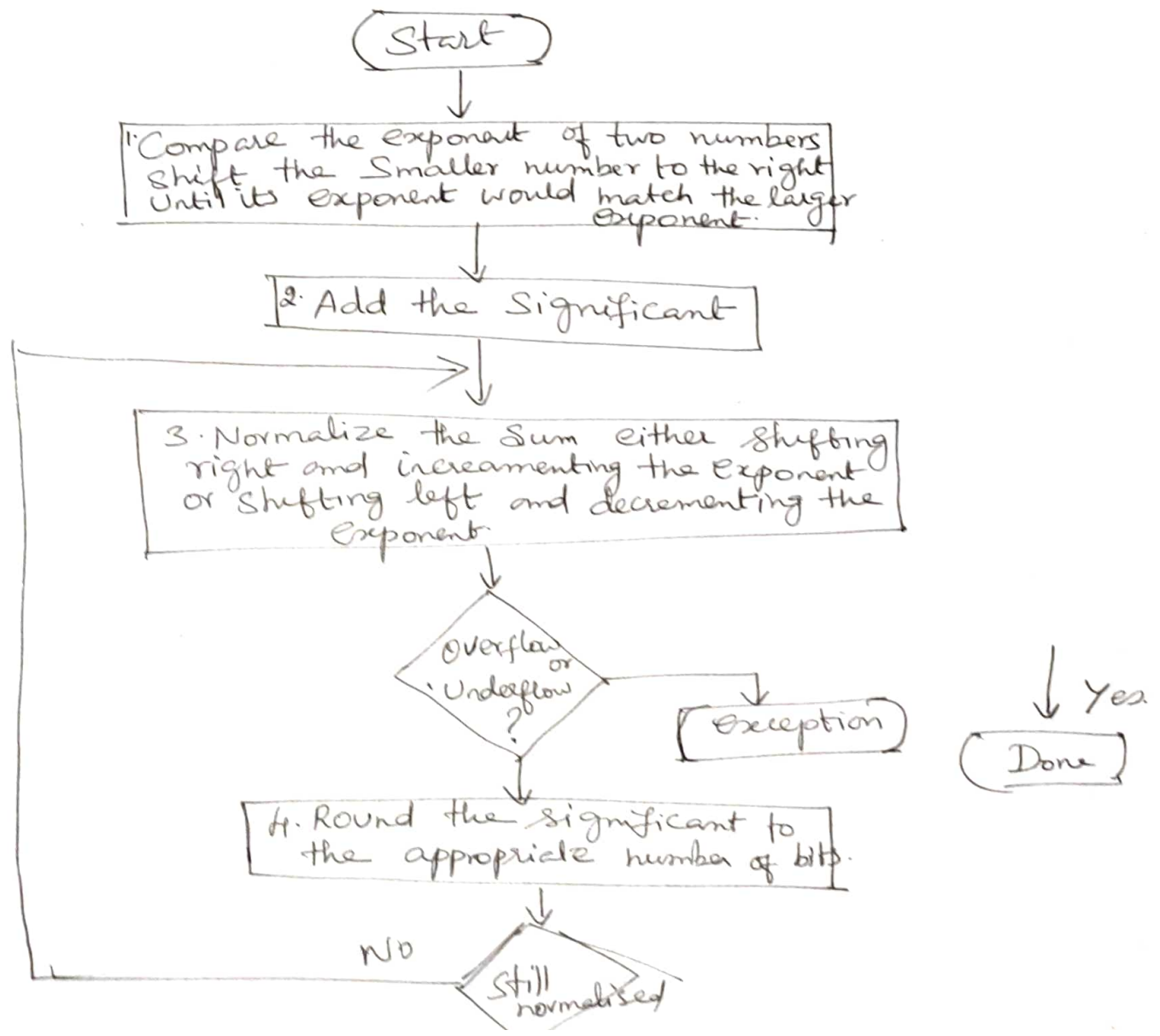
Step 2: Now add the significant.

$$\begin{array}{r} 9.999 \\ 0.016 \\ \hline 10.015 \end{array}$$

Step 3: This sum is not in normalized scientific notation, so we need to adjust it

$$10.015 \times 10^1 = 1.0015 \times 10^2 = 1.0015 \times 10^2$$

Step 4: Round the number to four digits in the significant to  $1.002 \times 10^2$ .





## Subword Parallelism.

→ Many Scientific Problems require operations on large arrays of numbers these numbers are usually formulated as vectors. A vector is an ordered set of a one-dimensional array of data items.

→ It has been observed that many scientific problems and graphics and audio applications perform the same operation of vectors.

→ Suppose the size of the data item within the vector is 128 bits. With this data, by partitioning the carry chains a processor can use parallelism to perform simultaneous operations on short vectors of sixteen 8-bit operands, eight 16-bit operands, four 32-bit operands or two 64-bit operands. Such parallelism which occurs within a wide word is known as subword parallelism.

→ The subword parallelism is also classified under the more general name of data level parallelism. It is also called vector or SIMD (for single instruction multiple data).

→ ARM processors have more than 100 instructions in the NEON multimedia instruction extension to support subword parallelism.



## floating-Point Representation (11)

→ To accommodate very large integers and very small fractions, a computer, must be able to represent numbers and operate on them in such a way that the position of the binary point is variable and is automatically adjusted as computation proceeds.

→ In this case the binary point is said to float and the numbers are called floating point numbers.

→ The floating point representation has three fields: sign, significant digits and exponent.

→ Let us consider the number  $111101.1000110$  to be represented in the floating point format.

→ It is important to  $111101.1000110$  →  $1.11101100110 \times 2^5$   
↑ Exponent 5  
↓ Scaling factor

→ The string of the significant digits is commonly known as mantissa.

In above example.

Mantissa =  $11101100110$ .

Exponent = 5.

## IEEE Standard for floating Point Numbers

→ The standard for representing floating point numbers in 32 bits and 64 bits have been developed by Institute of Electrical and Electronics Engineers (IEEE) referred to as IEEE 754 Standard.

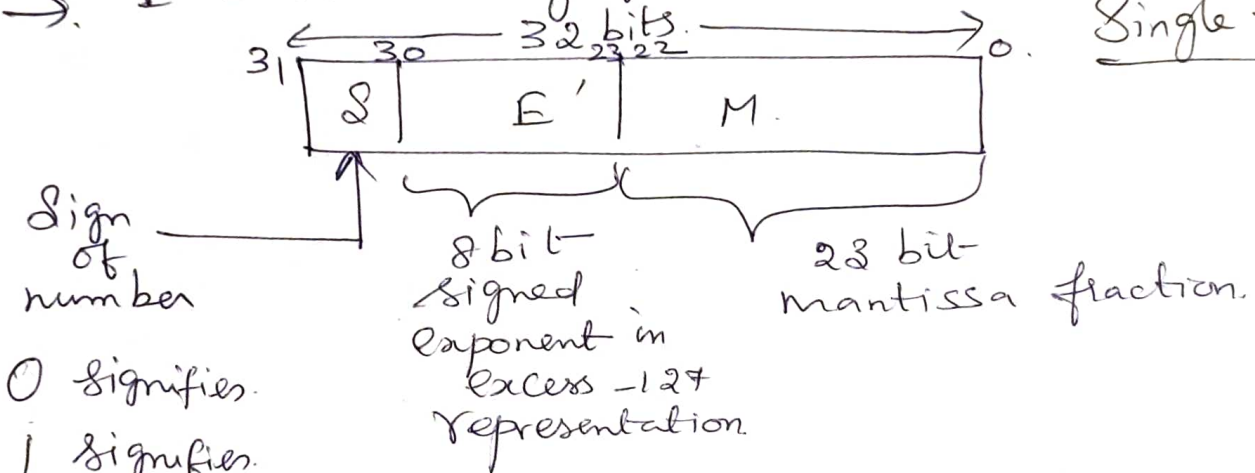


→ 32 bit Standard Representation is single precision. representation. The 32 bits are divided into three fields

- (field 1) Sign → 1-bit
- (field 2) Exponent → 8 bit
- (field 3) Mantissa → 23 bits

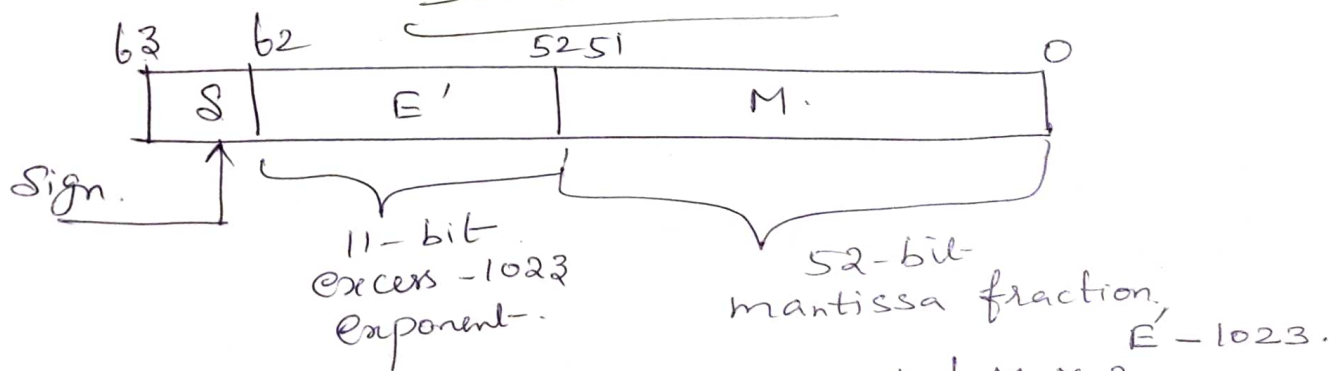
→  $E' = E$  (scaling factor) + bias.

Single-precision



Value Represented =  $\pm 1.M \times 2^{E' - 127}$

Double Precision.



Value Represented =  $\pm 1.M \times 2^{E' - 1023}$

→ In 32 bit floating point system (single precision) bias is 127. Hence  $E' = E$  (scaling factor) + 127. This representation of exponent is also called as the excess 127 format.



32 bits

$(AEB)_{16} \rightarrow \text{Hex. no}$   
 $\downarrow \quad \downarrow \quad \downarrow$   
 $010011101011 \rightarrow \text{Binary number.}$

(12)

1000	- 8
1001	- 9
1010	- 10 A
1011	- 11 B
1100	- 12 C
1101	- 13 D
1110	- 14 E

$(1259)_{10} = (10011101011)_2$

$(0.125)_{10} = (0.001)_2$

Binary number =  $10011101011 + 0.001$   
 $= 10011101011.001$

Step 2: Normalised the number.  $10011101011.001 = 1.0011101011001 \times 2^{10}$

Step 3: Single Precision representation:  
 for a given number  $S = 0, E = 10$ .

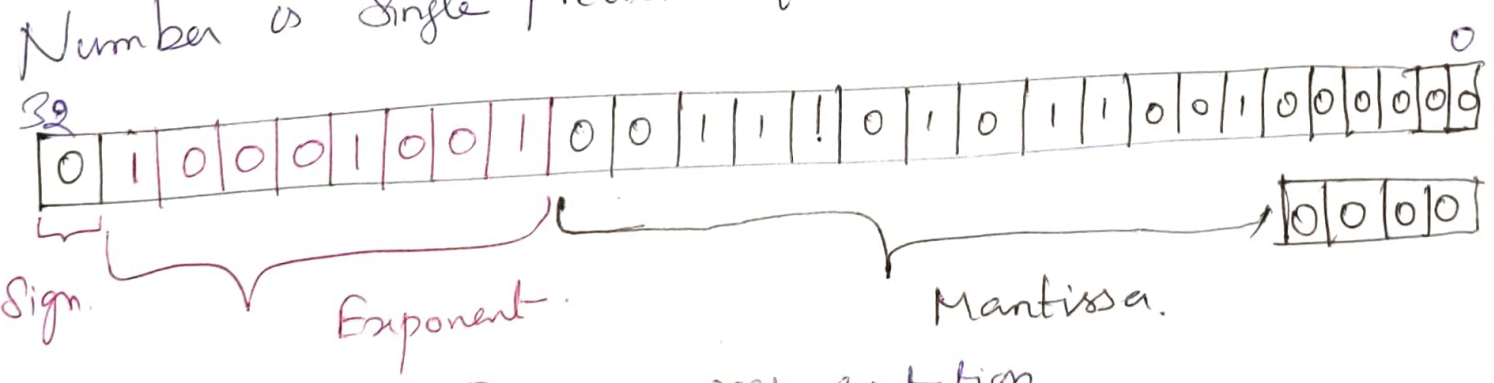
$M = 0011101011001$

Bias for single precision format is 127

$E' = E + 127 = 10 + 127 = 137_{10}$

$= 10001001_2$

Number in single precision format is given as.



Step 4: Double Precision representation  
 for given number.

$S = 0, E = 10$ , and  $M = 0011101011001$

Bias for double precision format is 1023

$E' = E + 1023 = 10 + 1023 = 1033_{10}$   
 $= 10000001001_2$



→ The range of  $E'$  for normal values in Precision is  $0 < E' < 255$ . This means that for representation the actual exponent  $E$  is in the range

$$-126 \leq E \leq 127.$$

→ The 64 bit standard representation is called a double Precision representation because it occupies two 32 bit words.

- 64 bits are divided into three fields:
  - field 1 sign → 1 bit
  - field 2 Exponent → 11 bit
  - field 3 Mantissa → 52 bits

→ In the double Precision format value actually stored in exponent field is given as

$$E' = E + 1023.$$

Here bias value is 1023 and hence it is also called excess 1023 format.

→ Thus range  $E'$  for normal values in double Precision is  $0 < E' < 2047$ . This means that for 64 bit representation the actual exponent  $E$  is in range  $-1022 \leq E \leq 1023$ .

Represent  $1259.125_{10}$  in single Precision and double Precision formats.

Step 1:

16	1259	11	↑	LSB
16	78	14	↑	MSB
	4			

$(1259)_{10} = (4EB)_{16}$

Fraction

$0.125 \times 2$	=	<span style="border: 1px solid black; padding: 2px;">0</span> .25
$0.25 \times 2$	=	<span style="border: 1px solid black; padding: 2px;">0</span> .50
$0.50 \times 2$	=	<span style="border: 1px solid black; padding: 2px;">1</span> .00
		0.001

↓

he



# Point Multiplication.

13

## Multiplying Decimal Numbers in Scientific Notation

$$1.110_{\text{ten}} \times 10^{10} \times 9.200_{\text{ten}} \times 10^{-5}$$

Assume that we can store only four digits of the significant and two digits of the exponent.

Step 1: Unlike addition we calculate the exponent of the product by simply adding the exponents of the operands.

$$\text{New exponent} = 10 + (-5) = 5$$

Let's do this with the biased exponents as well to make sure we obtain the same result.

$$\begin{array}{l} - 10 + 127 = 137 \\ - 5 + 127 = 122 \end{array} \left. \vphantom{\begin{array}{l} - 10 + 127 = 137 \\ - 5 + 127 = 122 \end{array}} \right\} \begin{array}{l} \text{New exponent} = 137 + 122 \\ = 259 \end{array}$$

This result is too large for the 8 bit exponent field.

$$\begin{aligned} \text{So New Exponent} &= 137 + 122 - 127 \\ &= 259 - 127 = 132 \Rightarrow (5 + 127) \end{aligned}$$

So to get correct biased sum when we add biased numbers we must subtract the bias from the sum.

Step 2: Next comes the multiplication of the significant.

$$\begin{array}{r} 1.110_{\text{ten}} \\ \times 9.200_{\text{ten}} \\ \hline 0000 \\ 0000 \\ 2220 \\ 9990 \\ \hline 10212000_{\text{ten}} \end{array}$$

The product is  $10.21200_{\text{ten}}$ .



Assuming that we can keep only three digits to the right of the decimal point.

$$10.212_{\text{ten}} \times 10^5.$$

Step 3: This product is unnormalized, so we need to normalize it

$$10.212_{\text{ten}} \times 10^5 = 1.0212 \times 10^6.$$

Step 4: Round off the product.  $1.021_{\text{ten}} \times 10^6.$

Step 5: The sign of the product depends on the sign of the original operands. If they are both the same, the sign is positive; otherwise it's negative. Hence the product is

$$+ 1.021_{\text{ten}} \times 10^6.$$

Ex Multiply the numbers  $0.5_{\text{ten}}$  and  $-0.4375_{\text{ten}}$  using floating point multiplication algorithm.

$$\text{Binary equivalent of } 0.5_{\text{ten}} = 1.000 \times 2^{-1}$$

$$-0.4375_{\text{ten}} = -1.110 \times 2^{-2}.$$

Step 1: Adding the exponents without bias.

$$-1 + (-2) = -3.$$

or Using biased representation

$$(-1 + 127) + (-2 + 127) - 127$$

$$= (-1 - 2) + (127 + 127 - 127) = -3 + 127 = 124$$

Step 2: multiplying the significand.



$$\begin{array}{r}
 1.000. \\
 1.110 \\
 \hline
 0000 \\
 1000 \\
 1000 \\
 1000 \\
 \hline
 1110000 \text{ two}
 \end{array}$$

The Product is  $1.110000 \times 2^{-3}$   
 but we use 4 bits, so it is  
 $1.110 \times 2^{-3}$ .

Step 3: Normalize the Product as per one example.

$$1.110_{\text{two}} \times 2^{-3}$$

Step 4: Rounding the Product makes no change.

$$1.110_{\text{two}} \times 2^{-3}$$

Step 5: Since the sign of the original operands differ make the sign of the product negative.

$$-1.110 \times 2^{-3}$$

Converting to decimal to check our result.

$$-1.110 \times 2^{-3} = -0.00111 = \frac{1}{8} + \frac{1}{16} + \frac{1}{32} = 0.21875_{10}$$

### Advantages of floating Point

- \* More bits can be stored in the instruction format.
- \* Floating point has separate registers.

### Disadvantages of floating Point

- \* Floating point operation requires much more hardware than integer operations.
- \* Implementing floating point arithmetic is more complex.







1. If  $i^{\text{th}}$  bit  $b_i$  is 0 and  $(i-1)^{\text{th}}$  bit  $b_{i-1}$  is 1, then take  $b_i$  as +1 (ie, for (0,1) pair). Multiplicand is added to the partial product.
2. If  $i^{\text{th}}$  bit  $b_i$  is 1 and  $(i-1)^{\text{th}}$  bit  $b_{i-1}$  is 0, then take  $b_i$  as -1 (ie for (1,0) pair), Multiplicand is subtracted to the partial product.
3. If  $i^{\text{th}}$  bit is 0 and  $(i-1)^{\text{th}}$  bit  $b_{i-1}$  is 0 then take bias 0.
4. If  $i^{\text{th}}$  bit  $b_i$  is 1 and  $(i-1)^{\text{th}}$  bit  $b_{i-1}$  is 1, then take  $b_i$  as 0 (ie for (0,0) and (1,1) pair), then neither addition nor subtraction is performed.
5. When least significant bit  $b_0 = 1$  (assume that it had  $b_{i-1}$  as 0 (implied zero)), thus take  $b_0 = -1(0)$ .

Multiplier		Version of Multiplicand Selected by bit $i$
Bit $i$	Bit $i-1$	
0	0	$0 \times M$
0	1	$+1 \times M$
1	0	$-1 \times M$
1	1	$0 \times M$

Booth Multiplier Recording Table.

0 0 1 0 1 1 0 0 1 1 1 0 1 0 1 1 0 0  
 0 +1 -1 +1 0 -1 0 +1 0 0 -1 +1 -1 +1 0 -1 0

Booth recoding of a Multiplier.



Example Recode the multiplier 101100 for Booth's multiplication.

1 0 1 1 0 0 [0], implied zero.

-1 +1 0 -1 0 0 Recoded multiplier.

Example Multiply 001100 (+12) and 010011 (+19)

Multiplier : 001100. Multiplicand : 010011

Recoded Multiplier : 0+10-100.

Multiplication:

0 1 0 0 1 1  
0 +1 0 -1 0 0

0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	0	1	1	0	1		
0	0	0	0	0	0	0	0	0			
0	0	0	1	0	0	1	1				
0	0	0	0	0	0	0					

if multiplier & multiplicand 6 digits, add six zeros = 12

2's Comp of Multiplicand.

0 0 0 0 1 1 1 0 0 1 0 0 (228)

The same algorithm can be applied for negative multiplier and negative multiplicand.



Multiply  $01110 (+14)$  and  $11011 (-5)$ . (10)

Multiplicand :  $11011$       Multiplier :  $01110$

Revised Multiplier :  $0-1+10-1$

$01110$

$0-1+10-1$

$111110010$

$000000000$

$00001110$

$1110010$

$0000000$

$111011010 = (-70)$

### Hardware Implementation of Booths algorithm.

The hardware consists of 32 bit register M for the multiplicand, 64 bit product register PR and a 1-bit register C, 32 bit ALU and Control. Initially M contains multiplicand, PR contains multiplier in lower half of register and initialize the upper half of product register to 0 and C contains bit 0. The algorithm is the following steps.

Repeat 32 times.

1. If  $(P_0 C)$  pair is

\* 10 :  $PR = PR - M.$

\* 01 :  $PR = PR + M.$

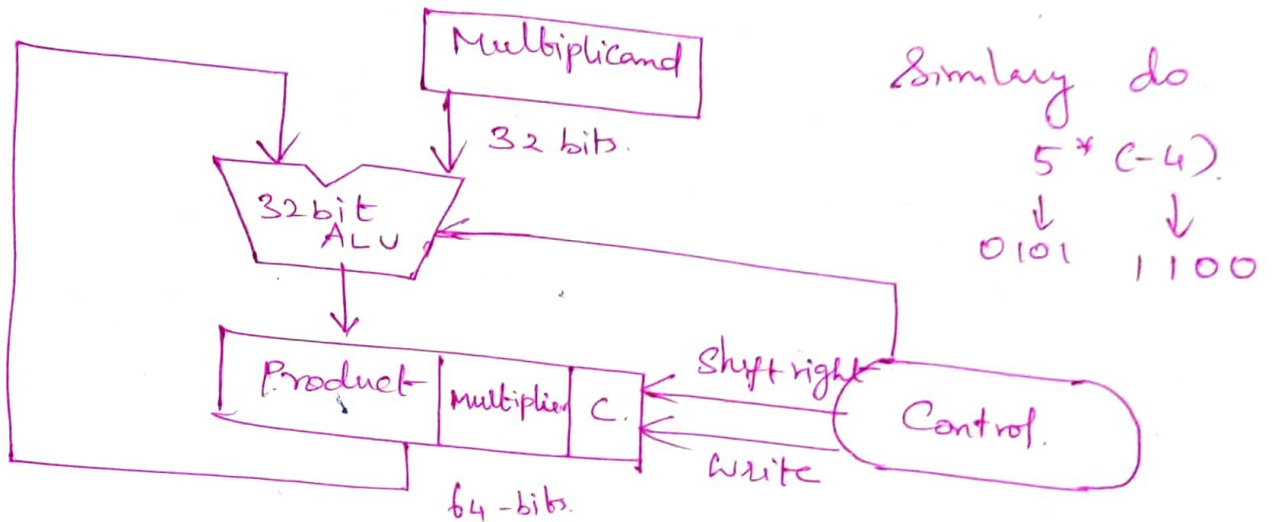
\* 00 and 11 : do nothing.

2. Arithmetic shift PR right 1 bit (Arithmetic shift preserves the sign of a two's complement number and it is equivalent to logical shift with sign extension)



The shift out bit gets into C.

In Booth's algorithm if the multiplicand and multiplier are n-bit two's Complement numbers. the result is considered as 2n-bit two's Complement value. The overflow bit @ is ignored.



Hardware for Booth's Multiplication

Case 1: Both Positive ( $5 \times 4$ )

Steps	PR	multiplier	C	Operation
Step 1	0000 0000	0100 0010	0	Initial
Step 2	0000	0001	0	Shift right
Step 3	1011 1101	0001 1000	0 1	Shift right PR ← PR - M
Step 4	0010 0001	1000 0100	1 0	Shift right PR ← PR + M
Result:	0001	0100	+do.	



THE PROCESSORIntroduction

\* MIPS (Microprocessor without Interlocked Pipeline Stages). Is a reduced Instruction Set Computer (RISC) developed by MIPS Technologies.

\* Multiple revisions of the MIPS instruction set exist, including MIPS I, MIPS II, MIPS III, MIPS IV, MIPS V, MIPS 32, and MIPS 64. The current revisions are MIPS 32 (for 32 bit implementation) and MIPS 64 (for 64-bit implementations).

Characteristics of MIPS

\* ALU Operations have 3 register operands.

\* Simple instruction set.

\* Uniform encoding:

- All instructions are 32 bit long.
- opcode is always in the high order 6 bits.
- 3 types of instruction format.

Registers for MIPS

→ MIPS 64 has 32 64-bit general purpose registers (GPRs) named R0, R1... R31.

→ GPR are also sometimes known as Integer registers.

→ Additionally there is a set of 32 floating point registers (FPRs) named F0, F1... F31.

→ Both single and double precision floating point operations (32-bit and 64 bit) are provided.

→ MIPS also include instructions that operate on two single precision operands in a single 64-bit floating point register.



# A basic MIPS Implementations

The basic MIPS implementation that includes a subset of the core MIPS instruction set.

- \* The Memory Reference Instruction load word (lw) and store word (sw).
- \* The arithmetic-logical Instruction add, sub, AND, OR and slt.
- \* The instructions branch equal (beq) and jump (j).

## Logic Design Convention

→ To design any computer, we must decide how the hardware logic implementing the computer will operate and how the computer is clocked.

→ The data path element in MIPS implementation consists of two different types of logic elements.

- a) Combinational Element
- b) State Element

## Combinational Elements

- \* Elements that operate on data values and elements that contain state.
- \* Output depends only on current inputs.
- \* Does not have internal storage.

Example AND gate, multiplexer and ALU.

## State Elements

- \* It contains internal storage.
- \* State element has at least two inputs and one output.
- \* State element also called sequential elements because o/p depends on current as well as previous input.



\* The inputs are data values written into the element and the clock, which determines when the data value is written. (19)

Example: D-flip flop  $\rightarrow$  which has exactly two inputs. (a value and clock) and one o/p.  
2. Memories and Registers.

### Clocking Methodology

\* A clocking Methodology defines when signals can be read and when they can be written.

\* It is important to specify the timing of reads and writes, because if a signal is written at the same time it is read, the value of the read would correspond to the old value and newly written value or even some mix of the two.

\* clocking Methodology is designed to make hardware predictable.

### Edge Triggered clocking

\* Edge triggered clocking methodology means any value stored in a sequential logic element are updated only on a clock edge.

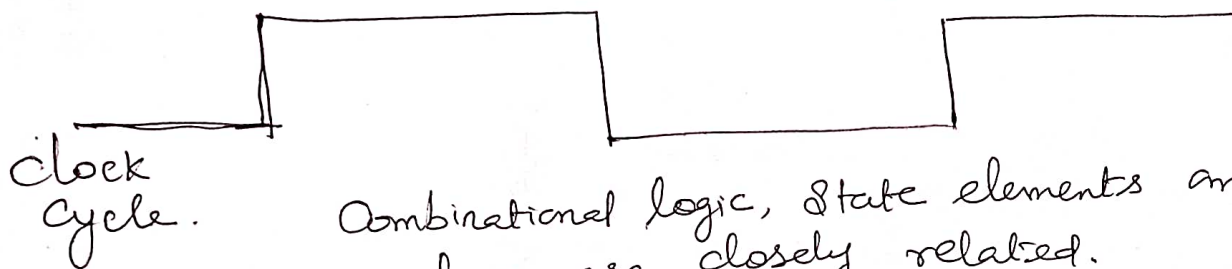
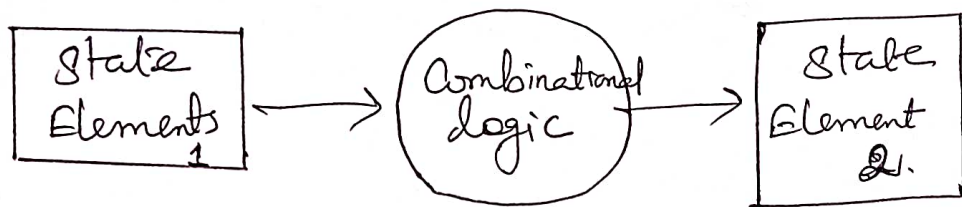
\* It is a quick transition from low to high or high to low.

\* Because only state elements can store a data value, any collection of combinational logic must have its inputs come from a set of state elements and its o/p written into a set of state elements.



## Typical Execution.

- \* Read contents of some state elements.
- \* Send values through some combinational logic
- \* Write results to one or more state elements.



Combinational logic, state elements and the clock are closely related.

\* It shows two state elements surrounding a block of combinational logic, which operates in a single clock cycle, all signals must propagate from state element 1, through combinational logic and to state element 2 in the time of one clock cycle.

\* The time required for the signals to reach state element 2 is called as length of the clock cycle.

\* If the state element is not updated on every clock then an explicit write control signal is required.

\* Both the clock cycle and write control signal are input and state element is changed only when the write control signal is asserted and a clock edge occurs.

\* Control signal is a signal used for multiplexer selection for directing the operation of a functional unit.

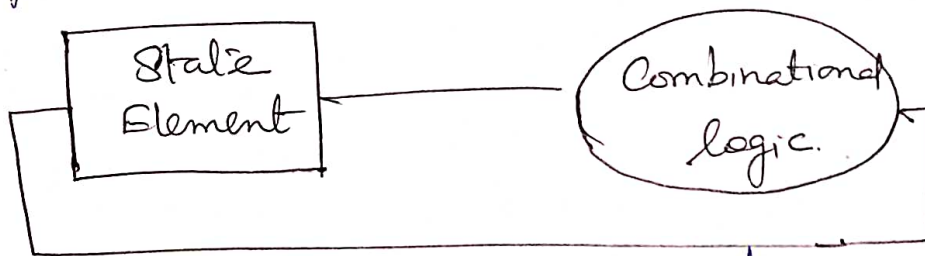
\* Two logical terms are used to represent the signal level.

20  
→ The term asserted to indicate a signal that is logically high and assert to specify that a signal should be driven logically high.

→ The term deasserted to represent logically low.

→ An edge triggered methodology allows us to read the content of register, send the value through some combinational logic and write that register in the same clock cycle.

example .



→ After reading process was done send the value through some combinational logic and write the register in the same clock cycle.

This allows the processor to read the register content send the value through some combinational logic and write that register in same clock under the assumption that the state elements are controlled by implicit clock cycles.

### Building Data Path

→ A data path is a collection of functional units organised in a manner to execute each class of instruction.



→ A data path is a representation of the flow of information (data and instruction) through the CPU implemented using Combinational and Sequential circuitry

→ Single cycle datapath that have separate instruction and memories because.

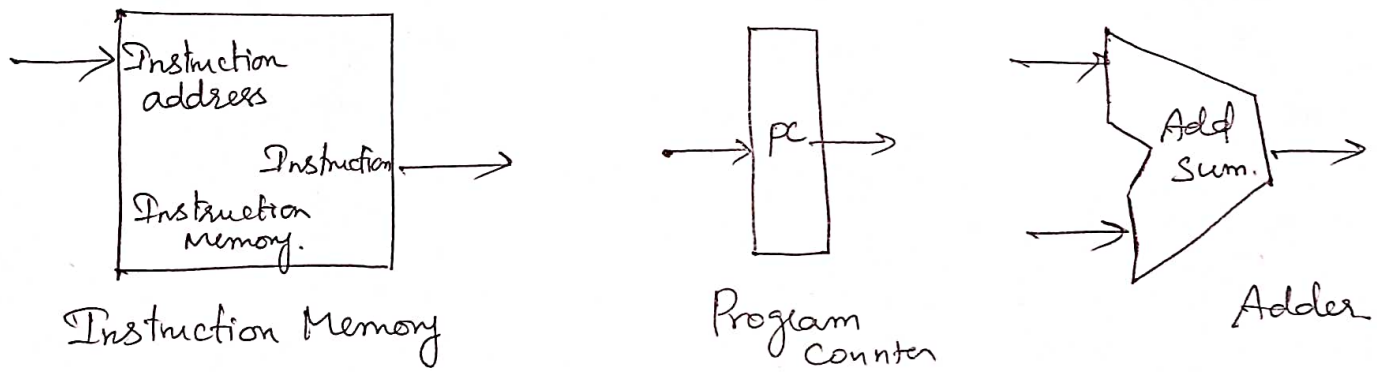
① The format of data and instructions is different in MIPS and hence different memories are needed. Having separate memories is less expensive.

②.

To design a data path first we should examine the major components required to execute MIPS instruction.

### Data path elements

→ Data path elements include. 1. Instruction memory. 2. Program Counter 3. Adder, 4. Register file. 5. ALU, 6. Data memory. 7. Sign extend. 8. Multiplexers.



Major Components to execute MIPS Instruction.

### Instruction Memory

→ It is first data path element, memory unit to store the instructions of a program and supply instructions given an address.

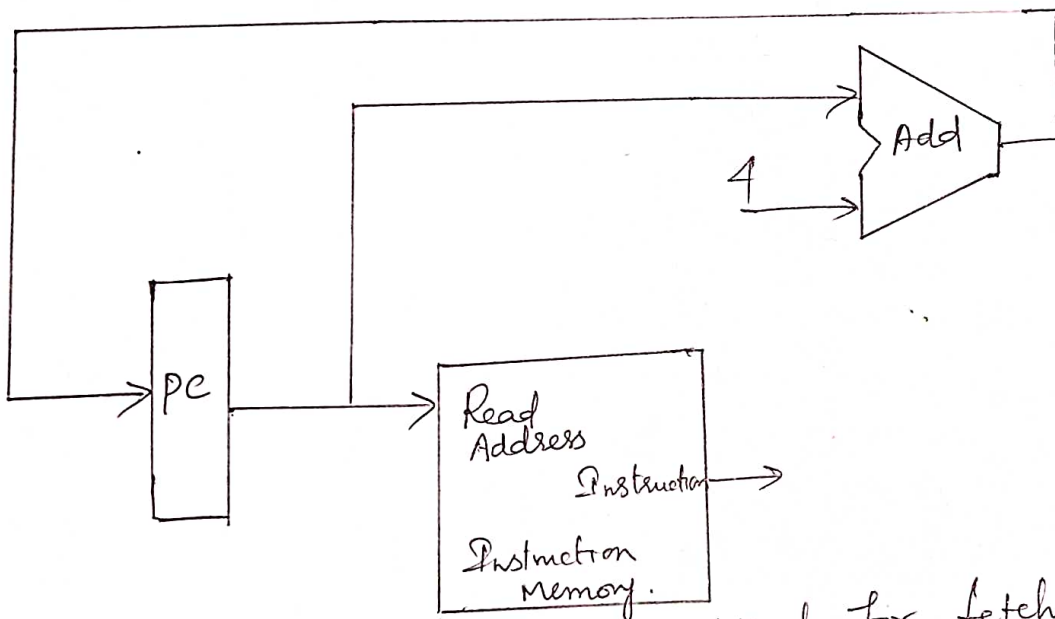
## Program Counter

Program Counter is used to hold the address of current instruction.

Adder: → It adds two 32 bits input and place the sum on its output. It is also increment the PC to the address of the next instruction.

## Fetch operation.

→ The fundamental operation in Instruction Fetch is to send the address in the PC to the instruction memory and obtain the specified instruction, and the increment the PC.



A portion of the datapath used for fetching instruction and incrementing the Program Counter.

## Register file.

\* The Processor's 32 general-purpose registers are stored in a structure called a register file.

\* A register file is a collection of registers in which any register can be read or written by specifying the number of the register in the file.



\* The register file contains the register state of the machine.

\* In addition we will need an ALU to operate on the values read from the registers.

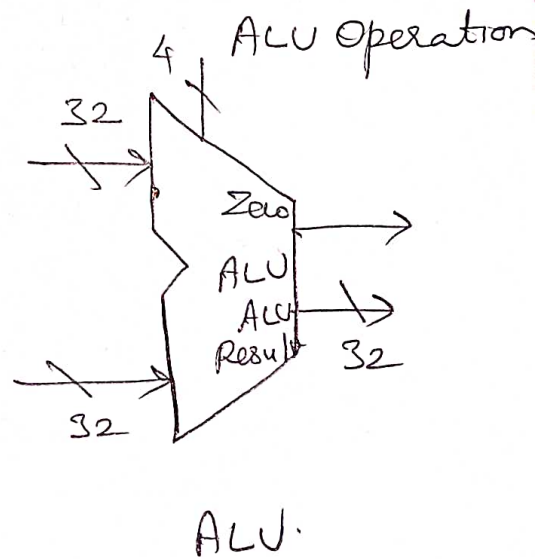
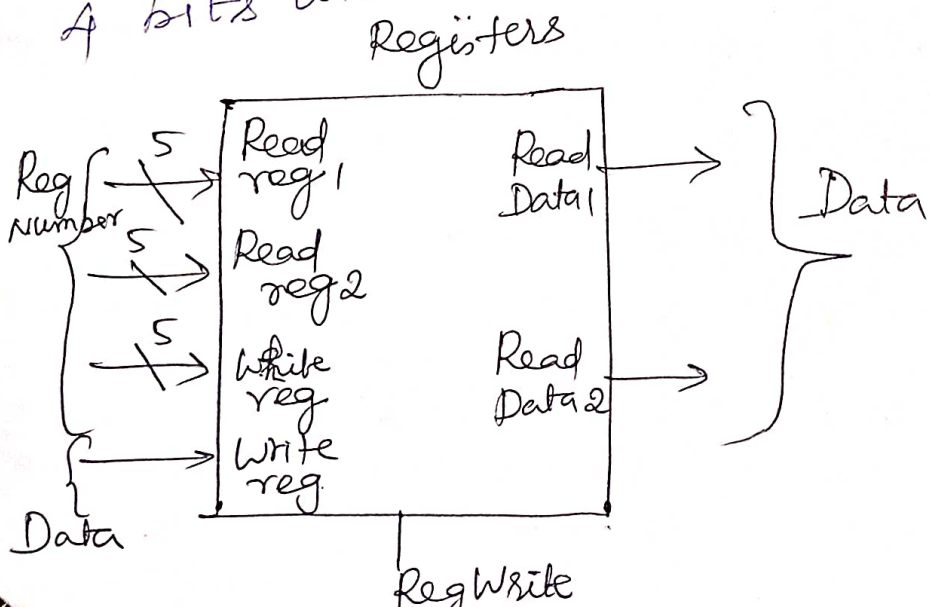
\* Because the R-format instructions have three register operands, we will need to read two data words from the register file and write one data word into the register file for each instruction.

\* For each data word to be read from the register we need an input to the register file that specifies the register number to be read and an output from the register file that will carry the value that has been read from the registers.

\* To write a data word, we will need two inputs; one to specify the register number to be written and one to supply the data to be written into the registers.

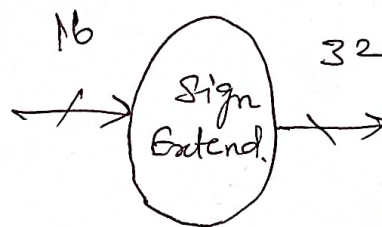
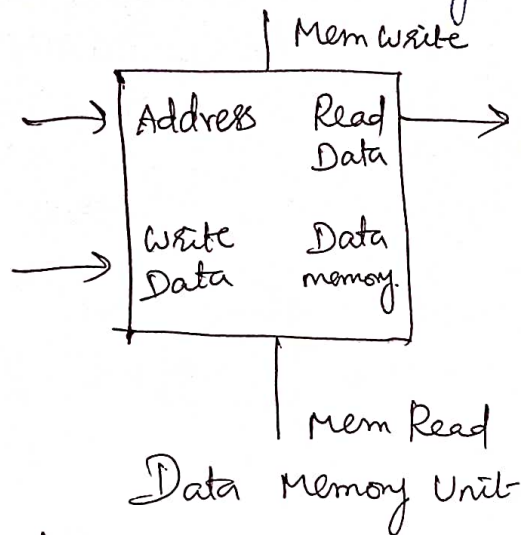
ALU → It shows the ALU, which takes two 32-bit inputs and produces a 32-bit result, as well as a 1-bit signal if the result is 0.

→ The operation to be performed by the ALU is controlled with the ALU operation signal, which will be 4 bits wide.



# Data Memory.

Data Memory must be written on Store Instruction hence. it has both read and writes control signals. an address input. as well as an input for the data to be written into memory.



Sign extension Unit

## Multiplexers

four multiplexers are used in simple Data path

### Multiplexer 1: Select second ALU operand.

It is used to select the second ALU operand. It is controlled by ALU<sup>src</sup> control signal. It selects any one of the following two values as the second input to the ALU.

1. Register value for R-type instruction.
2. 32 bit sign extended value for memory instruction.

### Multiplexer 2: Select input for write data.

→ It is used to select the data for write data input of register file. It is controlled by Mem to Reg. control signal. The two possible inputs are the following and this multiplexer select any one of the following two values.



① Data from memory for a memory instruction.

② ALU result for a R-type instructions.

Multiplexer 8: Selects the next PC value.

It is used to select the next value for the PC. It is controlled by PC Src control signal. The two possible cases are as follows as this this multiplexer is used to select any one of the following two values.

1.  $PC + 4$ .

2. Branch Address.

Multiplexer 4: Selects the destination register.

It is used to select the destination register. It is controlled by RegDst control signal. It select one of the two values.

① Destination comes from rt field of the instruction.

② Destination comes from rd field of the instruction.

Data Path for Branch Instruction.

Branch Target address calculation.

→ The beq instruction has three operands, two registers that are compared for equality, and a 16 bit offset used to compute the branch target address. relative to the branch instruction address.

→ The syntax for the branch instruction is beq \$t1, \$t2, offset.

\* To implement this instruction we must compute the branch target address by adding the sign extended offset field of the instruction to the PC.

→ Two important things in the definition of branch instructions.

①. The instruction set architecture specifies that the base for the branch address calculation is the address of the instruction following the branch (i.e., PC + 4, the address of next instruction).

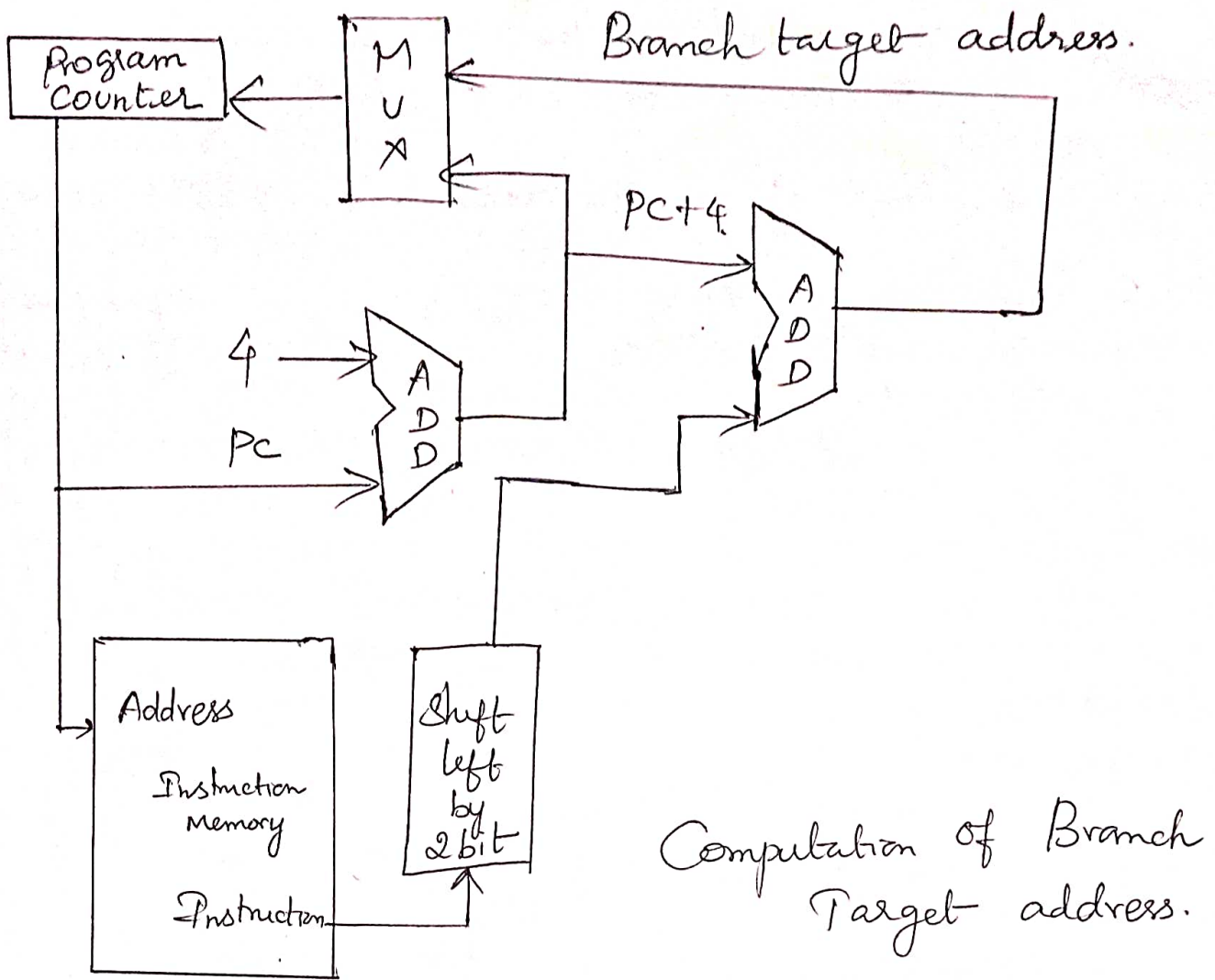
②. The MIPS architecture also states that the offset field is shifted left 2 bits so that it is a word offset; this shift increases the effective range of the offset field by a factor of 4. Then to shift the offset field by two.

→ Therefore branch target address is given by  
Branch Target address = PC + 4 + offset (shifted left 2 bits)

→ In addition to computing the branch target address, we must also see whether the two operands are equal or not. If two operands are not equal the next instruction is the instruction that follows sequentially (PC = PC + 4). In this case, we say that the branch is not taken.

→ On other hand if two operands are equal (i.e. condition is true), the branch target address becomes the new PC, and we say that the branch is taken.





## Jump Instruction

→ The Jump Instruction operates by replacing the lower 28 bits of the PC with the lower 26 bits of the instruction shifted left by 2 bits. This shift is accomplished simply by concatenating 00 to the Jump Offset, accomplishes this shift.

## Simple Implementation Scheme.

→ Control implementation scheme can be build using datapath and some implementation method.

→ This simple implementation covers load word (lw), store word (sw), branch equal (beq), and arithmetic logical instructions add, sub, and, or, and set on less than.

# The ALU Control

The MIPS ALU defines the six following combinations of four control inputs.

ALU Control lines	Function
0000	AND
0001	OR
0010	Add
0110	Subtract
0111	Set on less than
1100	NOR.

→ Depending on the instruction class, the ALU will need to perform one of these first five functions.

→ For load word and store word instructions, we use the ALU to compute the memory address by addition.

→ In case of the R type instruction, the ALU needs to perform one of the five actions (AND, OR, Subtract, add, or, Set on less than).

\* In case of Branch equal the ALU must perform a subtraction

We can generate the 4-bit ALU control input using a small control unit that has an input the function field of the instruction and 2 bit control field. which we call ALUOP.

ALUOP	Action
00	load and store
01	Subtract for beg.
10	The operation encoded in the function field
11	



→ Table shows how to Set the ALU Control Input based on the 2-bit ALU op. Control and 6-bit function code.

Instruction opcode	ALUOP	Instruction operation	funct. field	Desired ALU Action	ALU Control $\mu p$ .
		Loadword	XXXXXX	add	0010
LW	00				
		Storeword	XXXXXX	add	0010
SW	00				
		branch equal	XXXXXX	Subtract	0110
Branch equal	01				
		add	100000	add	0010
R type	10				
		Subtract	100010	Subtract	0110
R type	10				
		AND	100100	AND	0000
R type	10				
		OR	100101	OR	0001
R type	10				
		Set on less than	101010	Set on less than	0111
R type.	10				

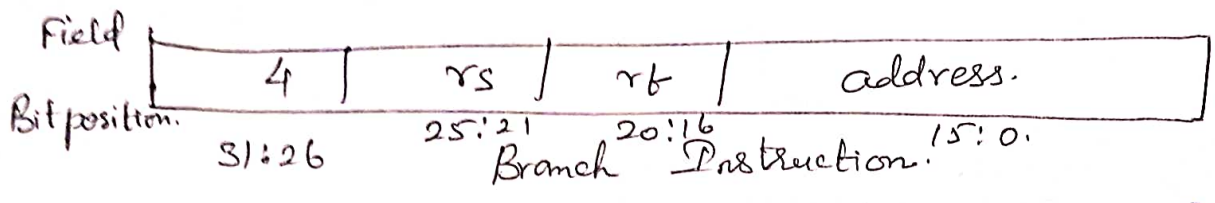
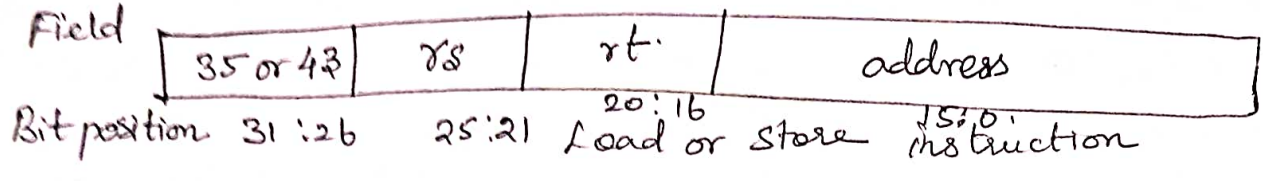
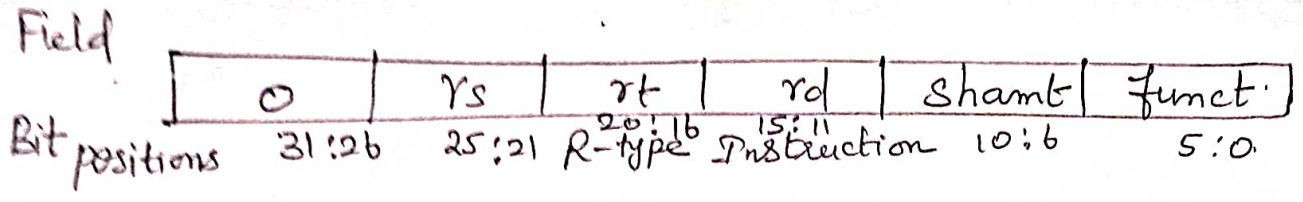
→ Here Multiple levels of decoding technique is Used.

- Advantages of Using multiple levels of decoding
- \* It reduces the size of the main Control Unit;
  - \* Use of several smaller Control Units may also potentially increase the speed of the control Unit.

Designing the Main Control Unit

To design a main Control Unit, first identify the fields of an instruction and control signals needed for the datapath. It is useful to review the format of the three instruction classes.

- ① R type Instruction
- ② Branch Instruction
- ③ Load/Store instruction.



- The op field also called the opcode. is always contained in bits 31:26.
- The two registers to be read are always specified by the rs and rt field. at position 25:21 and 20:16. This is true for R type instruction branch equal and for store.
- The base register for load and store instructions is always in bit position 25:21.

## An Overview of Pipelining.

### Pipelining

- Pipelining is an implementation technique in which multiple instructions are executed simultaneously by overlapping them in execution to save time and resource and improves the throughput.
- In a computer pipeline each step in the pipeline completes a part of an instruction. Different steps are completing different instructions in parallel. Each of these steps is called a pipeline stage (or) pipeline segment.
- Pipelining yields a reduction in the average execution time per instruction.



MIPS instructions Classically take five steps:

① Fetch instruction from memory.

② Read registers while decoding the instruction  
The format of MIPS instructions allows reading and decoding to occur simultaneously.

③ Execute the operation or calculate an address.

④ Access the operand in data memory.

⑤ Write the result into a register.

Single cycle Versus Pipelined Performance.

→ Consider a simple Program segment consists of eight instructions: load word (lw), store word (sw), add (add), subtract (sub), and (and), or (or), set less than (slt) and branch on equal (beq).

→ Compare the average time between instruction of a single cycle implementation, in which all instructions take 1 clock cycle, to a pipelined implementation.

Instruction class.	Instruction fetch	Register Read	ALU Operation	Data Access	Register Write	Total Time
Load word (lw)	200ps	100ps	200ps	200ps	100ps	800ps
Store word (sw)	200ps	100ps	200ps	200ps	100ps	700ps
R-format Cadd, sub, and or, slt)	200ps	100ps	200ps			600ps
Branch (beq)	200ps	100ps	200ps			500ps

Total time for each instruction calculated from the time for each component.

# Pipeline Speedup:

→ If the stages are perfectly balanced, then the time between instructions on the pipelined processor assuming ideal conditions - is equal to:

$$\text{Time between instructions pipelined} = \frac{\text{Time between instructions nonpipelined}}{\text{Number of Pipe stages}}$$

→ The above formula suggests that a five stage pipeline should offer nearly a five fold improvement over the 800 ps nonpipelined time or a 160 ps clock cycle.

→ Pipelining improves performance by increasing instruction throughput, as opposed to decreasing the execution time of an individual instruction.

## Pipeline Hazards.

→ There are situations in pipelining when the next instruction cannot execute in the following clock cycle. This means that it prevents the next instruction in the instruction stream from executing during its designed clock cycle. This is called Hazard.

→ There are three different types of hazards.  
① Structural Hazard    ② Data Hazard    ③ Control Hazard

### ① Structural Hazard.

→ First hazard is called a Structural Hazard. It means that the hardware cannot support the combination of instructions that we want to execute in the same clock cycle.



→ The Performance of pipelined processor depends on whether the functional units are pipelined and whether they are multiple execution units to allow all possible combinations of instructions in the pipeline. If for some combination, pipeline has to be stalled to avoid the resource conflicts then there is a structural hazard.

## Data Hazard

→ When either the source or the destination operands of an instruction are not available at the time expected in the pipeline and as a result pipeline is stalled, we say such a situation is a data hazard.

→ Consider a program with two instructions,  $I_1$  followed by  $I_2$ . When this program is executed in a pipeline, the execution of these two instructions can be performed concurrently. In such case the result of  $I_1$  may not be available for the execution of  $I_2$ .

→ If the result of  $I_2$  is dependent on the result of  $I_1$  we may get incorrect result if both are executed concurrently. For example  $A = 10$ . In the following two operations.

$$I_1: A \leftarrow A + 5$$

$$I_2: B \leftarrow A \times 2$$

→ In this case data used in the  $I_2$  depend on the result of  $I_1$ .

→ The hazard due to such situation is called data hazard or data dependent hazard.

## Solution for data Hazard

Adding extra hardware to retrieve the missing item early from the internal resource is called forwarding or bypassing.

- ① Operand forwarding hardware
- ② Reordering code (software).

Operand forwarding → It is also called bypassing.

→ It is method of resolving a data hazard by retrieving the missing data from internal buffer rather than waiting for the data to arrive from register or memory.

## Control Hazard.

→ Control hazard also called branch hazard. arises from the need to make a decision based on the result of one instruction while others are executing.

→ When the proper instruction cannot execute in the proper pipeline clock cycle because the instruction that was fetched is not the one that is needed; that is the flow of instruction addresses is not what the pipeline expected.

## Pipelined Data Path and Control.

\* The goal of pipelining is to allow multiple instructions execute at the same time. we may need to perform several operations in a cycle.

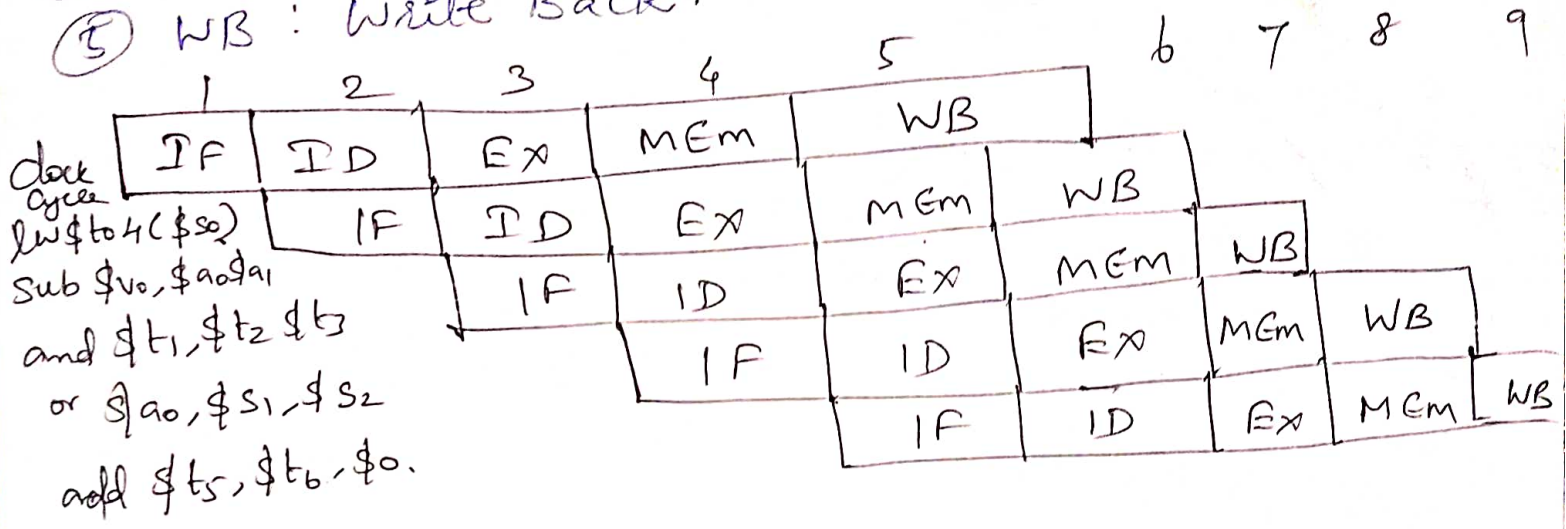
→ Increment the PC and add registers at the same time

→ Fetch one instruction while another one reads or writes data.

\* Pipelined data path is divided into five pieces, with each piece named corresponding to a stage of instruction execution.



- ① IF: Instruction fetch
- ② ID: Instruction decode and register file
- ③ Ex: Execution or address calculation
- ④ MEM: Data memory access.
- ⑤ WB: Write Back.

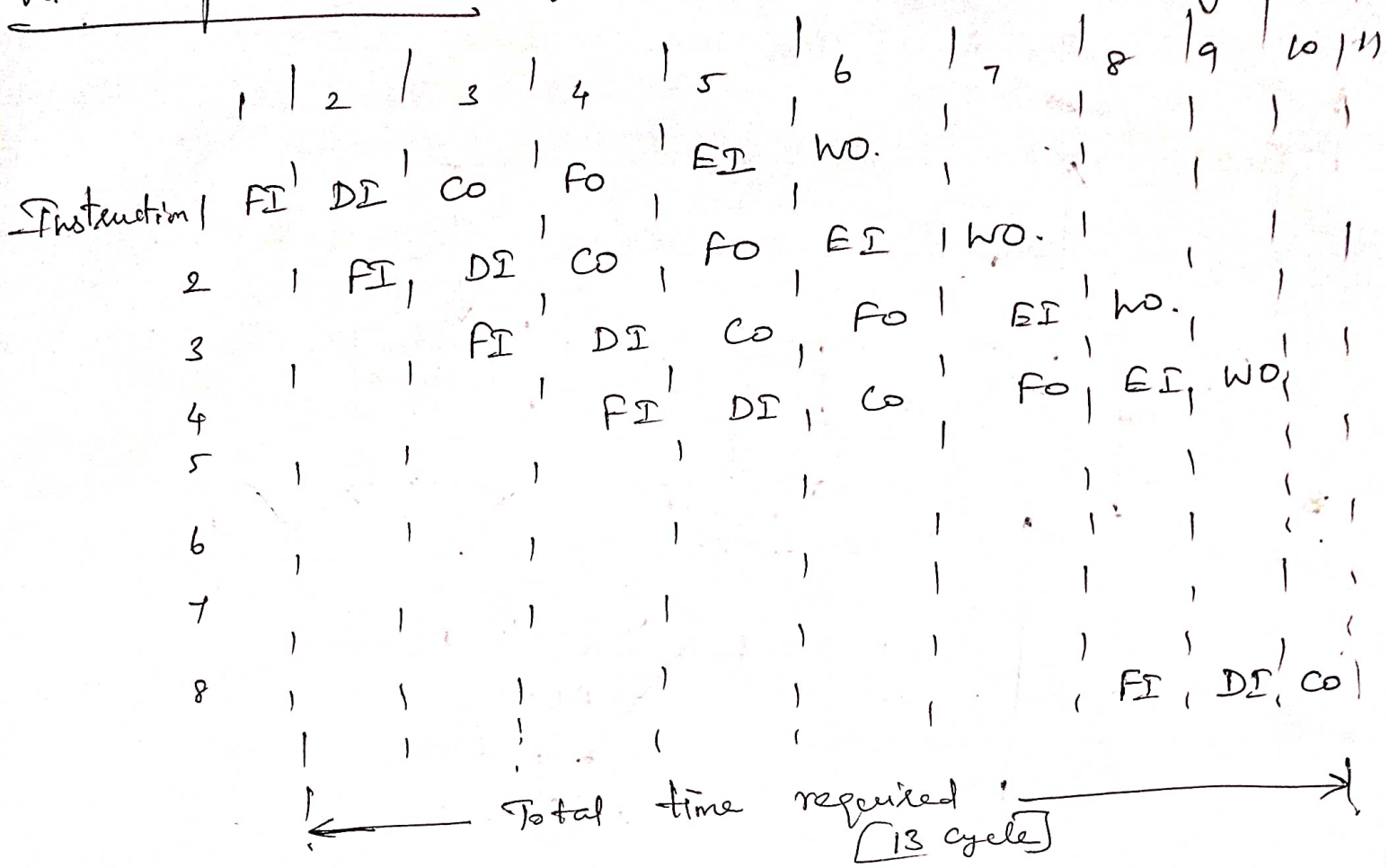


Example: Explain the function of a six segment pipeline and draw a space diagram for a six segment pipeline showing the time it takes to process eight tasks.

Solution Six stages in the pipelines.

- ① Fetch Instruction (FI): Read the next expected instruction into a buffer.
- ② Decode Instruction (DI): Determine the opcode.
- ③ Calculate Operands (CO): Calculate the effective address of each source operand.
- ④ Fetch operands (FO): Fetch each operand from memory.
- ⑤ Execute Instruction (EI): Perform the indicated operation and store the result, if any in the specified destination operand location.

Write Operand (W0) store the result in memory. (28)



Data Hazards Forwarding Versus Stalling:

→ Data Hazards arise when an instruction depends on the results of a previous instruction in a way that is exposed by the overlapping of instructions in the pipeline.

Operand forwarding

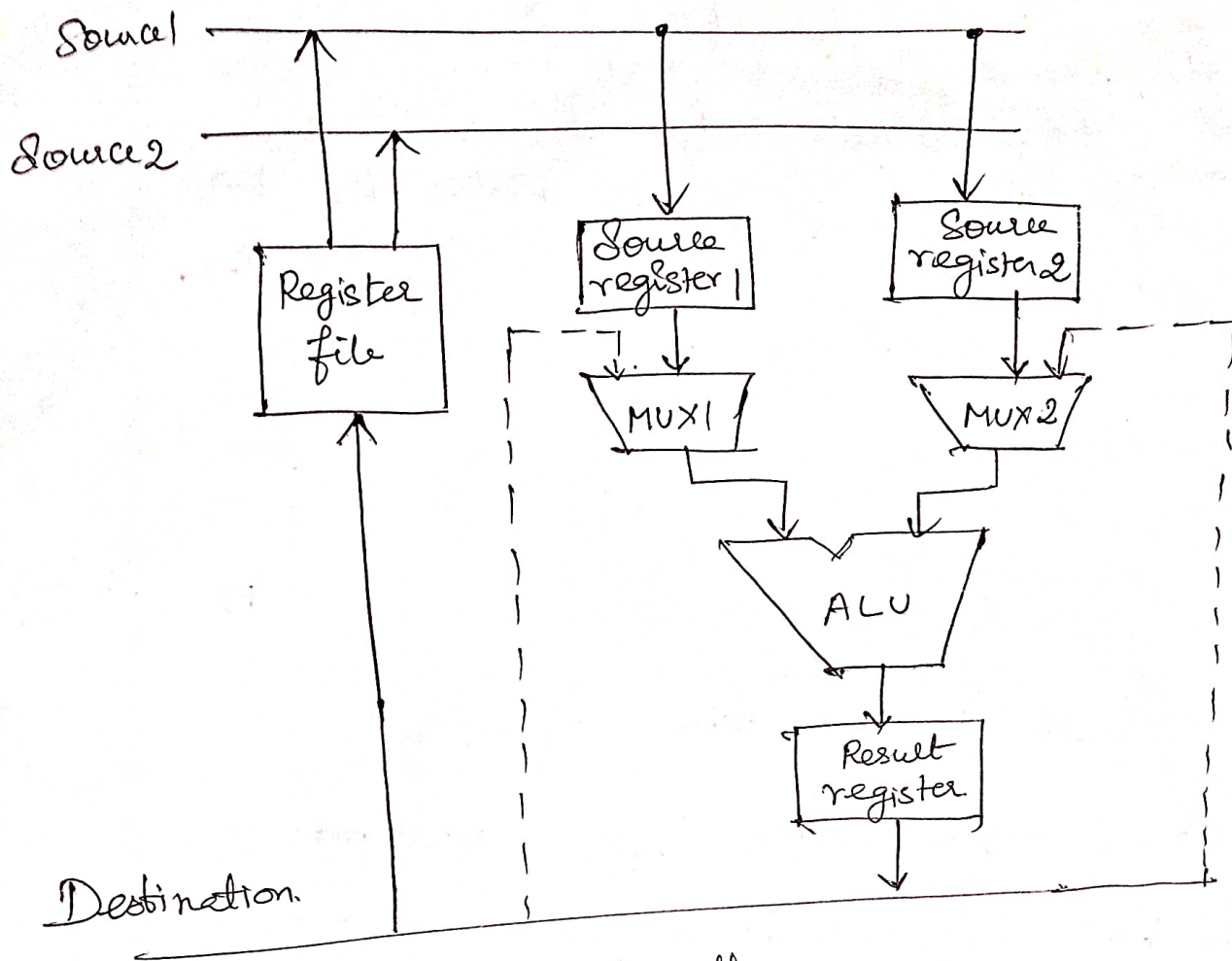
→ A simple hardware technique which can handle data hazard is called operand forwarding or register by passing. It is used to avoid the data hazards in pipeline process.

→ In this technique, ALU results are fed back as ALU inputs. When the forwarding logic detects the previous ALU operation has the operand for current instruction. It forwards ALU output instead of register file.

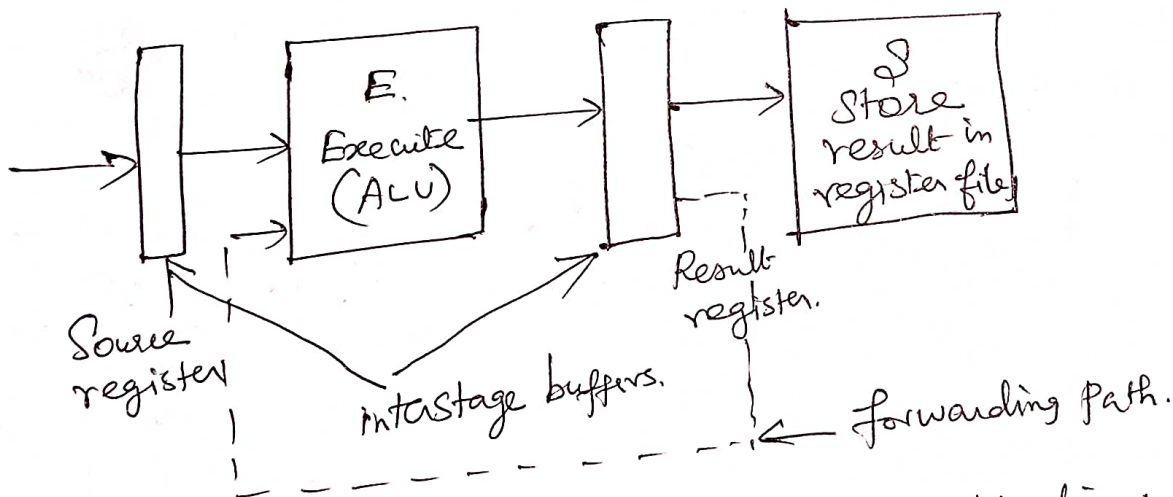
→ forwarding is



→ It shows a portion of the Processor datapath involving the ALU and register file.



(a) Data path



forwarding path in the processor pipeline.

## Data Hazards and Stalls.

- Hazards in pipeline can make the pipeline to stall.
- Eliminating a hazard often requires that some instructions in the pipeline to be allowed to proceed while others are delayed.
- when an instruction is stalled, instructions issued later than the stalled instruction are stopped. while the ones issued earlier must continue.
- No new instructions are fetched during the stall.
- In addition to a forwarding unit we need a hazard detection unit.
- It operates during the ID stage so that it can insert the stall between the load and its use.

## Control Hazards.

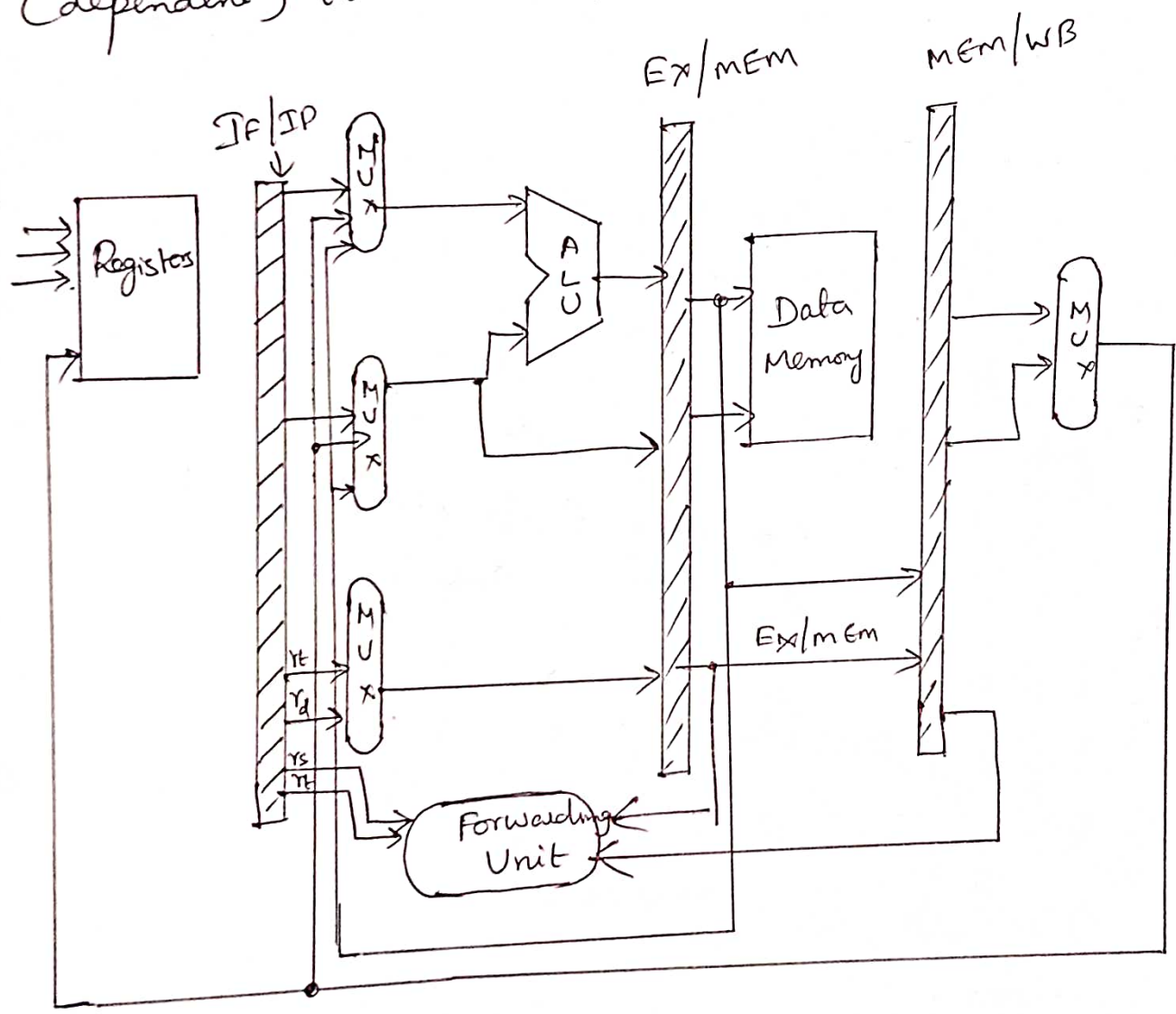
The next element to go into the pipeline may depend on currently executing instruction or we may have to wait until a stage is completed to determine the next stage.

(eg) Branch Instruction.

- An instruction must be fetched at every clock cycle to sustain the pipeline, yet in our design the decision about whether to branch doesn't occur until the MEM pipeline stage. This delay in determining the proper instruction to fetch is called control hazard or branch hazard.



- The data forwarding mechanism is indicated by dashed lines.
- The two multiplexers select the data for ALU either from destination bus or from Source 1 and Source 2 registers.
- When the forwarding logic detects data dependency it forwards ALU Output available in the result register using data forwarding path to the ALU for the next operation. Hence the execution of next (dependent) instruction proceeds without interruption.

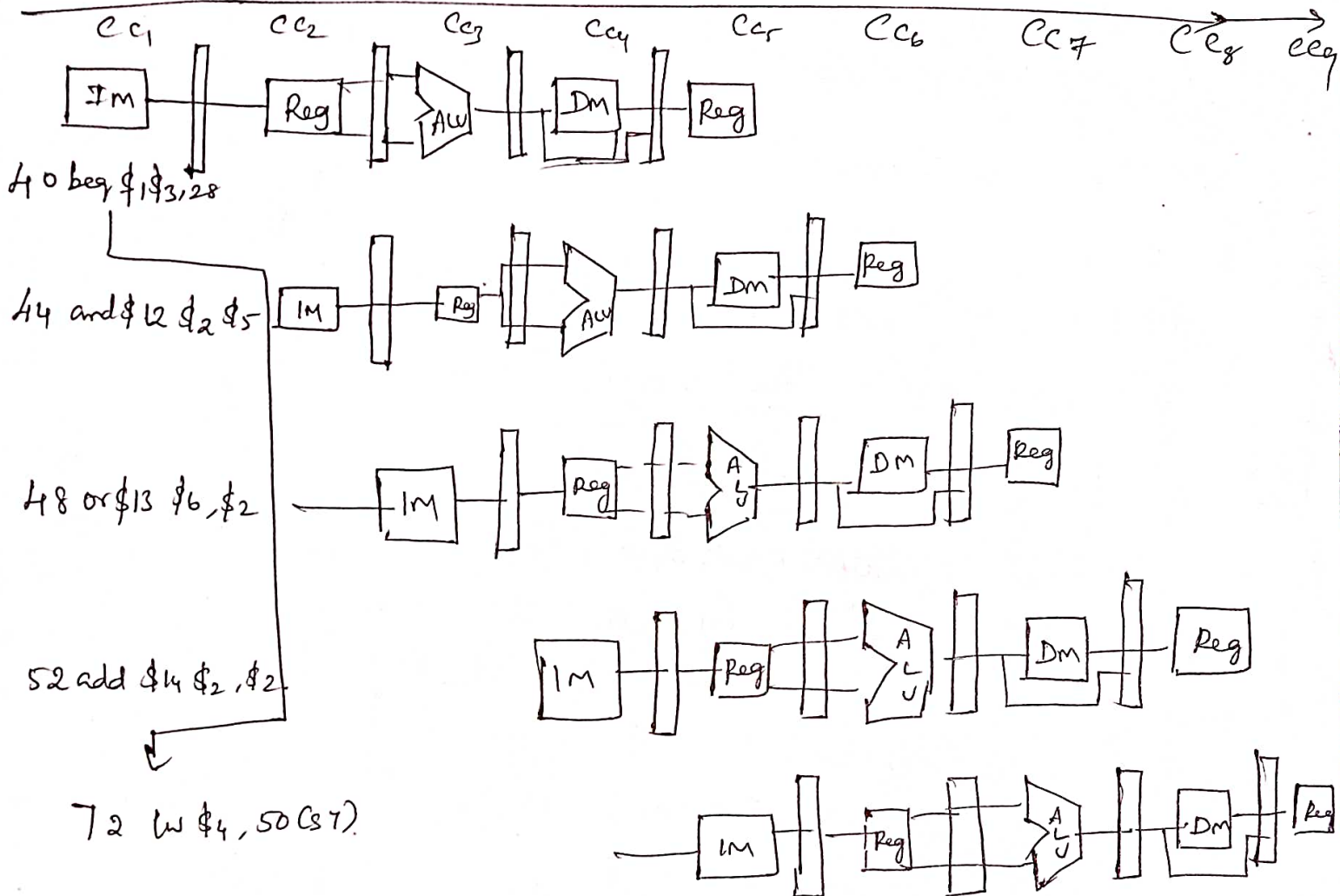


With forwarding.  
 ALU and Pipeline Registers. before adding forwarding

This section on Control hazard is shorter than the previous sections on data hazards because:

- a) Control hazards are relatively simple to understand.
- b) They occur less frequently than data hazards.

Time (in clock cycle)



The impact of the pipeline on the branch instruction.

### Reducing the Delay of Branches.

- One method to improve branch performance is to reduce the cost of the taken branch.
- The MIPS architecture was designed to support fast single cycle branches that could be pipelined with a small branch penalty.



→ The designers observed that many branches rely only on simple tests (equality or sign for example) and that such tests do not require a full ALU operation but can be done with at most a few gates.

→ When a more complex branch decision is required, a separate instruction that uses an ALU to perform a comparison is required.

→ Moving the branch decision up requires two actions to occur earlier.

(A) Computing the branch target address.

(B) Evaluating the branch decision

### Reducing Pipeline Branch Penalties:

→ There are many methods for dealing with the pipeline stalls caused by branch delay. There are two important schemes for reducing pipeline branch penalties.

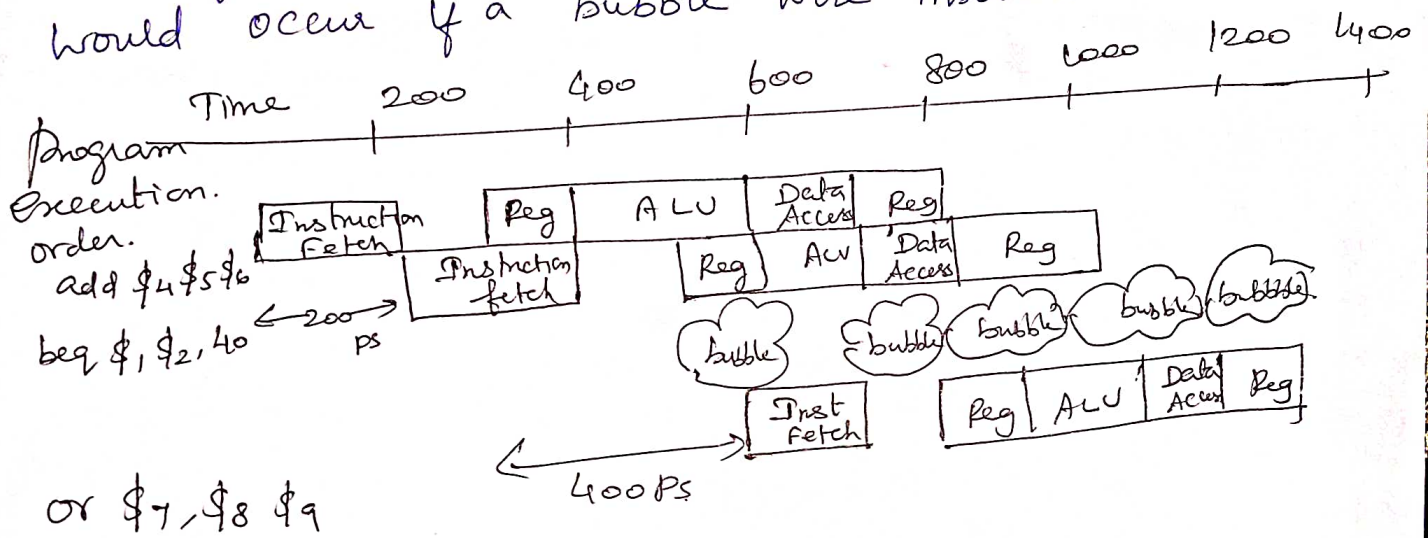
→ The first scheme is the static branch prediction. In this case branches are fixed for each branch during the entire execution. The software can try to minimize the branch penalty using knowledge of the hardware scheme and of branch behavior.

### Solution to Control Hazards:

→ Control hazard also called branch hazard, arises from the need to make a decision based on the results of one instruction while others are executing.

When the proper instruction cannot execute in the proper pipeline clock cycle because the instruction that was fetched is not the one that is needed.

Example. This example assumes the conditional branch is taken and the instruction at the destination of the branch is the OR instruction. There is a one stage pipeline stall, or bubble, after the branch. The effect on performance, however, is the same as would occur if a bubble were inserted.



Solution for control hazard

1. Branch Prediction
2. Delayed slot.



# UNIT-IV

## Memory and I/O Organization.

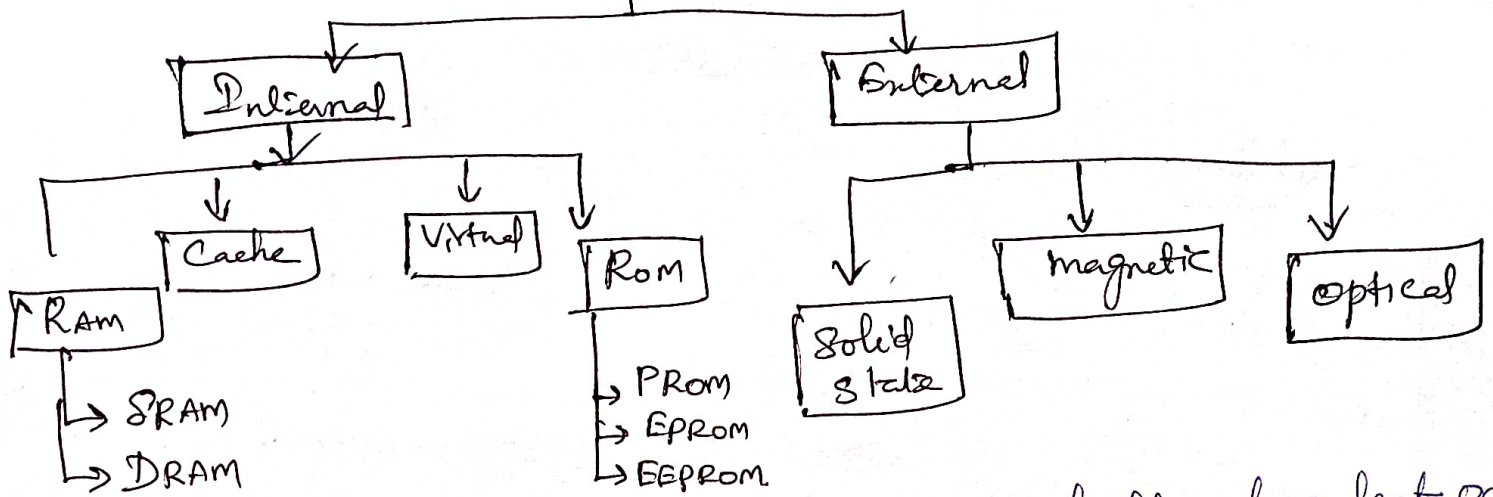
### Introduction

- Memories are made up of registers. Each register in the memory is one storage location also called memory location.
- Each memory location is identified by an address. The number of storage locations can vary from a few in some memories to hundred of thousand in others.
- Each register consists of storage elements (flipflops or capacitors in semiconductor memories and magnetic domain in magnetic storage). Each of which stores one bit of data. A storage element is called a cell.

### Characteristics of memory

- Location → CPU, Internal (main) & External (secondary).
- Capacity → word size, Number of words.
- Unit of Transfer → word, Block.
- Access Method → sequential access, Direct access, Random access, Associative access.
- Performance → Access time, Cycle time, Transfer rate.
- Physical Type → Semiconductor, magnetic, surface.
- Physical characteristics → Volatile / Nonvolatile.
- Computer employs different types of memory. basically it have two types.
  1. Internal
  2. External Memory.

# Computer Memory



→ Execution speed of programs is highly dependent on the speed with which instructions and data can be transferred between processor and the memory.

→ Memory would be fast, large and inexpensive. But is impossible to meet all the three requirements simultaneously.

→ Increased speed and size increased the cost, this problem is solved by developing another level structure.

→ This can be achieved by speed → cache.  
Virtual memory → increase apparent size of memory.  
Secondary memory → produces larger storage capacity.

## Memory Hierarchy.

→ An economical solution to that desire is a memory hierarchy which takes advantages of locality and cost performance of memory technologies.

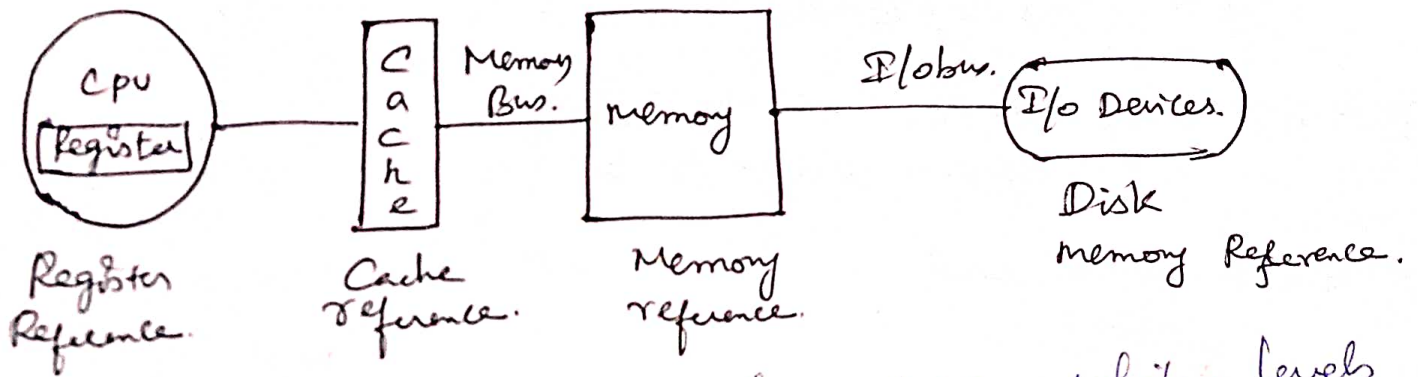
→ Principle of locality states that program access a relatively small portion of their address space at any instant of time.

→ A memory hierarchy consists of multiple levels of memory with different speed and sizes.



Faster memories are more expensive per bit than the slower memories.

→ There are three primary technologies used in building memory hierarchy..

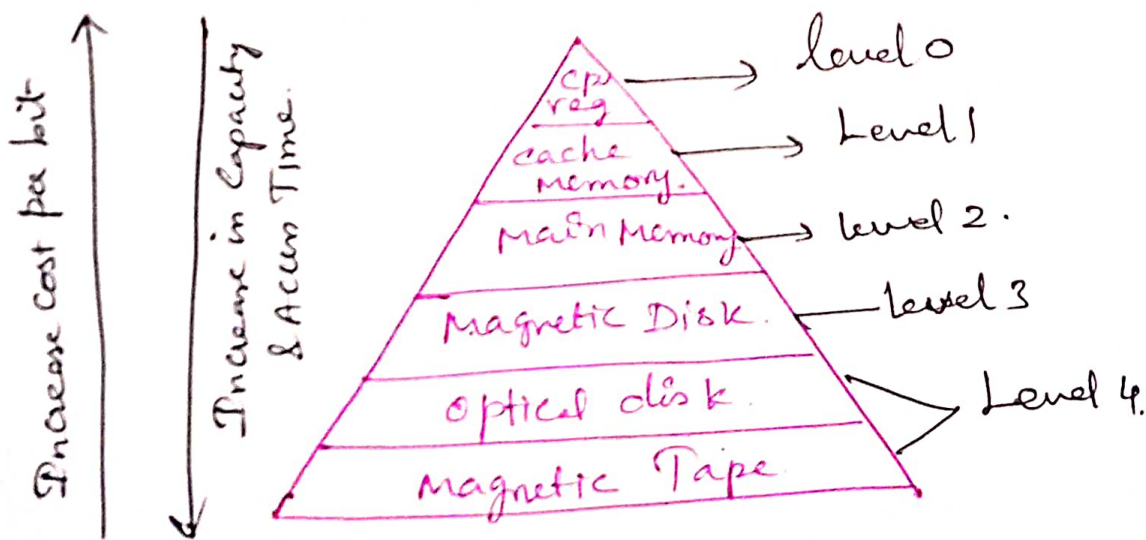


Main Memory is implemented from DRAM while levels closer to the processor (caches) use SRAM (Static Random Access Memory).

→ DRAM is less costly per bit than SRAM.

→ DRAM uses significantly less area per bit of memory and DRAM.

→ Third technology used to implement the largest and slowest level in the hierarchy is magnetic disk.



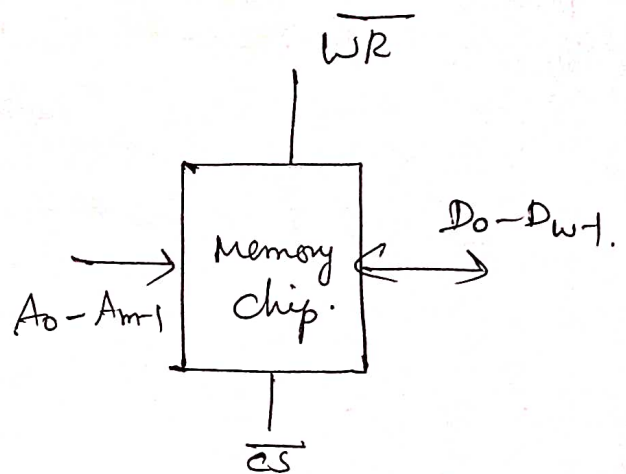
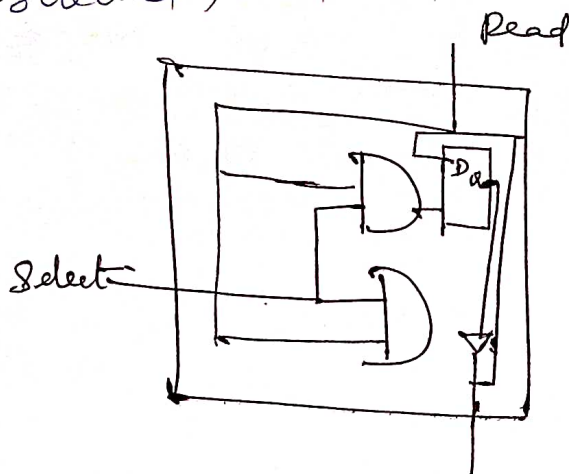
Memory Hierarchy Design.

# Memory chip Organization.

RAM → Random Access Memory.

→ Term Random means that any memory location can be accessed in the same amount of time.

→ fig represent the memory elements as D flip flop with additional control to allow the cell to be selected, read and written.



→ This topic mainly focus on how RAM cells are organized into chips and how this organization implements the Random access property.

→  $m$  bit address, having lines numbered from 0 to  $m-1$  is applied to pin  $A_0 - A_{m-1}$  while asserting chip select ( $cs$ ) and either  $\overline{WR}$  (for reading data from the chip).

→ The over bar on  $cs$  and  $\overline{WR}$  indicates that the chip is selected when  $cs=0$  and that write operation will occur when  $\overline{WR}=0$ .

→ When reading data from the chip after time period  $t_{SA}$  the  $w$  bit data word appears on data line  $D_0 - D_{w-1}$

→ Consider that a  $64 \times 1$  chip has 26 address lines ( $64M = 2^{26}$ ). This means that a conventional decoder would need  $2^{26}$ , 26 input and AND gates.



## Mapping function.

A particular block of main memory can be brought to a particular block of cache memory.

→ Mapping functions are used to map a particular block of main memory to a particular block of cache.

→ This mapping function is used to transfer the block from main memory to cache memory.

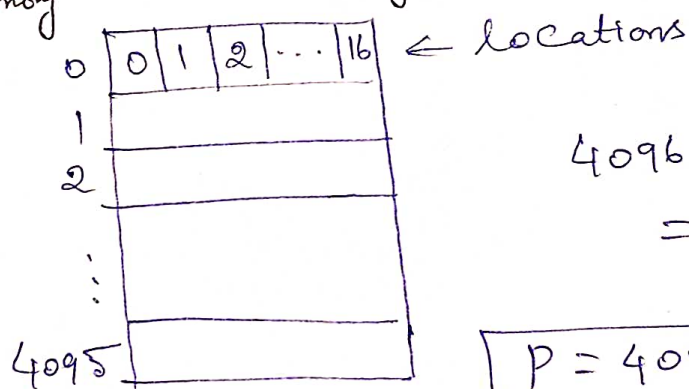
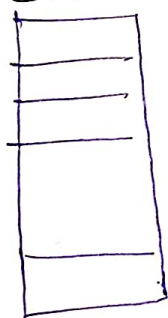
Three types of mapping function.

- ① Direct mapping.
- ② Associative mapping.
- ③ Set associative mapping.

## Direct mapping

A particular block of main memory can be brought to a particular block of cache memory.

Cache memory

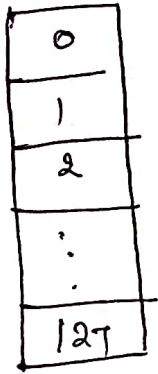


$$4096 \times 16 = 65536 = 2^{16}$$

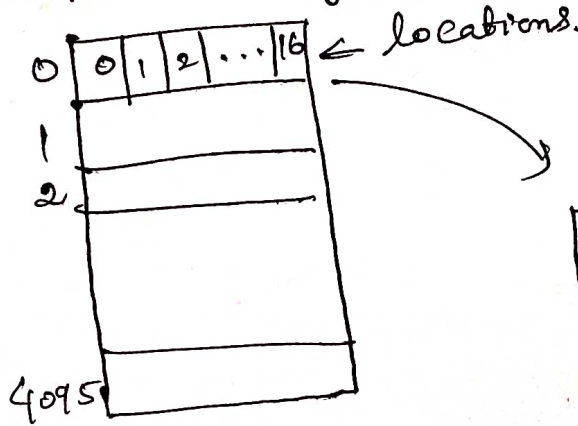
$$P = 4096$$

each page have 16 locations

Cache memory



Main Memory

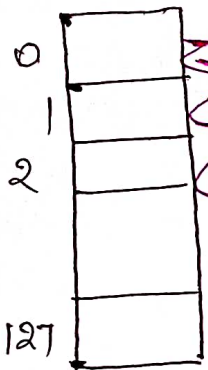


$C = 128$

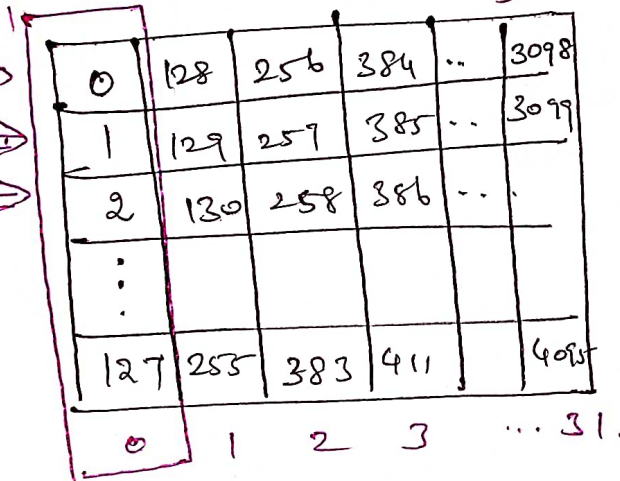
$P = 4096$

Each Cache frame size equal to each block size of main memory.

Cache memory



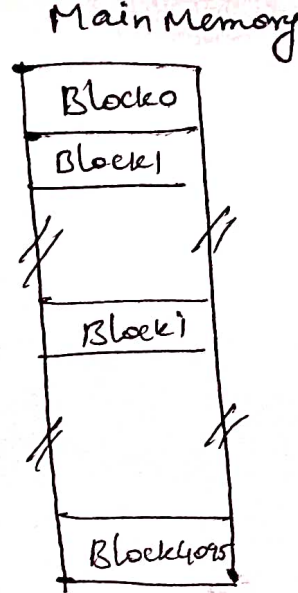
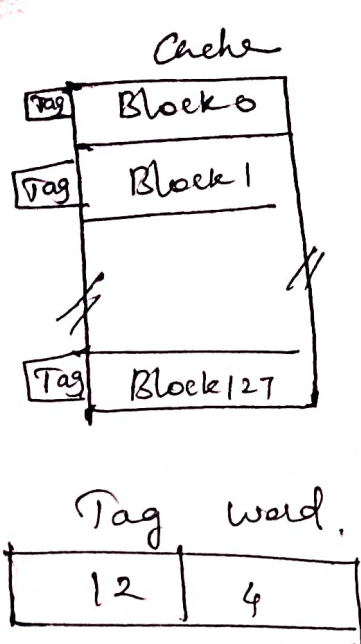
Main Memory



### Associative Mapping

- In this mapping function any block of main memory can potentially reside in any Cache block position.
- This is much more flexible mapping method.
- The tag bits of an address received from the Processor are compared to the tag bits of each block of the cache to see if the desired block is present. This is called the associative mapping technique.





### Associative mapped Cache.

Advantage. → space in the cache can be used <sup>efficiently</sup>

→ it gives full freedom in choosing a cache location to load memory block.

### disadvantage

→ cost of associative mapping cache is higher than Direct mapping cache.

### Set associative mapping

→ This mapping is a combination of the direct mapping and associative mapping techniques.

→ Blocks of the cache are grouped into sets and mapping allows a block of the main memory to reside in any block of a specific set.

# Virtual Memory Organization

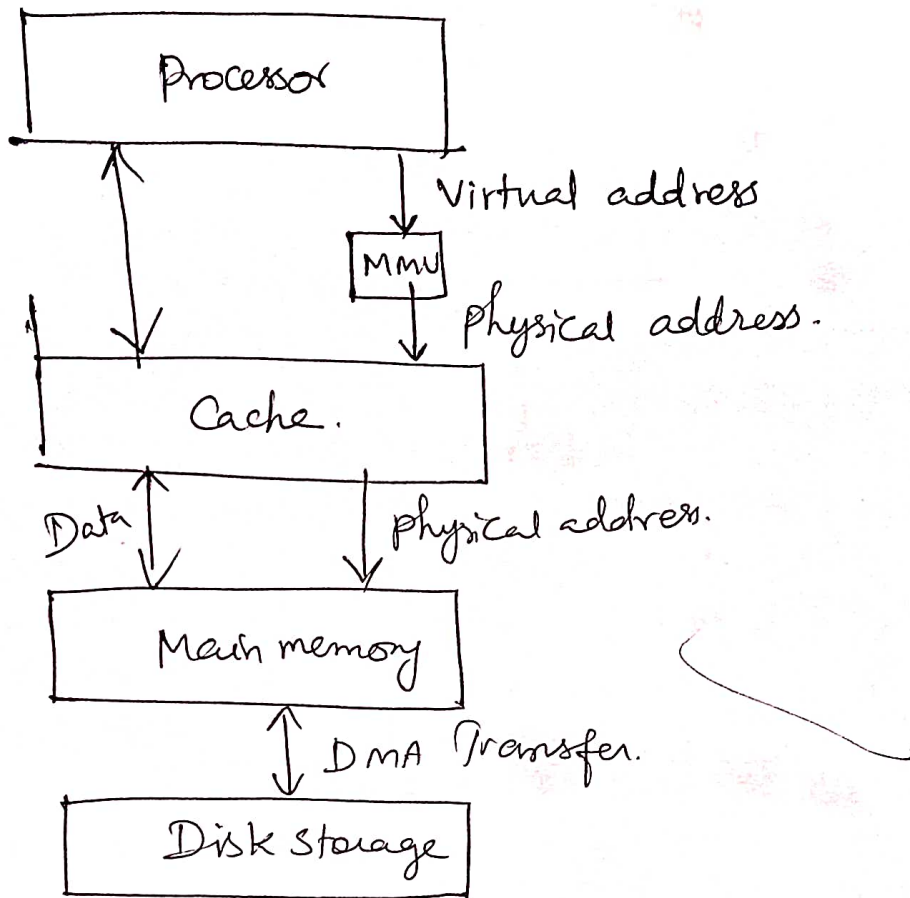
- The term virtual memory refers to something which appears to be present but actually it is not.
- The virtual memory technique allows user to use more memory for a program than real memory of computer.
- So virtual memory is the concept that gives the illusion to the user that they will have main memory equal to the capacity of secondary storage media.

## Need of virtual memory

- The program is stored in the secondary memory. The memory management unit (MMU) transfers the currently needed part of the program from the secondary memory to the main memory for execution.
- The segments which are currently being executed are kept in the main memory and remaining segments are stored in the secondary storage devices. Such as magnetic disk.
- If an executing program needs a segment which is not currently in the main memory, the required segment is copied from the secondary storage device.



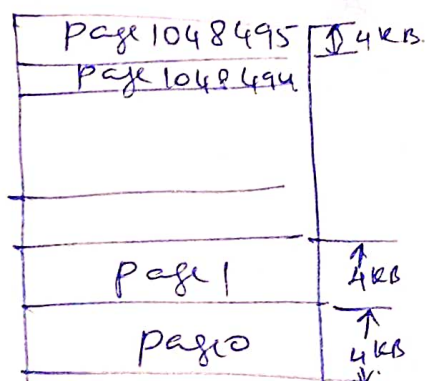
→ In modern Computers. the operating system moves program and data automatically b/w main memory and secondary storage.



### Virtual memory Organization.

Paging: The memory Management Unit (MMU) controls the virtual memory system. It translates virtual address into physical addresses.

→ A simple method for translating virtual address into physical addresses is to assume that all programs and data are composed of fixed length unit called pages.



Bus It is a Common path way that connects a number of devices.

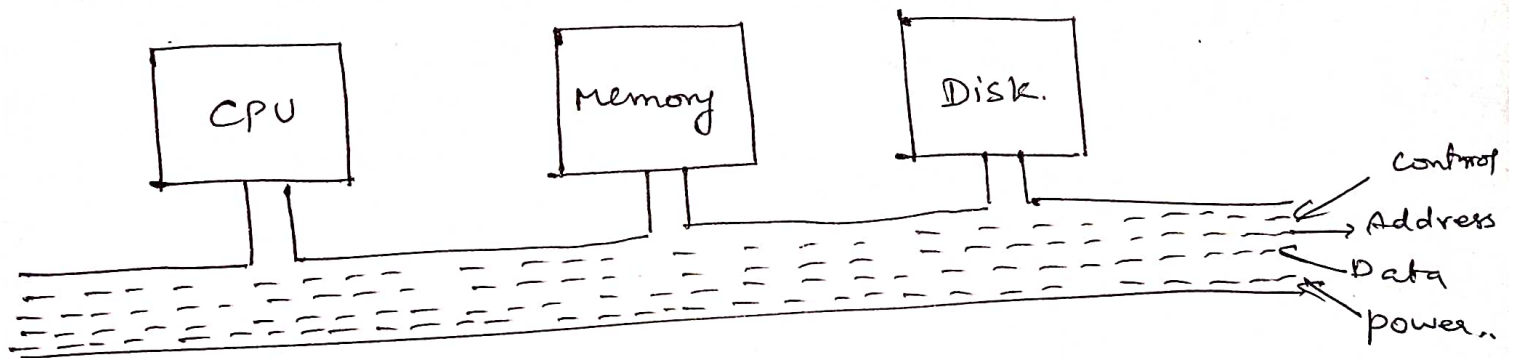
→ example of a bus can be found on the motherboard on the personal computer.

→ motherboard contains integrated IC such as CPU chips and memory chip board traces (wires) that connect the chip, a number of buses for chip and devices that need to communicate with each other and bridges that connect the buses.

### Bus structure, Protocol and Control

→ A bus consists of the physical parts like connectors and wires, and a bus protocol.

→ The wires can be partitioned into separate groups of control, address data and power.



Simplified illustration of a bus.

→ Single bus has few power lines such as GND at 0V. Positive voltage at +3.3V. positive and negative voltages at +5V and -5V.

→ All devices simultaneously listen but normally only one device receives.



# Address Translation

→ In virtual memory the address is broken into a virtual page number and page offset.

→ A simple method for translating virtual addresses into physical addresses is to assume that all program and data are composed of fixed length units called pages.

→ Pages commonly range from 2k to 16k bytes in length.

## Parallel Bus architecture

A computer system may contain many components that need to communicate with each other. In worst case scenario, all  $N$  components (CPU, graphics processor, CD ROM, magnetic disk, etc) may need to simultaneously communicate with every other component in which  $N \times (N-1) / 2$  links are needed for  $N$  components.

Example: Just  $N=10$  components. (CPU, graphics processor, a few memory modules, a few network interface cards, CD ROM) with a 64 bit datapath.

→ The number of wires needed to support this is  $\frac{10 \times 9}{2} \times 64 = 2880$ . which is enormous Degree of Complexity.

→ Only one device can be bus master at any given time and the remaining devices are slaves

## Advantages of Bus

To eliminate the need for connecting every device with every other device. Which avoids the wiring complexity.

## Disadvantages

① Slowdown introduced by the master/slave Configuration.

## Types of Bus

A bus can be classified into two types.

(a) Synchronous Bus

(b) Asynchronous Bus.

## Internal Communication Methodologies

→ Computer systems have a wide range of Communication tasks.

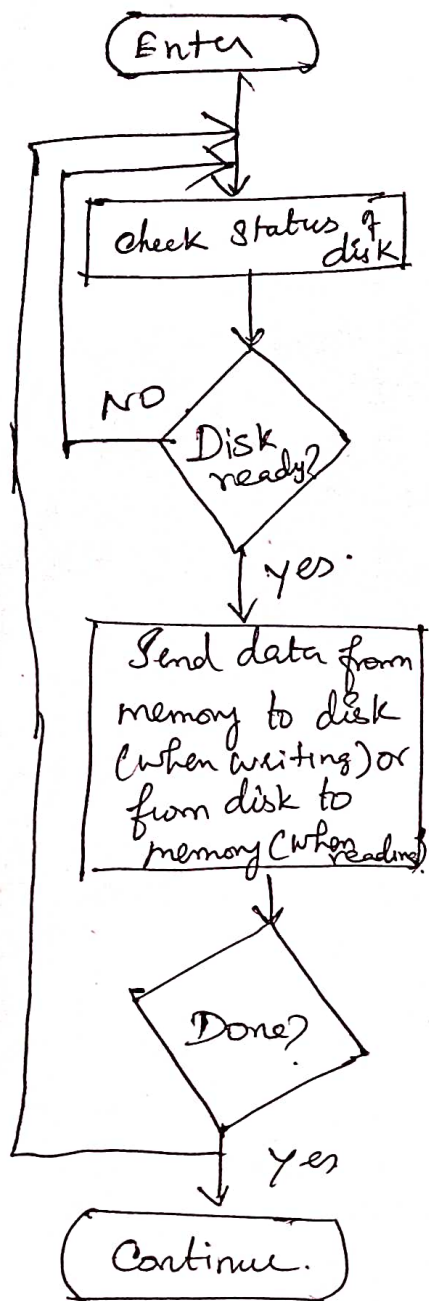
→ The CPU must communicate with memory and I/O devices such as keyboard, disk drives, and network interfaces. There are three different types of data transfer schemes are used.

(i) Programmed I/O. (ii) Interrupt I/O. (iii) Direct Memory Access (DMA).

## Programmed I/O (Polling)

→ Consider reading a block of data from a disk. In programmed I/O, the CPU polls each device to see if it needs servicing.





→ CPU first checks the status of disk.

→ If the disk is not ready to be read or written, then the process loops back and checks the status continuously until the disk is ready.

→ When the disk is finally ready, then a transfer of data is made between the disk and CPU.

### Advantages

- simple to implement
- Very little hardware support.

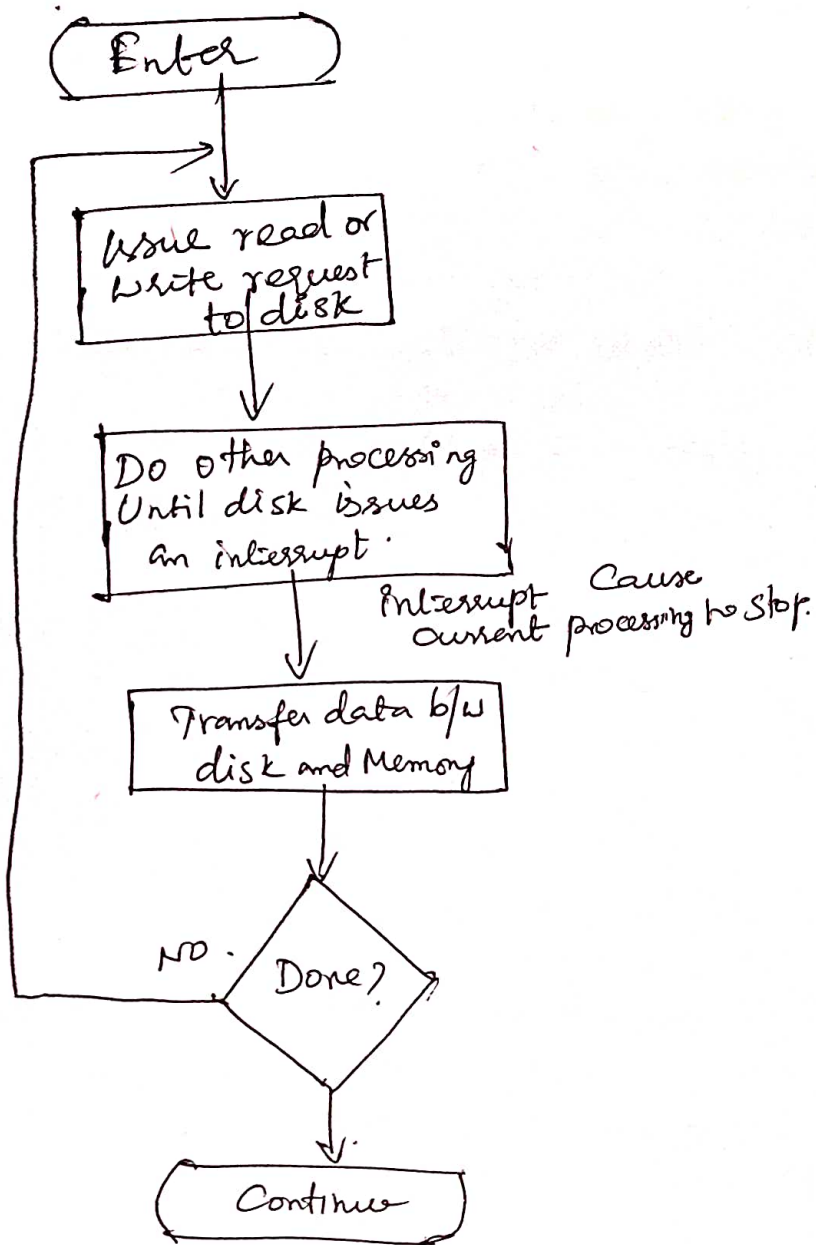
### Disadvantages

- CPU is kept busy for checking the status of the slower I/O device.
- busy waiting.

## Interrupt Driven I/O

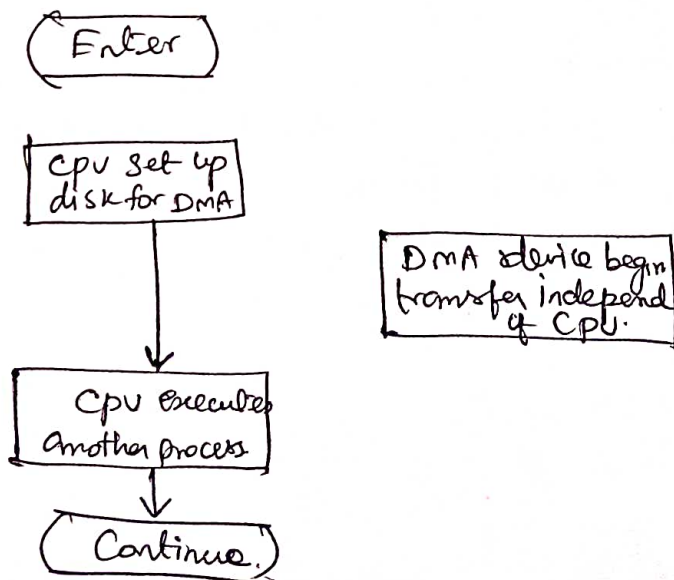
→ With interrupt driven I/O the CPU does not access a device until it needs servicing and so it does not get caught up in busy wait.

→ In interrupt driven I/O, the device request goes directly to the CPU through a special interrupt request line that



## Direct Memory Access

→ Interrupt driven I/O frees the CPU until the device requires services. The CPU is still responsible for making actual data transfer.





→ A direct memory Access (DMA) device can transfer data directly to and from memory rather than using the CPU as an intermediary and can thus relieve congestion on the system bus.

→ DMA Services are usually provided by a DMA Controller, which is itself a specialized processor whose specialty is transferring data directly to or from I/O devices and memory.

### Serial Bus architecture

→ Serial communication (bit by bit) have significant advantages over parallel communication.

→ Serial interconnects requires higher bit rates to compensate for the lack of parallelism.

→ It consumes less power, less space, and cost less money.

### Serial Communication Approaches

#### RS232

→ The first serial interconnection method is not a bus at all, but a simple point to point interconnect known as recommended standard 232C, or RS-232.

→ RS 232 standard commonly uses 9 pins and 25 pin connectors.

→ Pin 7 of 25 pin connector is for signal ground pin 2 send data and pin 3 receives data.

## Universal Serial bus

- USB and IEEE 1394 (fire wire) are two groups of standard for interconnecting peripheral devices.
- It initially supported data transfer rates of 12 Mbps with as many as 127 devices connected to a single host controller through special hub devices.
- The standard has evolved to USB 2.0 which supports data rates up to 480 Mbps. USB is typically used for mice, keyboard, modems and other devices.

## Mass Storage

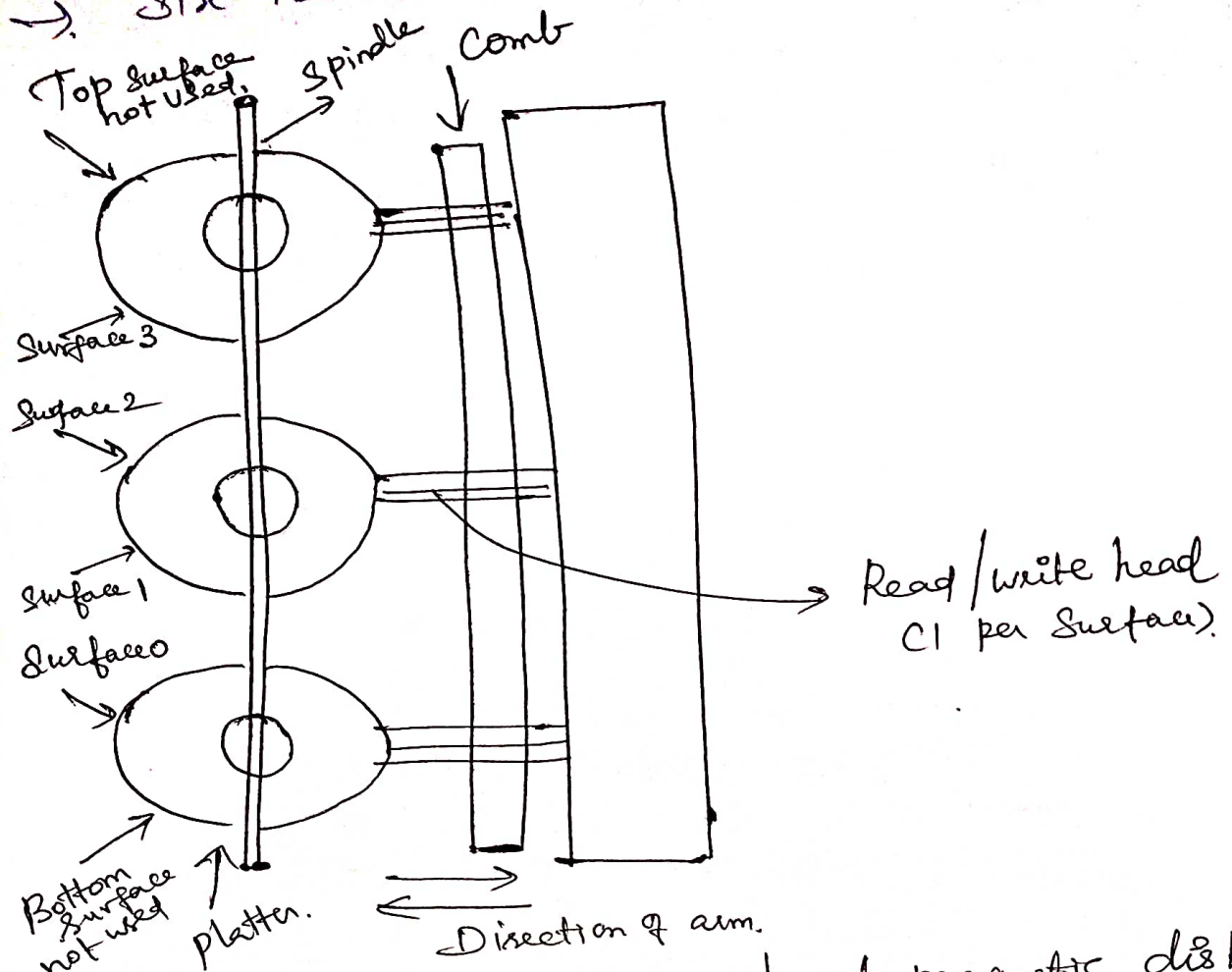
- Registers and RAM requires continuous power to retain their stored data, whereas media such as magnetic tapes and magnetic disks retain information indefinitely after the power is removed. which is known as indefinite persistence.
- This type of storage is said to be non-volatile.

## Magnetic Disks

- A magnetic disk is a device storing information that supports a large storage density and a relatively fast access time.
- A moving head magnetic disk is composed of a stack of one or more platters that are spaced close together and are connected via spindle.
- The most common size for desktop personal computer disks is 3.5" diameter.
- Laptop hard disk drives use 2.5" diameter or smaller.



- A single head is dedicated to each surface.
- Six heads are used in the example shown.



multiplatter hard magnetic disks.

- The heads are attached to a common arm (also known as comb), that moves in and out to reach different portions of the surfaces.
- In a hard disk drive, the platters rotate at a constant speed of typically 5400 to 15000 magnetic revolutions per minute (RPM).

Magnetic Tape:

→ The magnetic tape unit typically has a single read/write head, but may have separate heads for reading and writing.

## Input devices.

- Keyboard: → A keyboard is basically a board of keys. Keyboard is one of the primary input devices used with a computer.
- This keyboard layout is known as "QWERTY" layout, which gets its name from the first six letters across in the upper left hand corner of the keyboard.
- A keyboard layout using ECMA-23 Standard is shown.
- When a character is typed, a bit pattern is created that is transmitted to a host computer.
- ECMA-23 keyboards include additional of function keys (in row F, for example), and the addition of special keys such as tab, delete and carriage return.

## Tablets

- A digitizing tablet is an input device that consists of a flat surface and either a pen based stylus or puck.
- Tablet has an embedded two dimensional mesh of wires that detects an induced current created by the puck as it is moved about the tablet.
- The tablet transmits x-y (horizontal-vertical) positions and the state of the button on the puck (or stylus), either continuously, or for an event such as a key click or a element, depending on the control method.



## Touch Sensitive Pen-Based display.

Unlike a tablet which uses an active matrix that senses an electromagnetic field induced by a stylus or puck. Pen based Personal digital assistants (PDA) uses a passive matrix in which the pen can be anything that induces pressure on the screen.

→ Two transparent conducting layers are placed on the screen separated by a spacer dots. When the user applies pressure to the top layer as with a stylus or simply a finger the top and bottom layers make contact. The induced voltage at the edges varies according to the position of the stylus.

Joysticks → A joystick indicates horizontal and vertical position by the distance a rod that points from the base is moved.

→ Joysticks are commonly used in video and for indicating position in graphics systems.

## Output devices

Three common output devices: Laser Printer, Video display and LCD Panel.

Laser Printer. → A laser printer consists of a charged drum in which a laser discharges selected areas according to a bit-mapped representation of a page to be printed.

→ As the drum advances for each scan line the charged areas pick up electrostatically sensitive toner powder.

→ The drum continues to advance and toner is transferred to the paper which is heated to fix the toner on the page.

## Video displays

→ The display devices are known as O/P devices. The most commonly used output devices in a graphics video monitor.

→ The operations of most video monitors are based on the standard cathode ray tube (CRT).

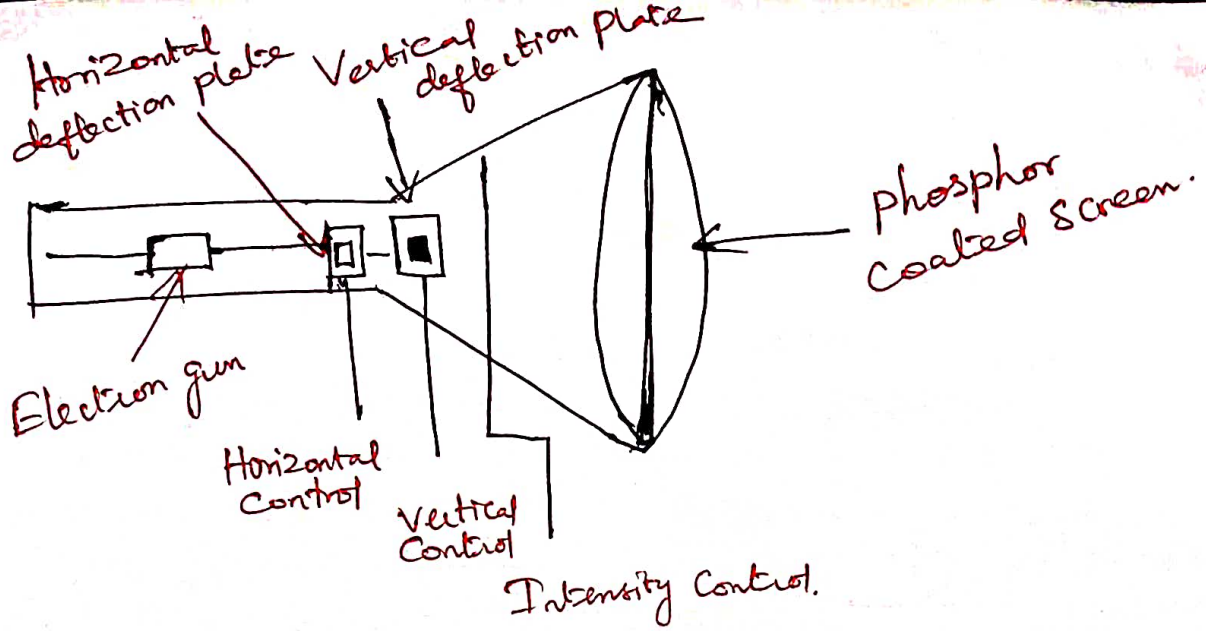
→ A video display, or monitor consists of a luminescent display device such as a cathode ray tube (CRT).

→ In a CRT, vertical and horizontal deflection plates steer an electron beam that sweeps the display screen in a raster fashion.

→ A configuration for a CRT is shown in fig.

An electron gun generates a stream of electrons that is imaged onto a phosphor-coated screen at positions controlled by voltages on the vertical and horizontal deflection plates. Electrons are negatively charged and so a positive voltage on the grid accelerates electrons towards the screen and negative voltage repels electrons away from the screen.





CRT with a single electron gun.

→ The colour produced on the screen is determined by the characteristics of the phosphor.

→ For a colour CRT there are three different phosphor types (red, green, blue) that are interleaved in a regular pattern, and three guns, which produce three beams that are simultaneously deflected.

### Liquid Crystal Displays (LCD's)

→ We always use devices made up of Liquid Crystal Displays (LCD) like computers, digital watches, also DVD and CD players.

→ Liquid Crystal Displays (LCD) have emerged as a prevalent choice of computer monitors. LCDs are small and light weight when compared with tube monitors, consume less power and dissipate less heat and are more shock resistant.

Advanced Computer Architecture

One method of improving the performance of a processor is to increase the clock rate. This will work up to a point because higher the clock rate consumes more power.  
→ An alternate approach is to increase the number of processors and decompose and distribute a single program onto the processors. This approach is known as parallel processing. In parallel processing a number of processors work collectively in parallel on a common problem.

Parallel Processing architectures and challenges.

At first the performance metrics for parallel processing is explained followed by parallel processing challenges and then specific types of parallel architectures are discussed.

Measuring Performance

A parallel architecture can be characterized in terms of three parameters.

- ① The number of processing elements.
- ② The interconnection networks among the PEs.
- ③ The organization of the memory.

→ Parallel time . It is the absolute time needed for a program to execute on a parallel processor.

→ Speedup It is the ratio of the time for a program to execute on a sequential processor to the time for that same program to execute on a parallel processor.



③ Efficiency It is the ratio of Speed up to the number of Processors used.

④ Throughput It is a measure of how much Computation is achieved over time.

### Parallel Processing Challenges.

High Performance can mean high throughput for independent task. Utilizing multiple Processors by running independent Programs simultaneously & called as task level parallelism. or Process level parallelism.

→ A single Program that runs on multiple Processors simultaneously is called as Parallel Processing Program.

→ Some of the challenges of Parallel Processing are

- \* Parallel Programming
- \* Complex algorithm.
- \* Memory issues.
- \* Communication costs.
- \* Load Balancing.

→ Some Problems do not have work efficient Parallel algorithms that exhibit massive parallelism.

That is we simply donot know how to solve these Problems with a large number of Parallel execution Units without significantly increasing the Computation Complexity.

→ Some Problems have known parallel algorithms with sufficient parallelism. but questionable numerical stability.

→ Some applications that have parallel algorithms are Overwhelmed by catastrophic load imbalance due to highly non Uniform data distribution.

## Hardware Multi-threading.

- (2)
- Multithreading is first introduced to exploit thread level parallelism within a processor. Hardware multi-threading allows multiple threads to share the functional units of a single processor in an overlapping fashion.
  - It tries to utilize the hardware resources efficiently. To allow this sharing, the processor must duplicate the independent state of each thread.
  - Hardware multi-threading increases the utilization of a processor by switching to another thread when one thread is stalled. Thread is a lightweight process.
  - All multi-threading approaches use similar resources partitioning system. Although there are similarities in the multi-threading implementations, there are still some differences. Two major differences are.
    - (1) Thread scheduling policy.
    - (2) Pipeline partitioning.
  - Based on above criteria there are three main approaches to multi-threading.
    - (1) Fine grained multi-threading (FGMT)
    - (2) Coarse grained Multi-threading [CGMT]
    - (3) Simultaneous multi-threading [SMT].

## Fine-grained Multi-threading.

- Fine grained multi-threading switches between threads after every instruction. It makes interleaved execution of multiple threads at the same time. This interleaving is done in round robin fashion on every clock cycle.



advantage → It hides throughput losses due to short and long stalls.  
→ Instructions from other threads can be executed when one thread stalls.

### Disadvantage

→ It slows down the execution of individual threads.  
→ Thread that is ready to execute without stalls will be delayed by instructions from other threads.

### Coarse-Grained Multithreading

Coarse grained multithreading is an alternative to fine grained multithreading. Coarse-grained multithreading won't switch out the executing thread until it reaches a situation that triggers a switch. This situation occurs when the instruction execution reaches either a long latency operation or an explicit additional switch operation.

#### advantages:

→ Faster, Reduces the penalty of high cost stalls.  
→ Less likely slow down the execution of an individual thread.

#### Disadvantage

→ Difficult to overcome throughput losses due to shorter stalls.  
→ When a stall occurs the pipeline must be emptied.

### Simultaneous Multithreading (SMT)

→ Hardware multithreading that uses the resources of multiple issues. dynamically scheduled pipeline processor is called Simultaneous Multithreading. Here both thread level and instruction level parallelisms are used.

#### Advantage

→ More functional unit parallelism available than single thread.  
→ Register renaming and dynamic scheduling are used.  
→ It does not switch resources for every cycle.

# Multicore and Shared Memory Multiprocessors (3)

Hardware multi-threading improved the efficiency of processors at modest cost, but the difficulty is how to run old programs in parallel hardware.

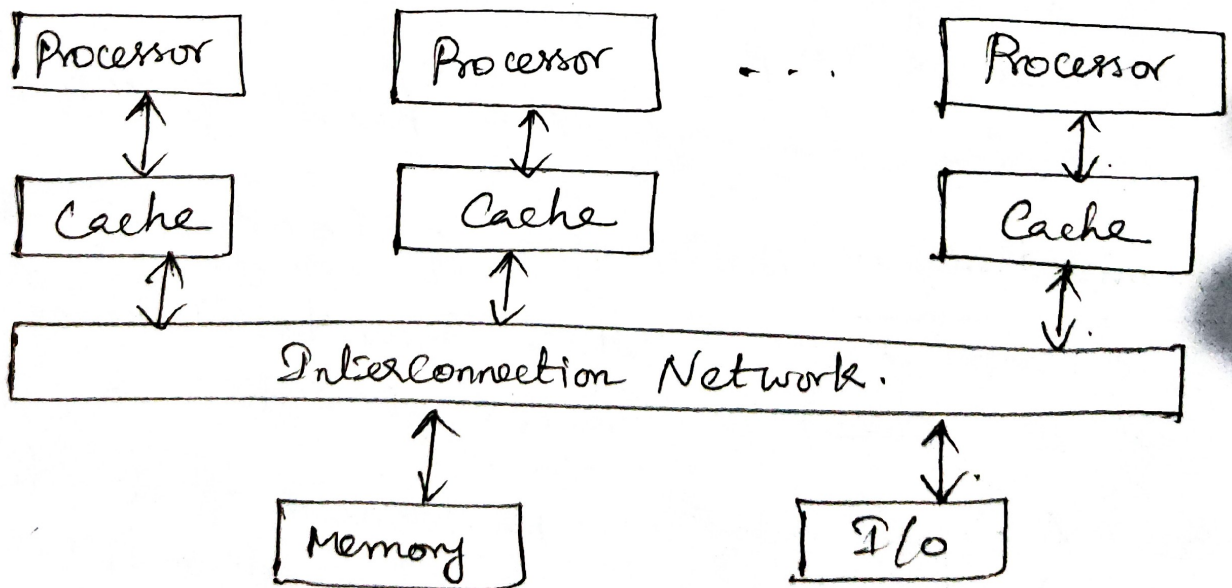
There are two solutions:

(i) Shared Memory Multiprocessor (SMP) Provides single address space to all processors. In this case, program need not to concern about where they are run. They may be executed in parallel. In this approach, all variables of the program must be made available at any time to any processor.

(ii) Clusters and other message passing multiprocessors. Separate address space is reserved for each processor. Here sharing must be explicit.

## Shared Memory Multiprocessor.

→



Organization of Shared memory multiprocessor



→ A parallel processor with single physical address space is called Shared Memory Multiprocessor.

→ Processors communicate through shared variable in memory. All the processors can access any memory location via loads and stores. Fig shows the organization of SMP.

→ There are two types of single address space multiprocessors.

### (i) Uniform Memory Access. (UMA)

Latency to a word in memory does not depend on which processor asks for it. That is the time taken to access a word is uniform for all processors.

### (ii) Non Uniform Memory Access. (NUMA)

→ Some memory access are much faster than others depending on which processor asks for which word.

→ Here the main memory is divided and attached to different microprocessors, or to different memory controllers on the same chip.

→ Programming challenges are harder for NUMA than UMA.

Synchronization. The process of coordinating the behavior of two or more processes which may be running on different processors, are called Synchronization.

Lock A synchronization mechanism that allows access to data to only one processor at a time is called lock. Other processor must wait until the original processor unlocks the variable.

# Introduction to Graphics Processing Units

- The main reason for Improving graphics Processing was the computer game industry. The fast growing game market encouraged many companies to make increasing investments in developing faster graphics hardware.
- The graphics Processing Improves at a faster rate than general purpose Processing in mainstream microprocessors. As the graphics Processors increased in power they earned the name Graphics Processing Unit or GPU to distinguish themselves from CPUs
- Anyone can buy a GPU, today which hundreds of Parallel floating point units at reasonable cost. It makes high performance Computing more accessible.

## GPU Vs CPU

There are a few key characteristics as to how GPU vary from CPU's

1. GPUs are accelerators that supplement a CPU, so they need be able to perform all the tasks of a CPU. This role allows them to dedicate all their resources to graphics.
2. The GPU problem sizes are typically hundred of megabytes to gigabytes where as the CPU problem size are ranging from hundreds of gigabytes to terabytes.

## An Introduction to the NVIDIA GPU Architecture.

- C inspired programming languages are used to write programs directly for the GPUs.
- Thousands of CUDA threads are combined together to utilize the various styles of parallelism within a GPU. Multithreading, MIMD, SIMD and instruction-level parallelism. These threads are blocked together and executed in groups of 32 at a time.

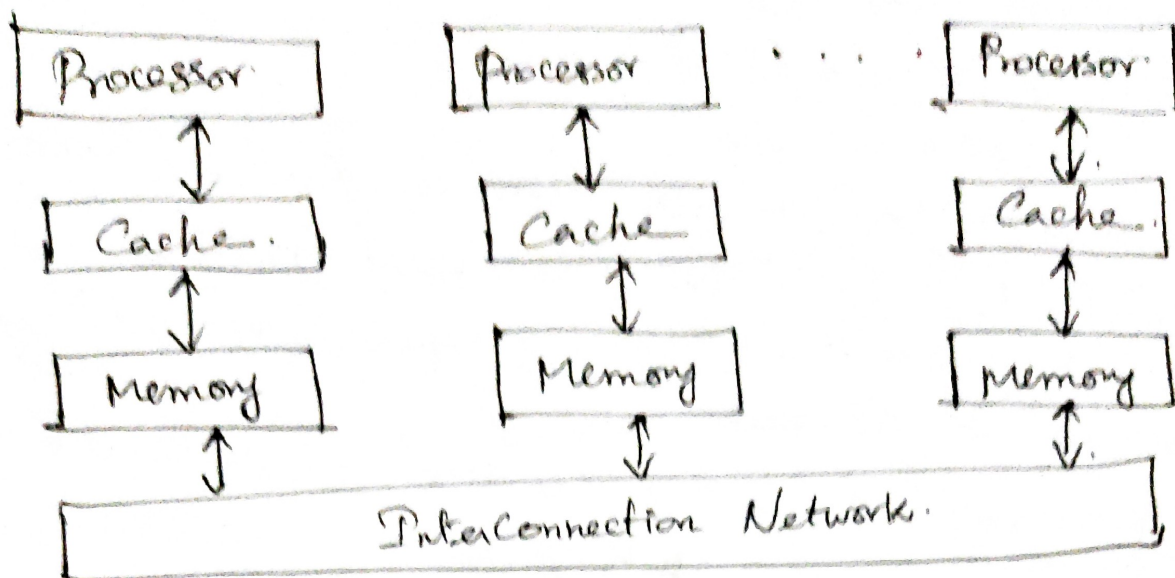


- Multicore Computers with SIMD instruction extensions have similarities with GPUs.
- SIMD Processors are also similar to vector processors. The multiple SIMD processors in GPUs act as independent MIMD cores, just as many vector computers have multiple vector processors.
- The biggest difference is multithreading which is fundamental to GPU and missing from most vector processors.

## Clusters and Warehouse Scale Computers.

### Cluster and other message passing multiprocessors.

- An alternative approach to sharing an address space is for the processors to have their own private physical address space. Figure shows the organization of a multiprocessor with multiple private address space. In this each processor have their own private physical address space.
- In this approach the multiprocessor communicate via explicit message passing scheme, hence the name message passing multiprocessors, communications b/w multiple processors is done by explicitly sending and receiving information.
- There are routine to send and receive message. Send message routine is used to send message to other processors in machines with private memories. Receive message routine is used to receive messages from another processors in machine with private memories.



Organization of message passing multiprocessor.

### Drawbacks of clusters.

- ① The cost of administering a cluster of  $n$  machines is about the same as the cost of administering  $n$  independent machines. But the cost of administering a shared memory multiprocessor with  $n$  processors is about the same as administering a single machine.
- ② The processors in a cluster are connected using the I/O interconnect of each computer. But the cores in a multiprocessor are usually connected on the memory interconnect of the computer.
- ③ The overhead in the division of memory, a cluster of  $n$  machines has  $n$  independent memories and  $n$  copies of the operating systems are needed.

### Warehouse scale Computers Vs Servers.

WSC have three major distinctions with servers.

1. More computational parallelism is not important
  - Most jobs are totally independent
  - Many independent jobs can proceed in parallel with little need for communication or synchronization

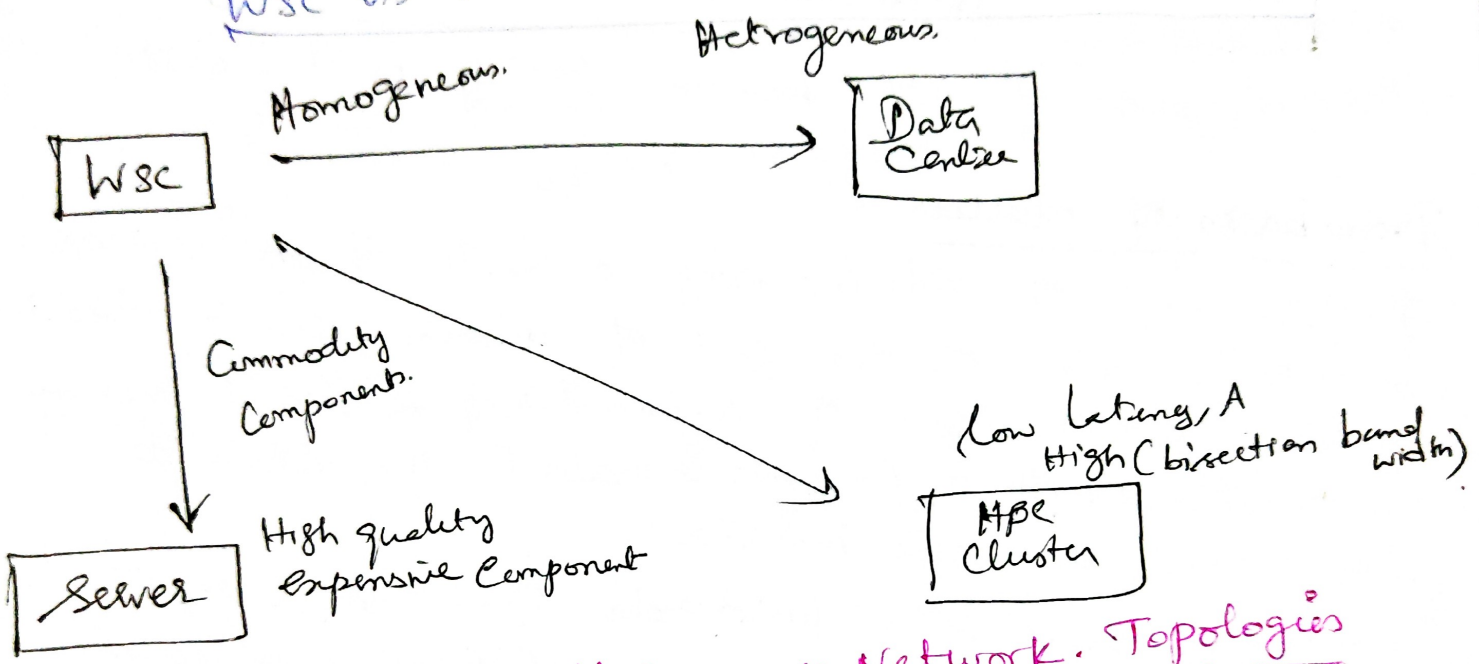


## ② Operational Cost Count

→ WSC have longer lifetimes; so that the Operational Costs increase.

③ Scale and the opportunities / Problems Associated with scale.  
To construct a single WSC, you must purchase 100,000 servers along with the supporting infrastructure, which means volume discounts.

## WSC Vs Data Center Vs HPC. Cluster Vs Server



## Introduction to Multiprocessor Network Topologies

→ MultiCore chips require on chip networks to connect cores together and clusters require local area networks to connect servers together. The performance metrics used to analyze the topologies are, network bandwidth and bisection bandwidth.

→ Network bandwidth is defined as the peak transfer rate of a network. It refers to the speed of a single link or the collective transfer rate of all links in the network.

→ Bisection bandwidth is defined as the bandwidth between two equal parts of a multiprocessor. This measure is for a worst case split of the multiprocessor.

→ There are number of different topologies available but only a few have been used in commercial parallel processors.

### → Multi stage network topologies

→ An alternative to placing a processor at every node in a network & to leave only the switch at some of these nodes.

→ The switches are smaller than processor-memory switch nodes and thus may be packed more densely. There by lessening distance and increasing performance.

→ Such networks are frequently called multistage networks to reflect the multiple steps that a message may travel.

→ Simply a network that supplies a small switch at each node & called as multistage network.

→ A network that allows any node to communicate with any other node in one pass through the network is called as cross bar network.

① Cross bar \* The crossbar uses  $n^2$  switches where  $n$  is the number of processors.

\* Here cross bar uses 64 switches versus 12 switch boxes.

\* It supports any combination of message b/w

Processors.

② Omega Network \* The Omega network uses  $2n \log_2 n$  of the large switch boxes.

\* Here Omega uses 48 switches.

\* It uses less hardware than the cross bar network.

\* It does not support any combination of messages between processors.