

## # Pipelining.

• Need of pipelining.

→ CPU performance must be high is today's requirement.

How it is possible?

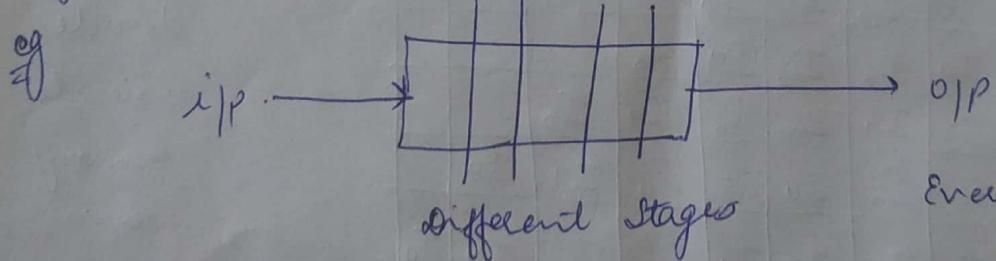
1) Circuit may be change and install the fastest circuit in CPU.

2) Hardware is arrange in such a way that multiple instns execute at a time.

If i<sub>3</sub> → i<sub>5</sub> → i<sub>7</sub> cost 4 sec

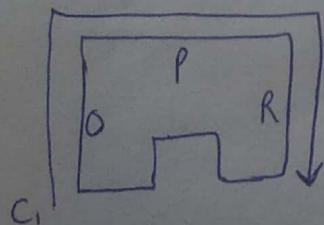
so from cost pt. of view, which method is used if such that cost will not ↑ se and existing hardware will change to ↑ se the performance. For this purpose, pipelining is used

Pipelining help to not to change the clk, clk must be designed as a pipe.



Every stage depends on each other.

Define through



$$\begin{array}{c} C_1 \rightarrow O | P | R \\ C_2 \rightarrow - | O | P | R \\ C_3 \rightarrow - | - | O | P | R \end{array}$$

$$\frac{5}{3} = 1.6$$

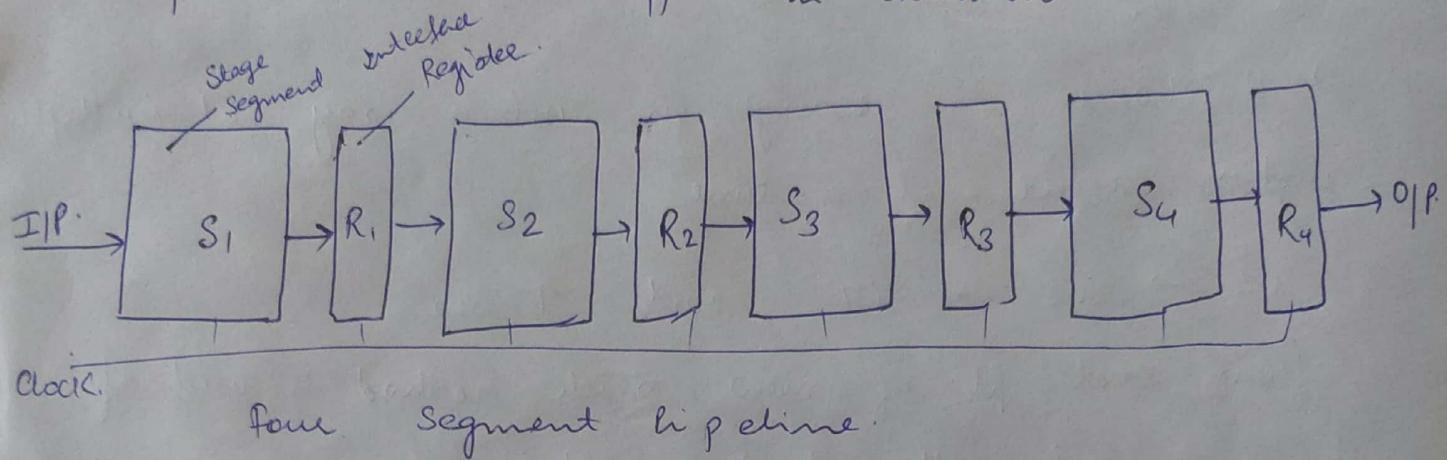
$$\begin{array}{c} C_1 \rightarrow O | P | R | \\ C_2 \rightarrow - | - | - | O | P | R | \\ C_3 \rightarrow - | - | - | - | - | O | P | R | \end{array}$$

$$\frac{9}{3} = 3$$

4.

## • Definition of Pipelining:

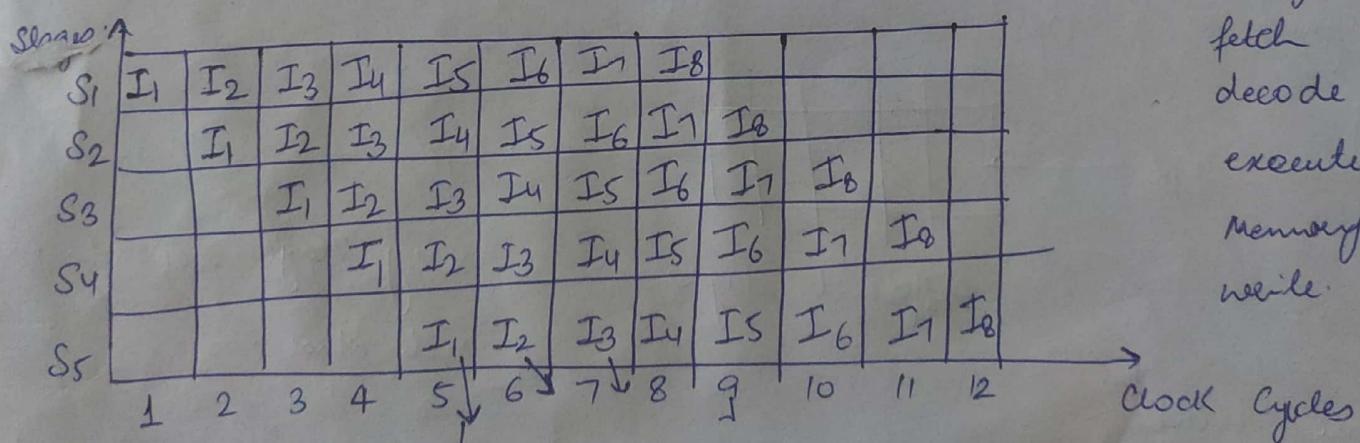
- Pipelining is a process of arrangement of hardware elements of CPU such that its overall performance is 1 sec.
- Simultaneous execution of more than one instruction takes place in pipelined processor.
- Multiple instrn are overlapped in execution.



Space time diagram:

Instruction execution in pipeline architecture.

$$P = 8 I \quad I_1 - I_8$$



$$1 \text{ instrn} - 5$$

$$8 - 5 \times 8 = 40 \text{ clock cycles}$$

Total Instr = n.

No. of stages = K.

First instr will take =  $1 \times K$

Reset  $(n-1)$

Total No. of clock cycle required =  $K + (n-1)$

$$5 + (6-1) = 12$$

Clock per Instr

CPI  $\cong 1$  AIM of Pipelined.

of 1000 instr.

$$5 + (1000 - 1) = 1004 \approx 1$$

$t_n t_p \rightarrow$  clock cycle time

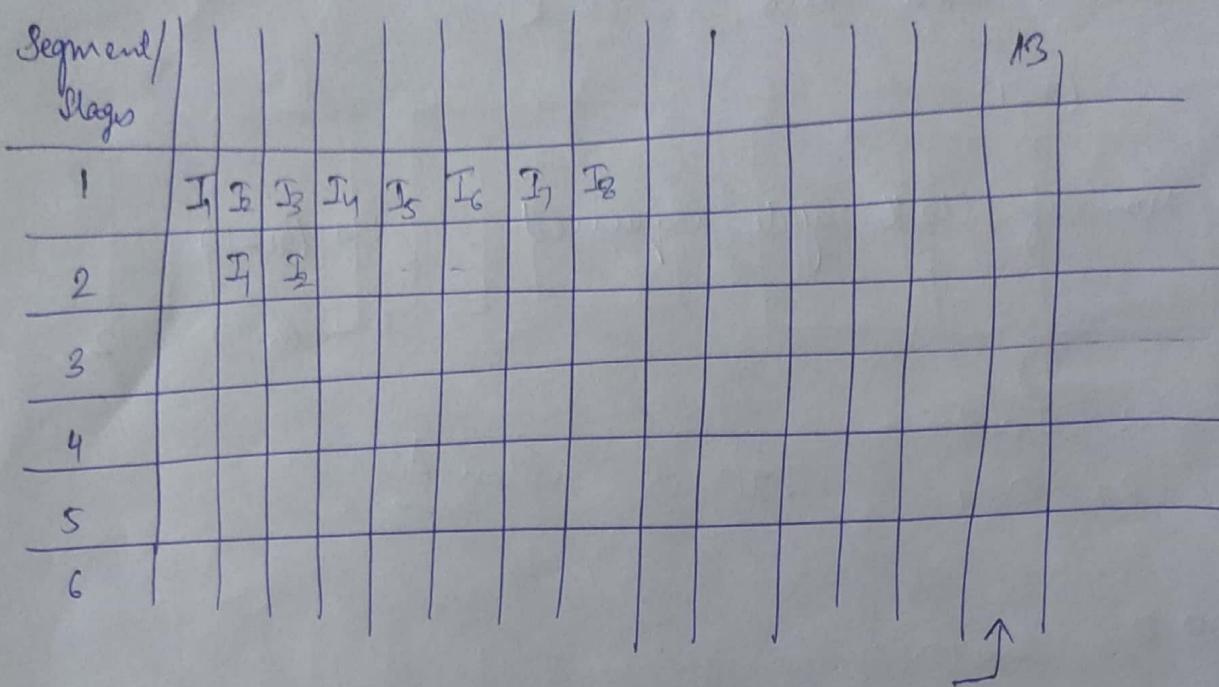
$$\text{Speed up} = \frac{NP}{P} = \frac{n t_n}{(K+n-1) t_p}$$

$$= \frac{40}{12} = 3 \quad \checkmark$$

$$\text{Efficiency} = \frac{\text{Used}}{\text{Total Box}}$$

$$= \frac{40}{60} = \frac{2}{3}. \quad \checkmark$$

D Draw a space-time diagram for six-segment pipeline showing the time it takes to process eight task.



$$K + m - 1 = 6 + 8 - 1 = 13 \text{ cycles}$$

D A non-pipelined system takes 50 ms to process a task. The same task process in six-segment pipeline with clock cycle of 10 ns. Determine speed up ratio of pipeline for 100 tasks

$$t_n = 50 \text{ ms}, \quad K = 6, \quad t_p = 10 \text{ ns}, \quad m = 100$$

$$S = \frac{m \cdot t_n}{(K+m-1) \cdot t_p} = \frac{100 \times 50}{(6+100-1) \cdot 10} = 4.76.$$

$$S_{\max} = \frac{t_n}{t_p} = \frac{50}{10} = 5$$

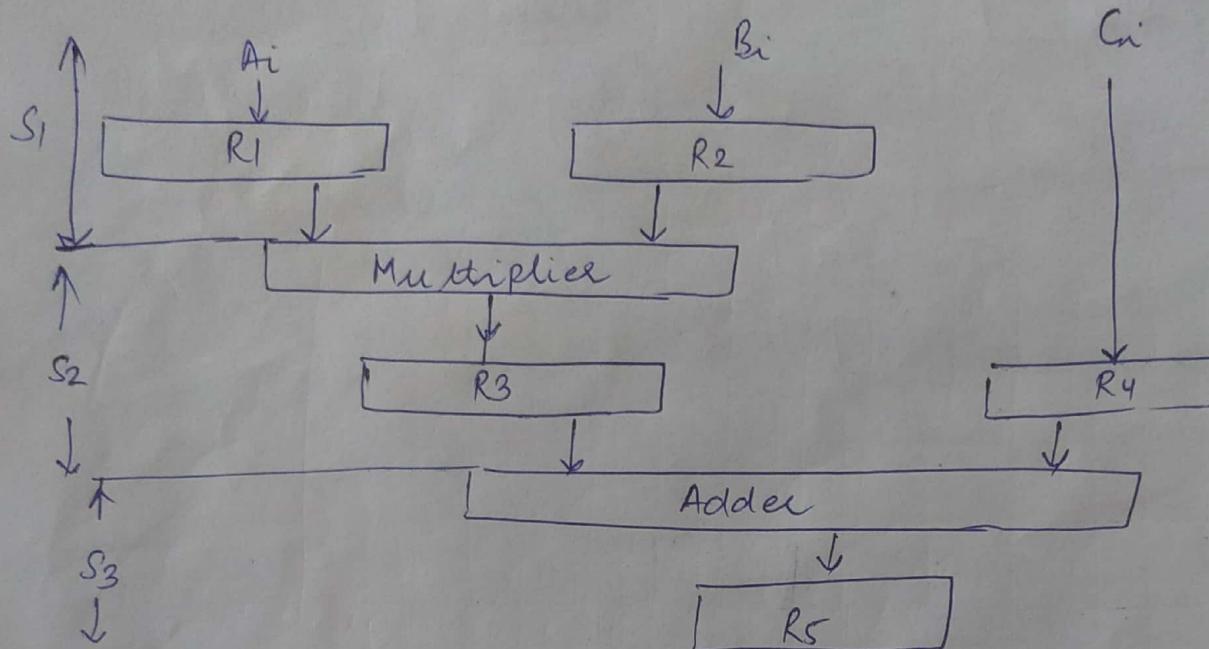
eg of Pipeline.

$$A_i * B_i + C_i, \quad i=1, \dots, 7.$$

$$R_1 \leftarrow A_i, \quad R_2 \leftarrow B_i$$

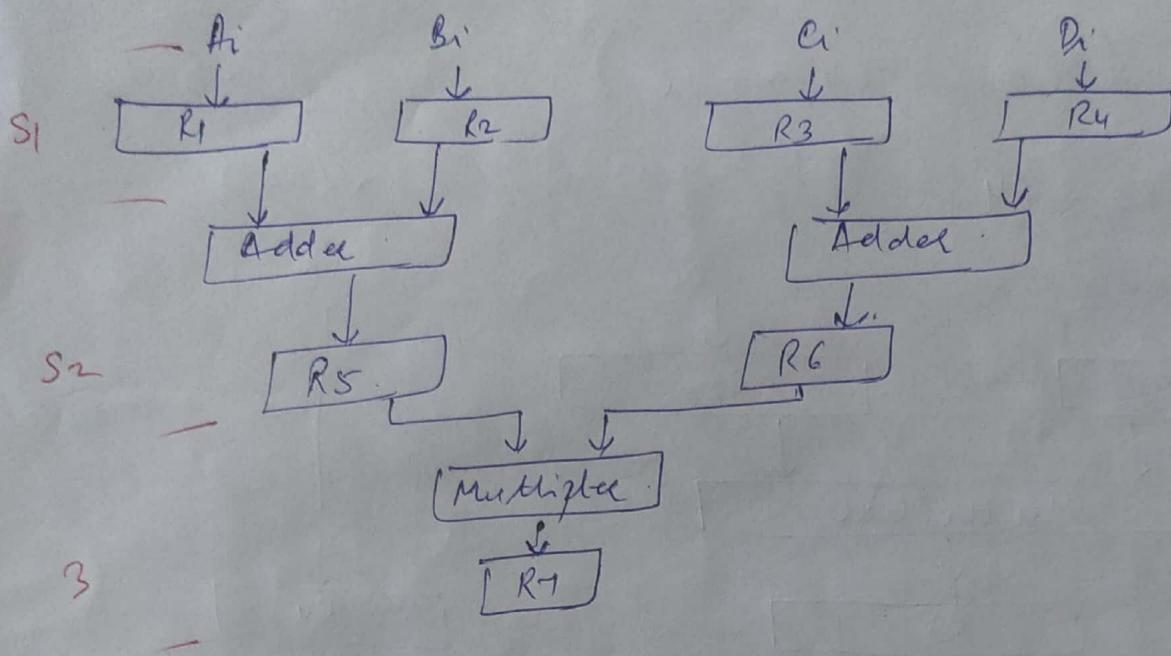
$$R_3 \leftarrow R_1 * R_2, \quad R_4 \leftarrow C_i$$

$$R_5 \leftarrow R_3 + R_4$$



Clock Pulse	Segment 1		Segment 2		Segment 3	
	R1	R2	R3	R4	R5	
1	A <sub>1</sub>	B <sub>1</sub>				
2	A <sub>2</sub>	B <sub>2</sub>	A <sub>1</sub> *B <sub>1</sub>	G <sub>1</sub>		
3	A <sub>3</sub>	B <sub>3</sub>	A <sub>2</sub> *B <sub>2</sub>	C <sub>2</sub>	A <sub>1</sub> *B <sub>1</sub> + C <sub>1</sub>	
4	A <sub>4</sub>	B <sub>4</sub>	A <sub>3</sub> *B <sub>3</sub>	C <sub>3</sub>	A <sub>2</sub> *B <sub>2</sub> + C <sub>2</sub>	
5	A <sub>5</sub>	B <sub>5</sub>	A <sub>4</sub> *B <sub>4</sub>	C <sub>4</sub>	A <sub>3</sub> *B <sub>3</sub> + C <sub>3</sub>	
6	A <sub>6</sub>	B <sub>6</sub>	A <sub>5</sub> *B <sub>5</sub>	C <sub>5</sub>	A <sub>4</sub> *B <sub>4</sub> + C <sub>4</sub>	
7	A <sub>7</sub>	B <sub>7</sub>	A <sub>6</sub> *B <sub>6</sub>	C <sub>6</sub>	A <sub>5</sub> *B <sub>5</sub> + C <sub>5</sub>	
8	A <sub>8</sub>	B <sub>8</sub>	A <sub>7</sub> *B <sub>7</sub>	G <sub>7</sub>	A <sub>6</sub> *B <sub>6</sub> + C <sub>6</sub>	
9						

Q In certain specific computation it is necessary to perform arithmetic operation  $(A_i + B_i) * (C_i + D_i)$ . Specify the pipeline configuration to carry the task.



# Arithmetic Pipeline

9212999888 25

- To speed up the performance.

e.g. Two floating pt. numbers.

Steps to follow

- Compare the exponent.
- Align the mantissa
- Add or subtract
- Normalise the result.

e.g.  $X = 0.9504 \times 10^3$

$Y = 0.8200 \times 10^2$

After Alignment

$X = 0.9504 \times 10^3$

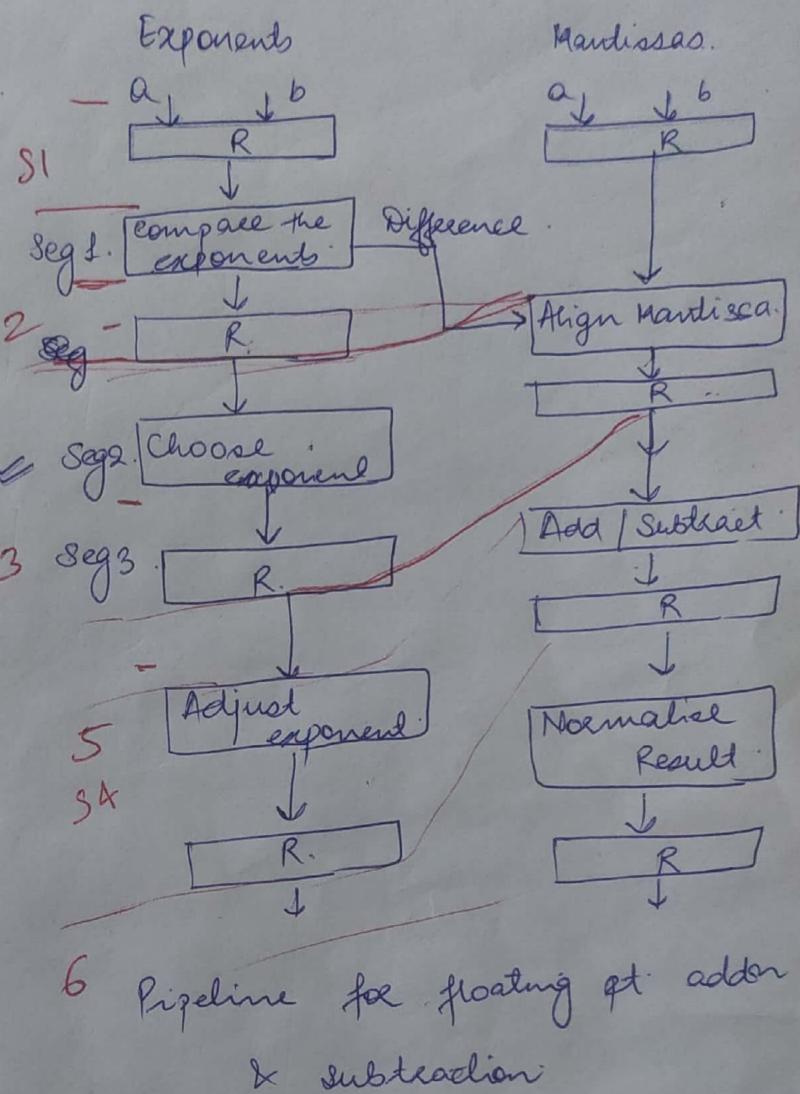
$Y = 0.08200 \times 10^3$

$Z = X + Y$

$= 1.0324 \times 10^3$

Normalise the result.

$= 0.10324 \times 10^4$



6 Pipeline for floating pt. addition & subtraction

# Flynn's Classification:-

27

(Parallel Processing in co)

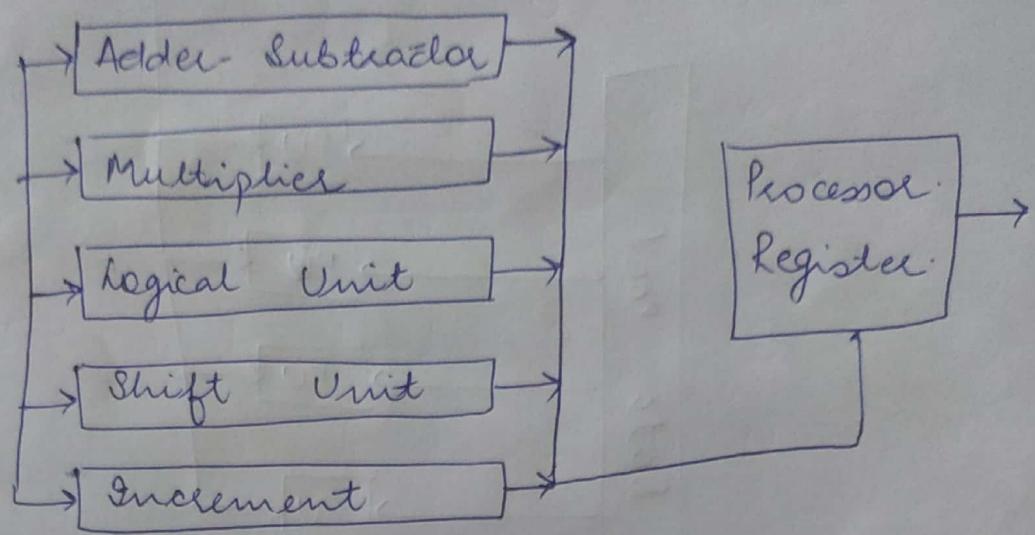


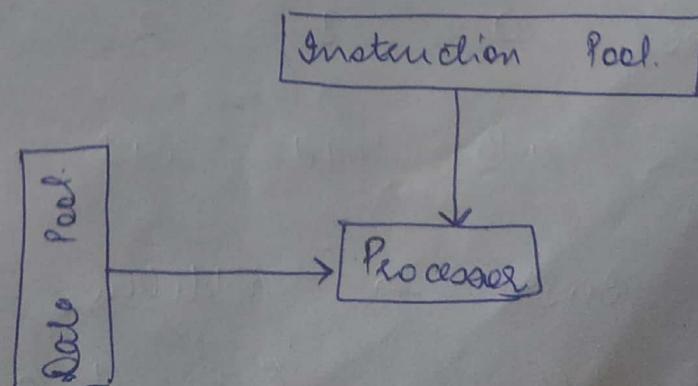
Fig: Processor with multiple Unit (ALU)

- Michael Flynn
- Instruction & Data Stream

## Classification

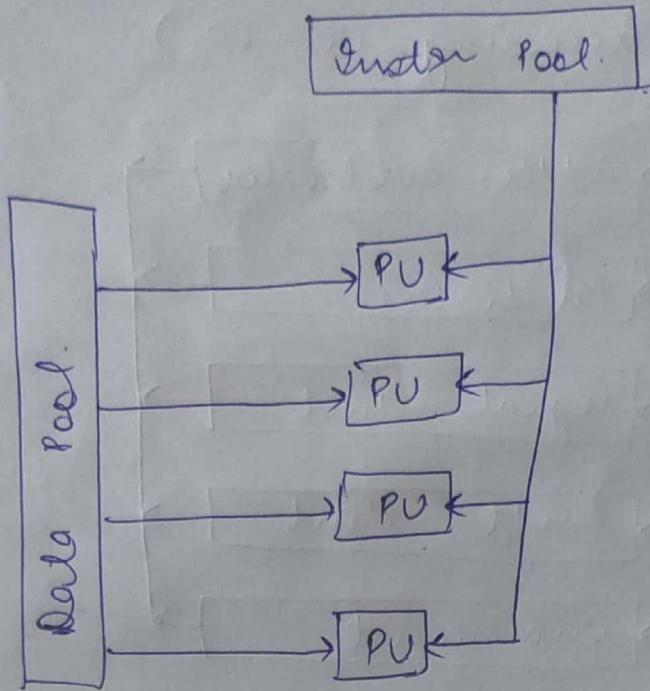
1. S1 SD      Single Instruction      Single Data
2. SI MD      Multiple
3. M1 SD
4. M1 MD

S1 SD

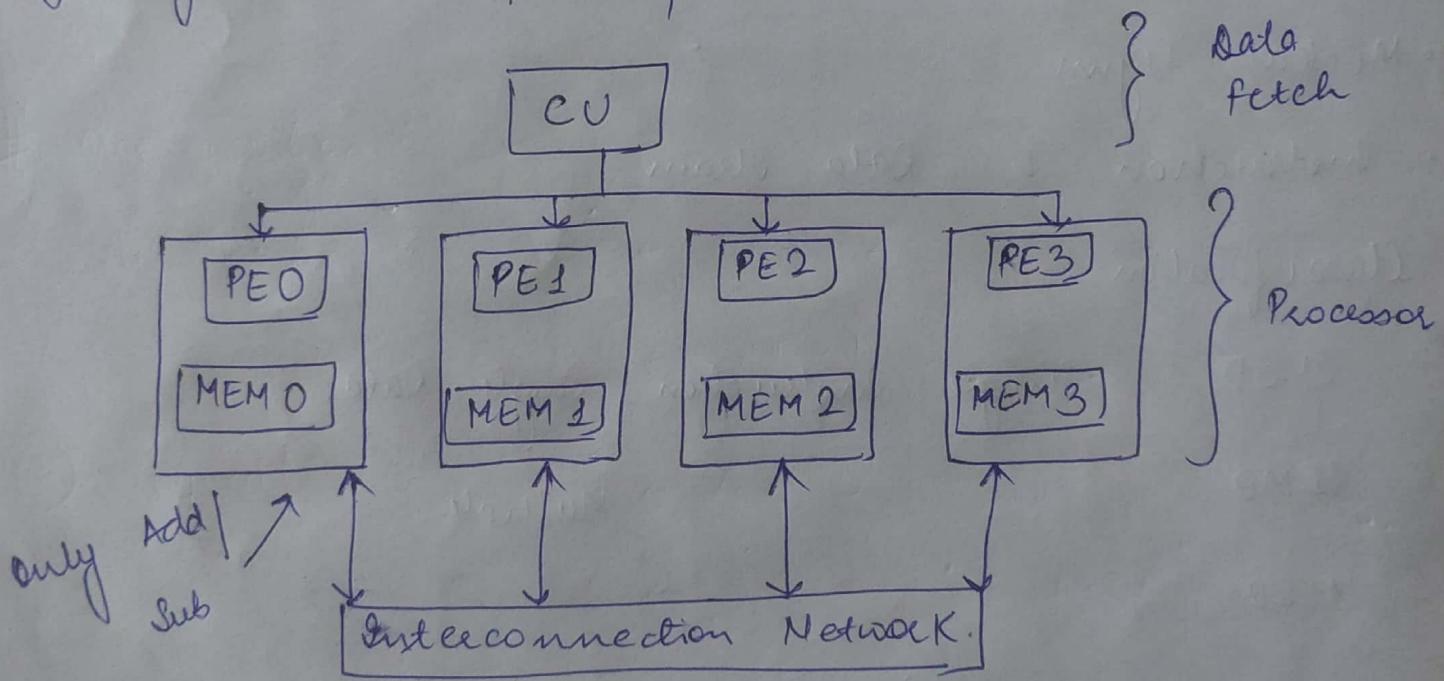


IBM  
1<sup>st</sup> Comp.

## 2. SIMD Architecture



eg Synchronous / Multiple Processor



Processor	0	1	2	3
	a+b	c+d	e+f	g+h.

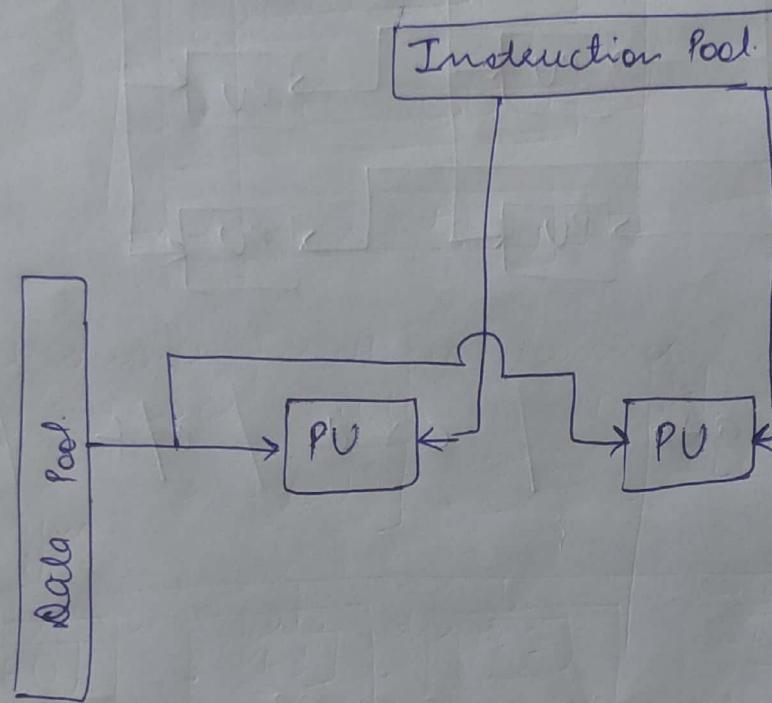
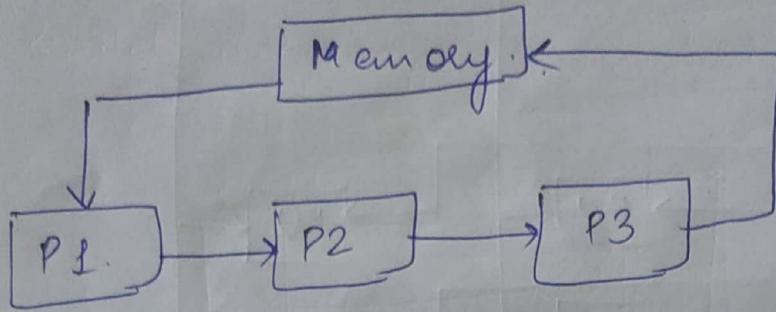
Single Index → Addition

Multiple data → a, b, c, d, e, f, g, h

eg  
ILLIAC-IV  
→ 1952  
→ 1st parallel computer

### 3. MISD Architecture

28



- Rarely Used for Specified Applic.

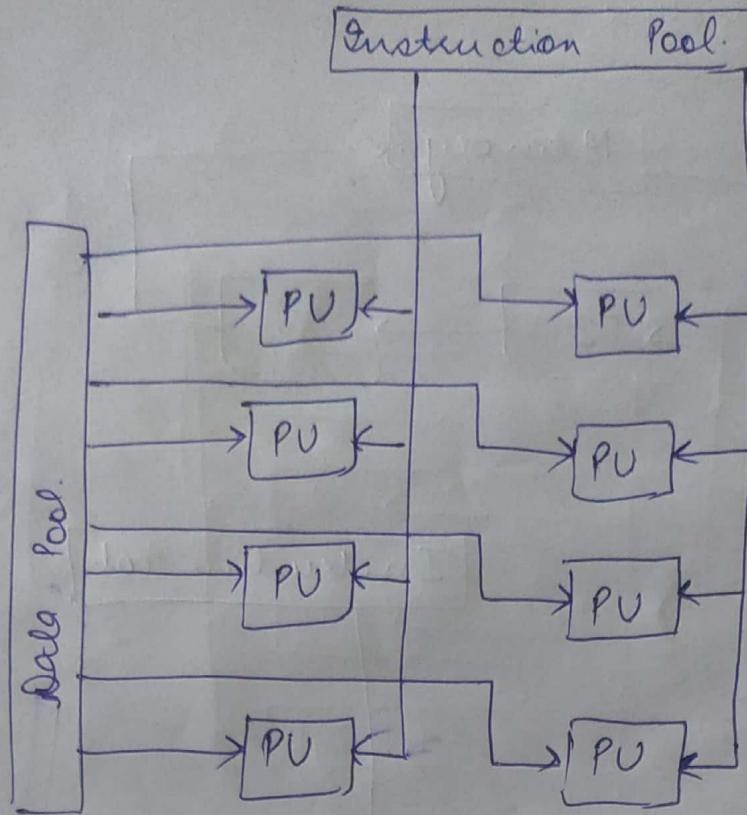
### 4. MIMD

- On Parallel Computer in real sense
- All multiprocessor system fall under this class

eg Burrough's Comp.

→ 1962

→ 60 programs compiled at a time



Processor  $\rightarrow [a+b] / [c \times d] / [e/f] / [m-n]$

## # Micro operation

① Arithmetic  $R_3 \leftarrow R_1 + R_2$

~~$R_1 \leftarrow R_1 - 1$~~

② Logical  $F \leftarrow A \wedge B$  AND.

~~$F \leftarrow A$  Transfer A~~

③ Shift  $R \leftarrow \text{shl } R$  shift left register R

$R \leftarrow \text{she } R$  — right " "

④ Logical Shift  $\text{shl}$  logical shift left  
 $\text{she}$  " " right

## ⑤ Circular Shift

cil      circular shift left  
 cir      "      "      right

### Application

- Manipulate bits or register.

eg Register A holds the 8 bit binary 11011001. Determine the B operation and logical microop<sup>n</sup> to be performed in order to change value in A to 01101101, 11111101.

- Selective Set

→ Bit pattern in B is used to set certain bits in A (Insert)

$$\begin{array}{r}
 1100 \\
 1010 \\
 \hline
 1110
 \end{array}
 \quad A \leftarrow A + B \quad (\text{OR opn})$$

- Selective Complement

→ Bit pattern B is used to complement bits in A

$$\begin{array}{r}
 1100 \\
 1010 \\
 \hline
 0110
 \end{array}
 \quad A \leftarrow A \oplus B \quad (\text{XOR})$$

- Selective Clear

→ Bit pattern in B is used to clear certain bits in A.

$$\begin{array}{r}
 1100 \\
 1010 \\
 \hline
 0100
 \end{array} \quad A \leftarrow A \cdot \bar{B}$$

- Mask operation

→ Bits of A are cleared only if there are corresponding 0's in B.

$$\begin{array}{r}
 1100 \\
 1010 \\
 \hline
 1000
 \end{array} \quad (A \leftarrow A \cdot B) \quad \text{AND opn}$$

Eg A register contain eight bit 0 110 1010. Replace the left most bits by the value 1001.

Soln:

A	0 110 1010	
	0 000 1111	B [AND bcoz clear]
	0 000 1010	

Insert new value

$$\begin{array}{r}
 0\ 0\ 0\ 0\ 1\ 0\ 1\ 0 \\
 + 1\ 0\ 0\ 1\ 0\ 0\ 0 \\
 \hline
 1\ 0\ 0\ 1\ 1\ 0\ 1\ 0
 \end{array}
 \quad A \text{ before} \quad (B \text{ Insect}) \quad A+B.$$

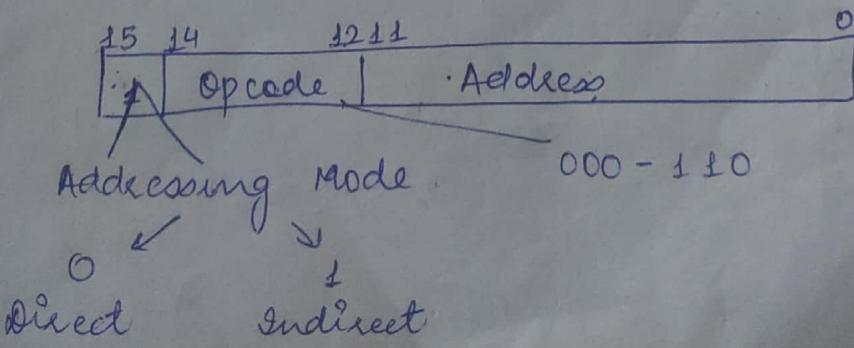
$$\begin{array}{r}
 \text{Solve (i)} \quad 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \quad A \ (\text{Given}) \\
 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \\
 \hline
 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \quad \text{XOR.}
 \end{array}$$

$$\begin{array}{r}
 \text{(ii)} \quad 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \\
 0 \ 0 \ 1 \ 0 \ 0 \ 1 \ 0 \ 0 \\
 \hline
 1 \ 1 \ 1 \ 1 \ 1 \ 1 \ 0 \ 1
 \end{array} \quad \text{xor}$$

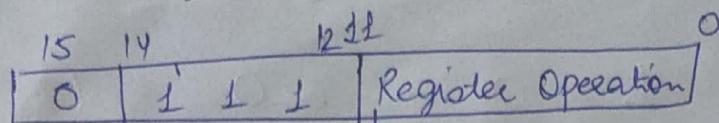
## # Computer Instructions :-

- The basic computer has three instruction code format.

## (a) Memory- Reference Instruction



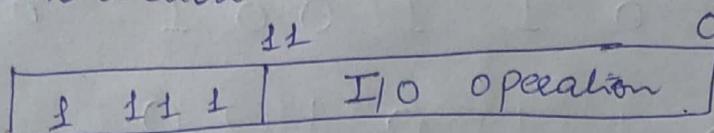
## 2. Register Reference



→ Operation perform on register

Reg. Ref. { 1 1 1      Opcode  
          0      Addressing mode

## 3. I/O Instruction



Hexadecimal Code

Symbol.	D=0	S=1	Description
AND	0 XXX	8 XXX	AND Memory Word to AC
ADD	1 XXX	9 XXX	Add      u      u      u      u
LDA	2 XXX	10 XXX	Load     u      u      u      u
STA	3 XXX	11 XXX	Store content of AC in memory
BUN	4 XXX	12 XXX	Branch Unconditionally
BSA	5 XXX	13 XXX	Branch & Save Return Address
ISZ	6 XXX	14 XXX	Increment & Skip if zero
CLA	7800		Clear Accumulator
CLE	7400		"      E
CMA	7200		Complement AC
CME	7100		"      E
CIR	7080		Circulate Right
CIL	7040		"      Left
INC	7020		Increment AC
SPA	7010		Skip Next instruction if AC is +ve
SNA	7008		"      "      "      "      "      -ve
SZA	7004		"      "      "      "      "      -zero
SZE	7002		E is zero
HLT	7001		Halt Computer

INP	F 800	Get to AC
OUT	F 400	O/P →
SKI	F 20 0	Skip on I/P flag
SKO	F 1 00	u u O/P flag
ION	F 0 80	Interrupt ON
IOP	F 0 40	u OFF

## # Instruction Cycle.

- A program is stored in the memory unit of the computer consist of sequence of instructions.
  - The program is executed in the computer by going through a cycle for each instruction.
  - Each instruction cycle is subdivided into sequence of phases. as following
- 1) Fetch an instruction from memory.  
(Retaining the information)
  - 2) Decode the instruction.  
(Which operation to be performed)
  - 3) Read the effective address from memory.  
(If the instruction has indirect address)
  - 4) Execute the instruction

Upon completion of step 4, control goes back to step 1 to fetch, decode and execute the next instruction.

This process continues until a HALT instruction encountered.

PC - contains the address of next instruction to be executed.

SC - If apply a clock pulse on sequence counter we get time To

Step 1 → Start the program  $SC \leftarrow 0$  (Sequence Counter) } fetch data from memory

Step 2 → Initial time  $T_0$  :  $AR \leftarrow PC$  (Program Counter placed in Address Register)

Step 3 →  $T_1$  :  $IR \leftarrow M[AR]$ ,  $PC \leftarrow PC + 1$   
 (Instruction is read from memory) (IR → Instruction Register)  
 (Whatever the instn, decode it)

Step 4 →  $T_2$  : decode opcode IR (12-14)

$\begin{array}{|c|c|} \hline \text{opcode} & \text{Address} \\ \hline 12-14 & \\ \hline \end{array}$ 
} Decode.

$AR \leftarrow IR(0-11)$

$I \leftarrow IR(15)$

Whatever the data in the memory is decoded in this step

Step 5.  $T_3$  : decision taking cycle  $\rightarrow MR / RR / IO$

→ Indirect exception at  $T_4$   
 direct " " "  $T_3$

Step 6.  $T_4$

Program  $\rightarrow$  Memory Unit



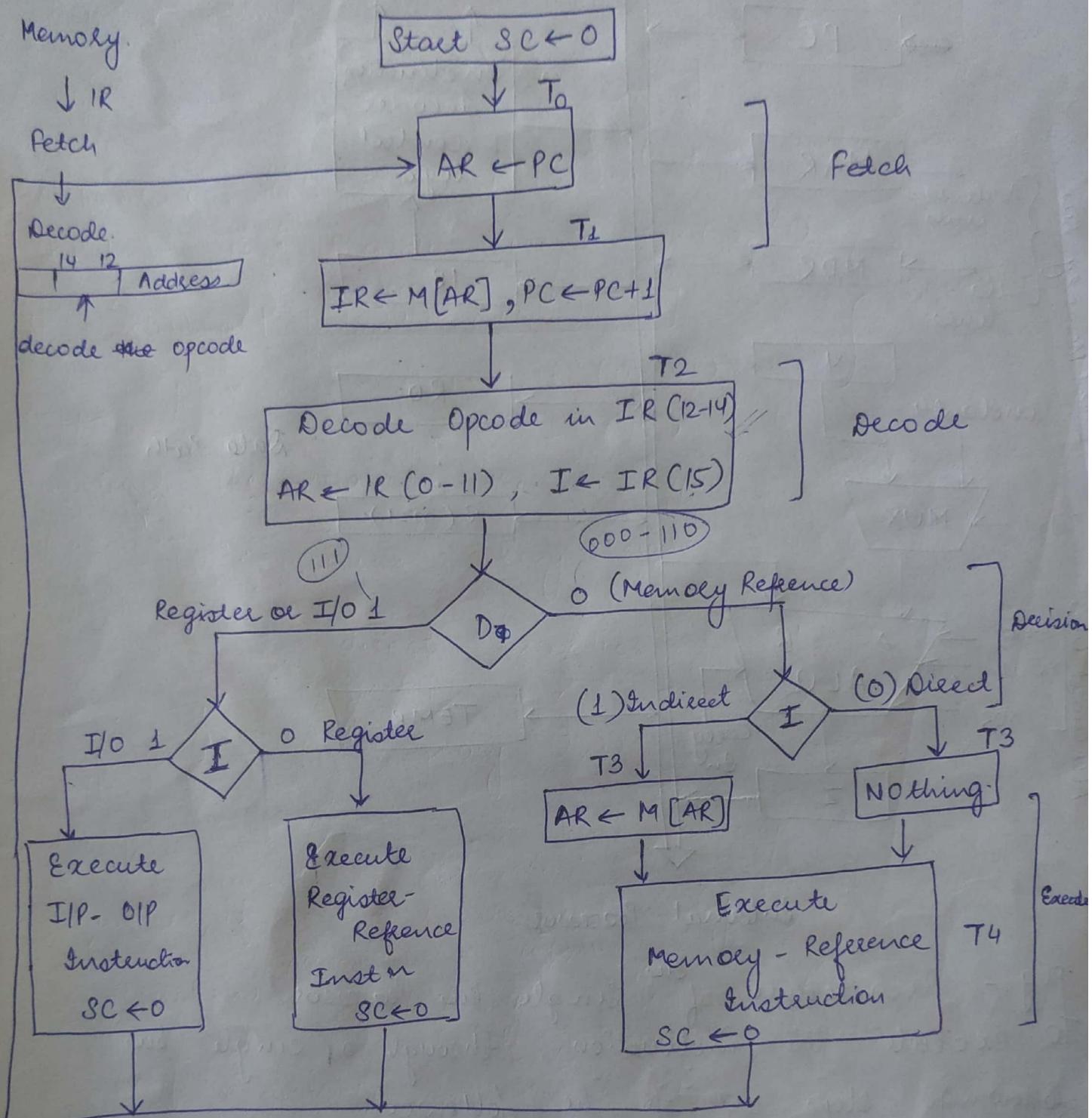
Sequence of Instructions  $\rightarrow$  Instruction Cycle.

$\begin{array}{cccc} \text{Fetch} & \text{Decode} & \text{Decision} & \text{Execute} \\ \swarrow & \downarrow & \searrow & \\ \end{array}$

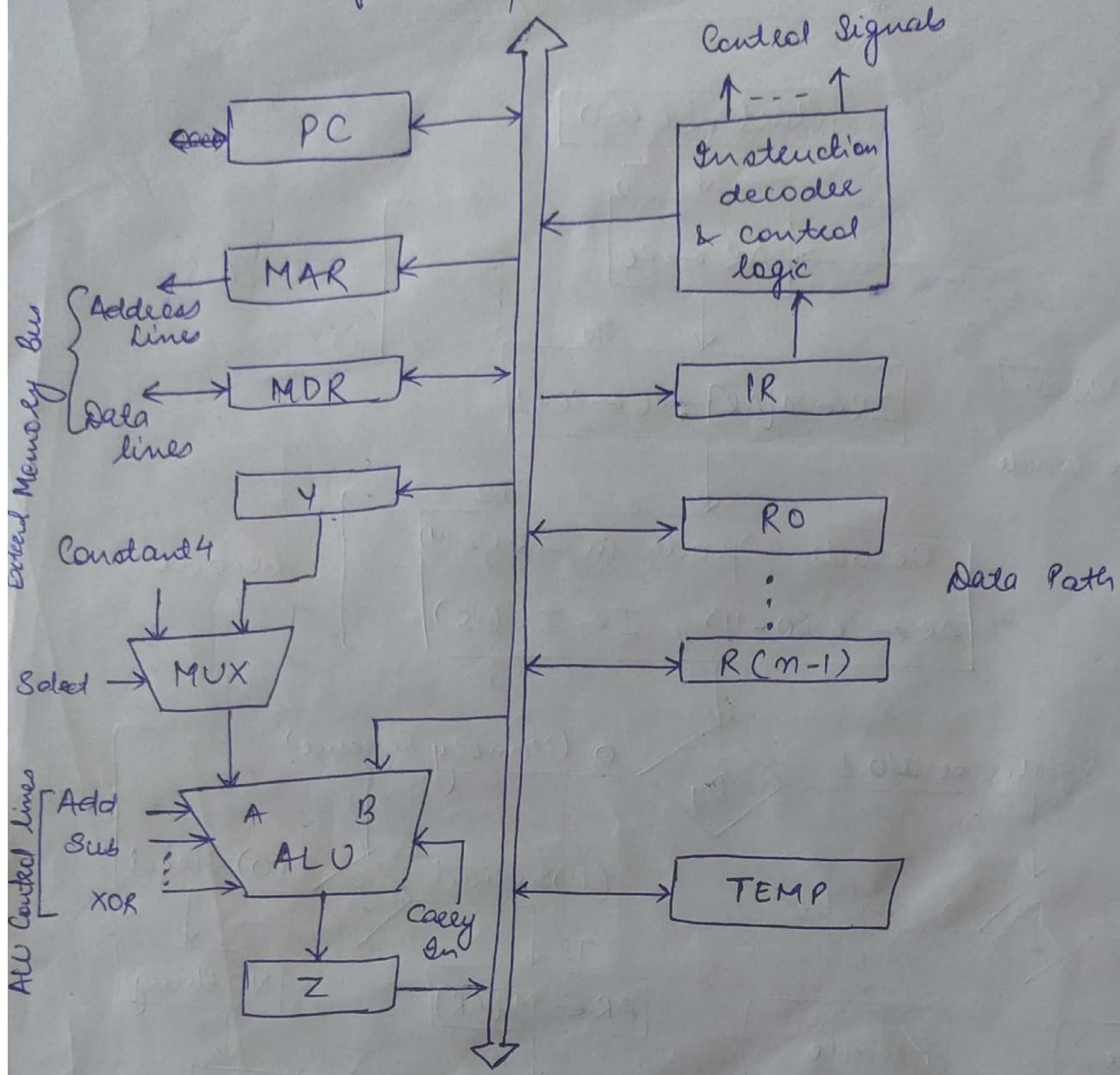
$T_{0.1}$	$T_2$		
Fetch	Decode	EA	Execute Decision

# Flow Chart of Instruction Cycle.

32



# # Execution of Complete Instruction



Internal Processor

Fig:- Architecture of <sup>Bus</sup> Single Bus Processor  
To execute the instruction, through a single bus following steps are to be followed.

- ① Fetch the data

In figure

- ① R<sub>0</sub> to R<sub>n-1</sub> are general purpose register
- ② IR - Instruction / Indexed Register  
It holds the instruction currently being executed
- ③ Instruction decoder - tells about the instruction to be performed  
for fetching :- (i) AR ← PC
- ④ Y, Z & TEMP are temporary registers used by processor during the execution of some instruction
- ⑤ PC : contains the address of next instruction to be executed.
- ⑥ MAR - Memory Address Register
- ⑦ MDR - All the data stored in Memory Data Register

eg ADD (R3), R1  
 ↑      ↑  
 Indirect AM    Direct AM

R3 → 1008 ,

R1 → \$0

# Memory

Address	Content
1000	Add(R3), R1
1008	20

PC 1000

R3 1008.

R1 20

IR Add(R3), R1

- 1) Fetch the instruction from memory
- 2) Fetch the first operand from memory.
- 3) Perform addition operation.
- 4) Load the result into R1.

## Control Sequence

- 1) PC<sub>out</sub>, MAR<sub>in</sub>, Read, Select 4, Add, Z<sub>in</sub>
  - 2) Z<sub>out</sub>, PC<sub>in</sub>, WMFC
  - (Wait for memory function to complete)
  - 3) MDR<sub>out</sub>, IR<sub>in</sub>
- ] Fetching instruction in IR

- 4) R3<sub>out</sub>, MAR<sub>in</sub>, Read
  - 5) R1<sub>out</sub>, A<sub>in</sub>, WMFC
  - 6) MDR<sub>out</sub>, Select 4, Add, Z<sub>in</sub>
  - 7) Z<sub>out</sub>, R<sub>in</sub>, End.
- Mux → 4  
Select  
A = 4  
B = 1000  
Z = 1004  
PC → 1004 Inst<sup>rn</sup>  
MDR<sub>out</sub> → ADD(R3), R1

### Sequence of operation

Add content of register R1 to those of R2 and store result in R3.

1.  $R_1$  out,  $Y$  in

2.  $R_2$  out, Select  $Y$ , Add,  $Z$  in

3.  $Z$  out,  $R_3$  in

fetching a word from memory.

e.g  $MOV(R1), R2$

$\rightarrow MAR \leftarrow [R1]$

$\rightarrow$  Start Read operation on memory bus.

$\rightarrow WMPC$

$\rightarrow$  Load MDR from memory bus

$\rightarrow R2 \leftarrow [MDR]$

Control Sequence

1.  $R_1$  out,  $MAR$  in, Read

2.  $MDR$  in,  $WMPC$

3.  $MDR$  out,  $R2$  in

~~eg~~ Add the immediate number (NUM) to register(R1)

CS

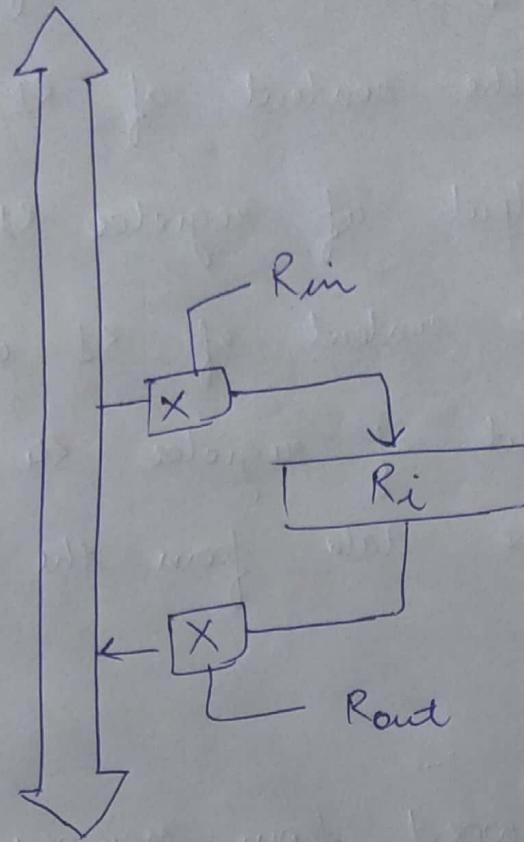
- ① PC<sub>out</sub>, MA<sub>in</sub>, Read, Select 4, Add, Z<sub>in</sub>
- ② Z<sub>out</sub>, PC<sub>in</sub>, WMFC
- ③ MDR<sub>out</sub>, IR<sub>in</sub>
- ④ IR NUM<sub>out</sub>, Y<sub>in</sub>
- ⑤ R<sub>1</sub> out, Select 4, Add, Z<sub>in</sub>
- ⑥ ~~Zout~~ R<sub>1</sub> in, End

~~eg~~ Add the content of memory location NUM to register R<sub>1</sub>.

CS

[Ans]

# Register Transfer



- I/P & O/P gates for register  $R_i$  are controlled by signals  $R_{in}$  &  $R_{out}$ .
- If  $R_{in}$  is set to 1, data available on common bus are loaded into  $R_i$ .
- If  $R_{out}$  is set to 1, content of register are placed on bus.
- If  $R_{out}$  is set to 0, bus can be used for transferring data from other registers.

of Data transfer b/w two registers.

→ Transfer the content of R1 to R4.

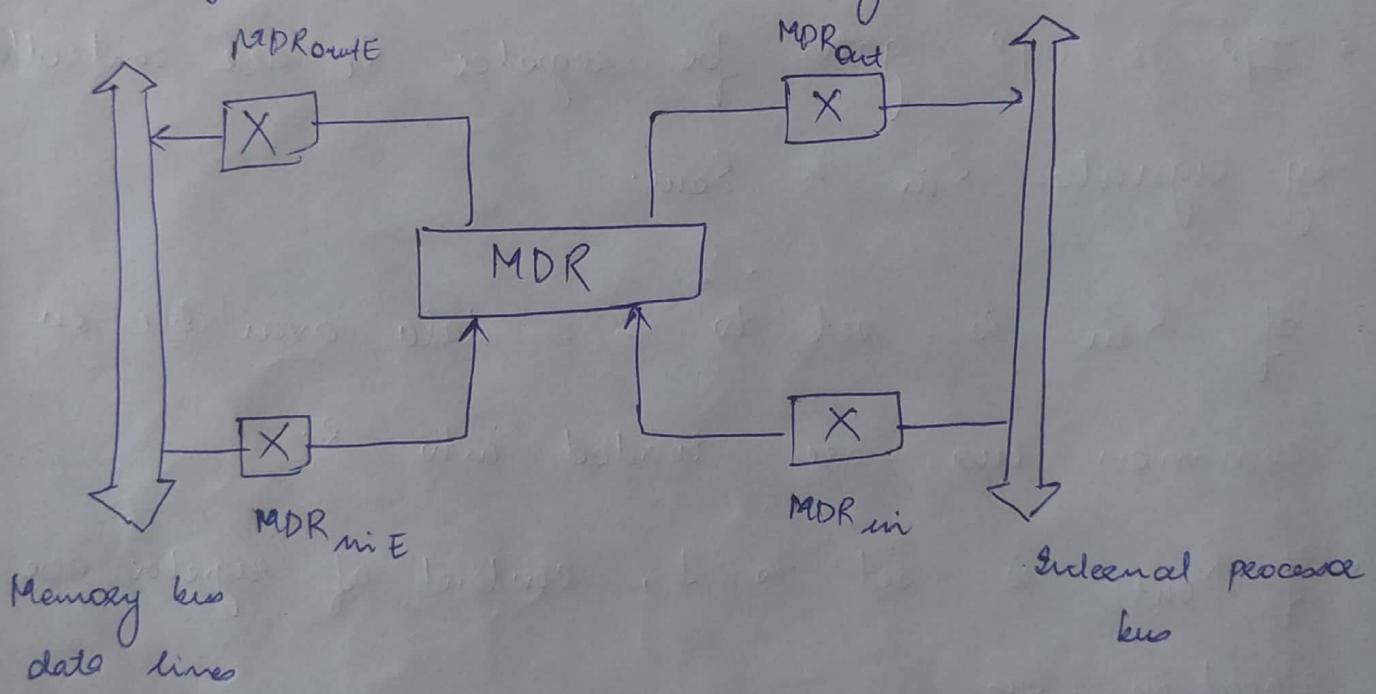
1) Enable output of register R1 by setting  $R_{1\text{out}} = 1$ .

This place the content of R1 on the processor bus.

2) Enable input of register R4 by setting  $R_{4\text{in}} = 1$ .

This loads the data from the processor bus into register R4.

# fetching a word from memory.



## # Memory Hierarchy

- The memory unit is an essential component in any digital computer since it is needed for storing programs and data.
- Not all the accumulated information is needed by the processor at the same time. Therefore, it is more economical to use low-cost storage devices to serve as backup for storing the information that is not currently used by the CPU.
- The main memory unit the communicates directly with CPU is called main memory. Only programs and data currently needed by processor reside in main memory.
- Devices that provide back up storage are called auxiliary storage memory.

→ e.g Magnetic disk / hard disk (TB)

Magnetic Tape / Tape drives (PB / EB)

→ AM are used for storing programs, large data files, and other back up information.

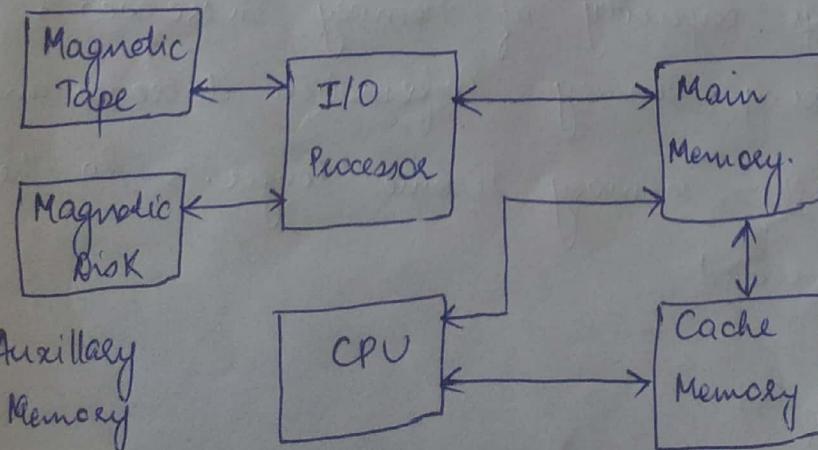


Fig:- Memory Hierarchy in computer

The total memory capacity of computer can be visualized as hierarchy of components

- Magnetic Tapes - used to store removable files.
- Magnetic Disk - used as backup storage.
- Main memory occupies a central position to communicate directly with CPU and auxiliary memory.

When programs in main memory are not needed by the CPU, they are transferred to auxiliary memory, to provide space for recently used program and data.

### Cache Memory-

A special very-high memory called a cache is used to increase the speed of processing by making current programs and data available at the CPU at a rapid rate.

The cache is used for storing segments of programs currently being executed in the CPU and temporary data frequently needed in the present calculations.

Note:- As the storage capacity of memory increases, the cost per bit for storing binary information decreases and the access time of the memory become longer. So the memory hierarchy is needed.

Auxiliary Memory - large storage capacity, relatively inexpensive 37  
low access speed compared to Main Memory

Cache Memory - small, relatively expensive,  
- very high access speed. 24  
53

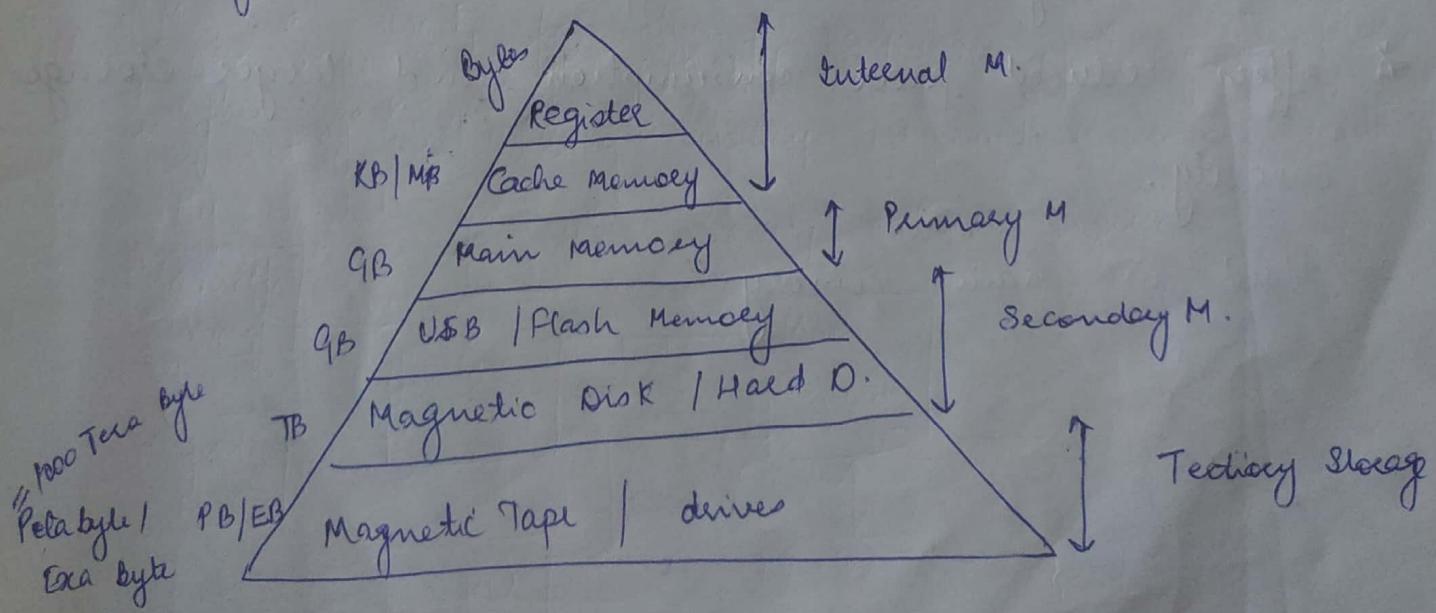
Therefore as the memory access speed increases, its cost increases.

## # Multiprogramming

Many operating system are designed to enable the CPU to process number of independent programs concurrently. This is known as multiprogramming.

## # Memory Management System

The part of computer system that supervise the flow of information b/w auxiliary memory & main memory is known as memory management system.



1EB = 1000 Peta Byte

1PB = 1000 Terra Byte

## # Main Memory (RAM)

- Central storage in computer system.
- Relatively large and fast memory.
- Used to store programs and data during computer operation.
- Integrated circuit RAM chips are available in two possible operating modes, static and dynamic.

→ Static RAM -

- \* Consist of internal flip-flop that store the binary information.
- \* Stored information remain valid as long as power is applied to the unit.

→ \* Used in cache memory.  
→ Dynamic RAM -

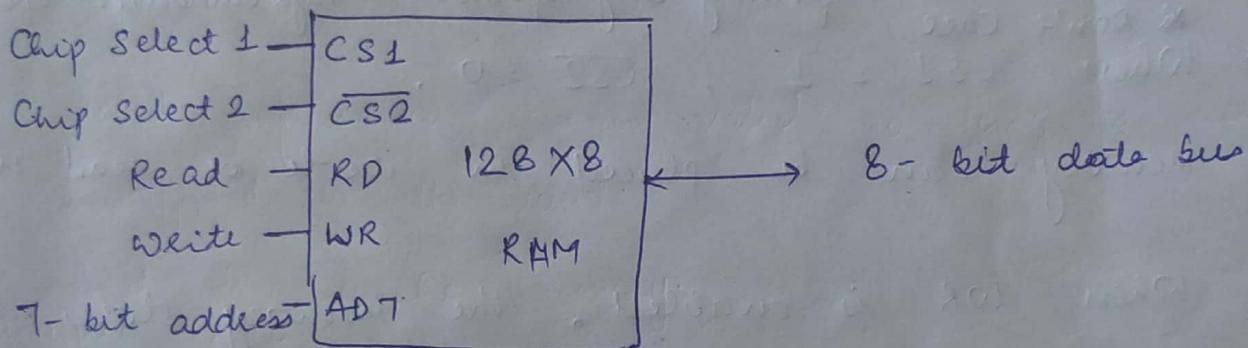
- \* Stores the binary information in the form of electric charges that are applied to capacitors.
- \* offers reduced power consumption and larger storage capacity.
- \* Used in main memory.

## RAM & ROM Chips.

\* Read - data transfer b/w main memory-to CPU. 38

\* Write - data transfer b/w CPU-to main memory.

- RAM chip is better suited for communication with the CPU if it has one or more inputs that select the chip only when needed.
- RAM has bidirectional data bus that allows the data transfer either from memory-to CPU i.e. read operation or CPU-to memory i.e. write operation.
- Block diagram of RAM chip



- The capacity of the memory is 128 words of eight bits (one byte) per word.
- This requires 7-bit address and 8-bit bidirectional data bus.
- The read and write input specify the memory operation and the two chip select (CS) control input are used for enabling the chip only when selected by processor.

\* 2 Chip Select →  $2^2 = 4$  RAM can be selected

## • Functional Table of RAM

CS1	$\overline{CS2}$	RD	WR	Function	State of Bus
0	0	x	x	Inhibit	High Impedance
0	1	x	x	Inhibit	"
1	0	0	0	Inhibit	"
1	0	0	1	Write	Input data to RAM
1	0	1	0	Read	Output data from RAM
1	1	x	x	Inhibit	High Impedance

\* Inhibit - Not performing any operation      \* High Impedance - Open Circuit

\* Don't Care

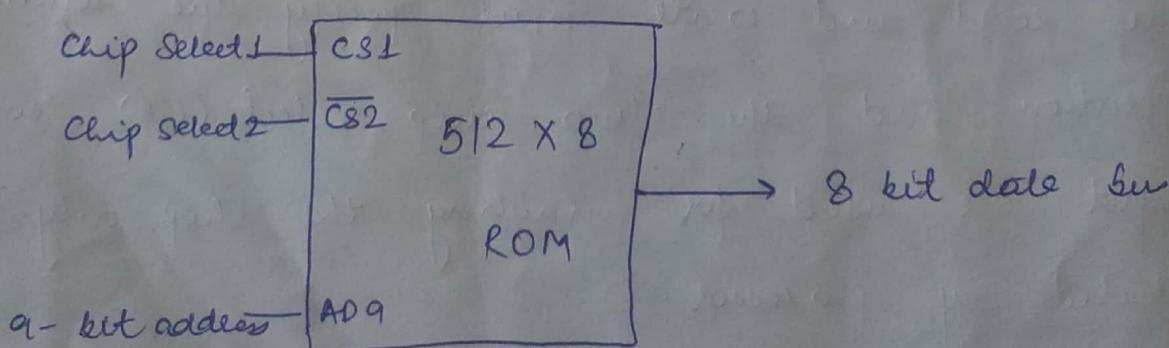
→ When  $CS1 = 1$ ,  $\overline{CS2} = 0$ ,

Memory can be placed in write or read mode.

→ When WR is enabled, the memory stores a byte from data into location

→ When RD is enabled, the content of selected bus is placed into the data bus.

## • Block Diagram of ROM



→ Since ROM can only read, data bus can only be in output mode.

Note:- For the same-size chip, it is possible to have more bits of ROM than of RAM, because the internal binary cells of ROM occupy less space than RAM. For this reason, ROM specifies 512 bytes, while RAM has only 128 bytes.

→ 9 address lines specifies 512 bytes.

→  $CS1 = 1$  and  $CS2 = 0$  for operation

Otherwise data bus is in high impedance state.

### Memory Address Map.

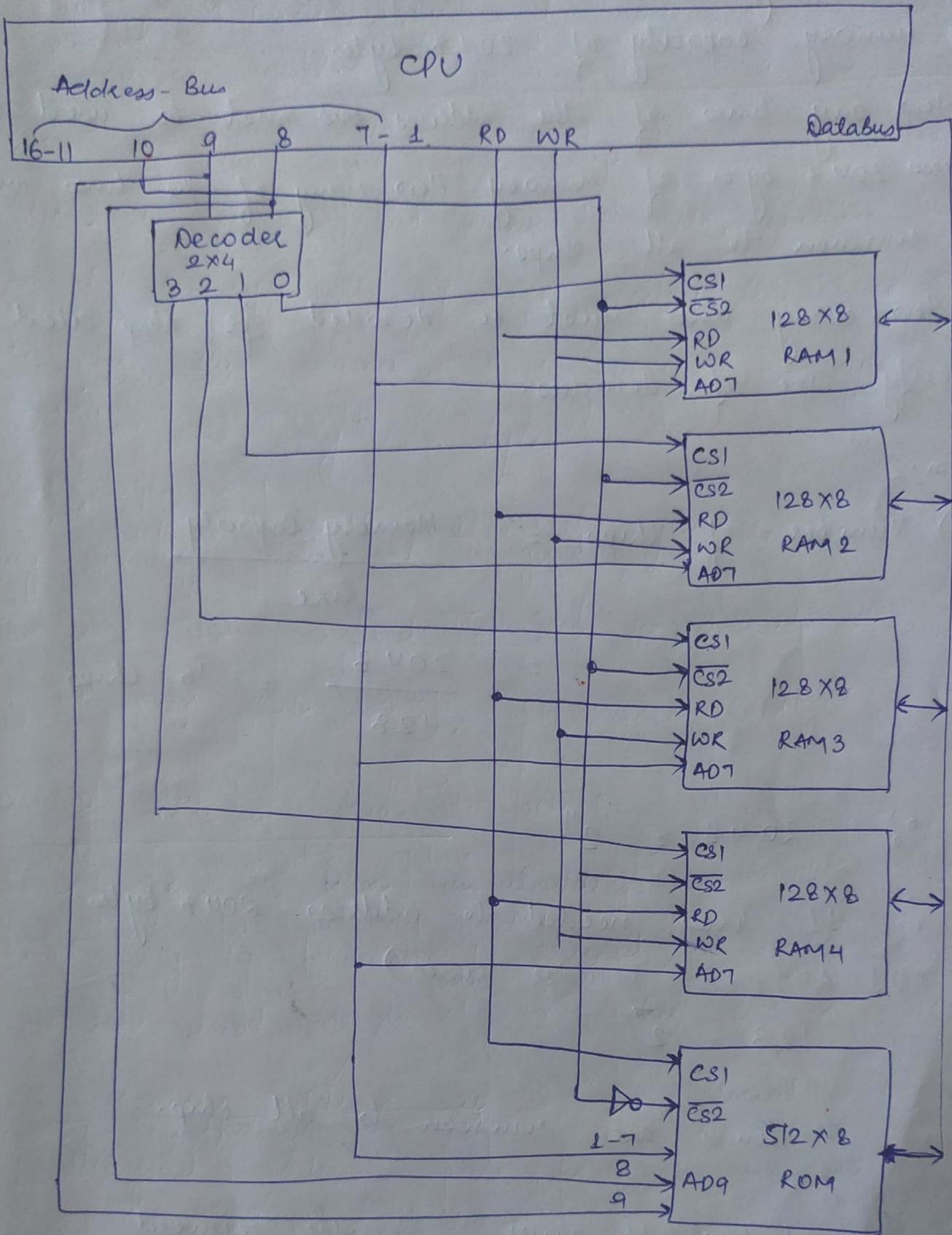
- The designer of the computer system must calculate the amount of memory required for particular application and assign it to either RAM or ROM.
- The Memory address map is a pictorial representation of assigned address space for each chip in the system
- eg 512 bytes RAM & 512 bytes ROM.

Component	Hexadecimal Address	Address Bus									
		10	9	8	7	6	5	4	3	2	1
RAM 1	0000-007F	0	0	0	X	X	X	X	X	X	X
RAM 2	0080-00FF	0	0	1	X	X	X	X	X	X	X
RAM 3	0100-017F	0	1	0	X	X	X	X	X	X	X
RAM 4	0180-01FF	0	1	1	X	X	X	X	X	X	X
ROM	0200-03FF	1	X	X	X	X	X	X	X	X	X

- The X's are assigned to low-order bus lines.
- When line 10 = 0, CPU selects RAM
- When line 10 = 1, it selects ROM.
- There are 16 lines in address bus, the table shows only 10 because in this example other 6 are not used.

### Memory Connection to CPU

Note:- The example shows gives an indication of interconnection complexity that can exist between Memory chips and CPU.



• 2x4 Decoder

Input                  Output

0 0	RAM1
0 1	RAM2
1 0	RAM3
1 1	RAM4

- Q) a) How many  $128 \times 8$  RAM chips are needed to provide a memory capacity of 2048 bytes.
- b) How many lines of the address bus must be used to access 2048 bytes of memory. How many of the lines will be common to all chips.
- c) How many lines must be decoded for chip select? Specify size of decoders.

Soln:-

$$\text{a) Number of chips} = \frac{\text{Memory Capacity}}{\text{Size}}$$

$$= \frac{2048}{128} = 16 \text{ chips}$$

$$\text{b) } 2048 = 2^n$$

$\therefore$  11 lines needed to address 2048 bytes.

(AD1 - AD7 + Decoder lines)

$$128 = 2^7$$

$\therefore$  7 lines are common to all chip.

c) For selecting 16 chips, decoder required is

$4 \times 16$  decoder

Q for 4096 bytes of RAM and 4096 bytes of ROM.<sup>41</sup>  
list the memory address map and indicate decoder size.

Solar:-

$$\text{Number of chips} = \frac{4096}{128}$$

= 32 RAM chips

$$= \frac{4096}{512}$$

= 8 ROM chips.

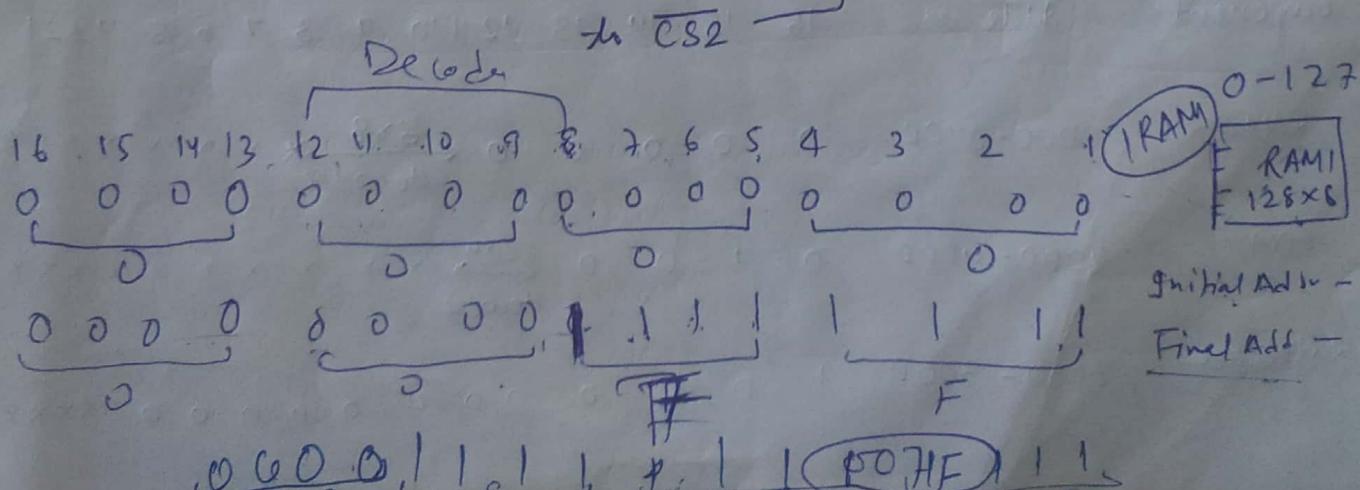
$$\text{Also, } 4096 = 2^{12}$$

12 common address lines & 1 to select  
between RAM & ROM.

Component Address 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1

RAM      0000-0FFF      0 0 0 0 ←  $\frac{5 \times 32}{\text{decoder}}$  → x x x x x x x x

ROM 1000-1FFF 0 0 0 1  $\xleftarrow[decoder]{3 \times 8}$  x x x x x x x x x



A computer employs RAM chips of  $256 \times 8$  and ROM chips of  $1024 \times 8$ . The system needs 2K bytes of RAM, 4K bytes of ROM and four interface units each with four registers. A memory mapped I/O configuration used. The two highest order bits of the address bus are assigned 00 for RAM, 01 for ROM and 10 for interface registers.

- How many RAM & ROM chips are needed.
- Draw a memory address map for the system.
- Give the address range in hexadecimal for RAM, ROM & Interface.

Soln:- a) RAM chips =  $\frac{2048}{256} = 8$  chips

ROM chips =  $\frac{4096}{1024} = 4$  chips

b)  $2048 = 2^{11}$ ,  $256 = 2^8$

$4096 = 2^{12}$ ,  $1024 = 2^{10}$

Interface  $4 \times 4 = 16$  Registers =  $2^4$

c) Component Address      16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1

RAM	0000-0FFF	0 0 0 0 0 0	$\xleftarrow[decode]{3 \times 8}$	x x x x x x x x
-----	-----------	-------------	-----------------------------------	-----------------

ROM	4000-4FFF	0 1 0 0 0 0	$\xleftarrow[decode]{2 \times 4}$	x x x x x x x x
-----	-----------	-------------	-----------------------------------	-----------------

Interface	8000-800F	1 0 0 0 0 0 0 0 0 0 0 0 0 0 0		xxxx
-----------	-----------	-------------------------------	--	------

## # Cache Memory.

Locality of Reference - Analysis of large number of typical programs has shown that the references to memory at any given interval of time tend to be confined with in a few localized areas in memory. This phenomena is referred as locality of reference.

Cache Memory - If the active portion of the program and data are placed in fast small memory, the average memory access time can be reduced, thus reducing the total execution time of the program. Such a fast small memory is referred as cache memory.

Hit Ratio - The performance of cache memory is frequently measured in terms of quantity called hit ratio.

- When the CPU refers to memory and finds the word in cache, it is said to produce a hit.
- If the word is not found in cache, it is in main memory and it counts as a miss.

The ratio of the number of hits divided by the total CPU references to memory is the hit ratio.

## Types of Mapping

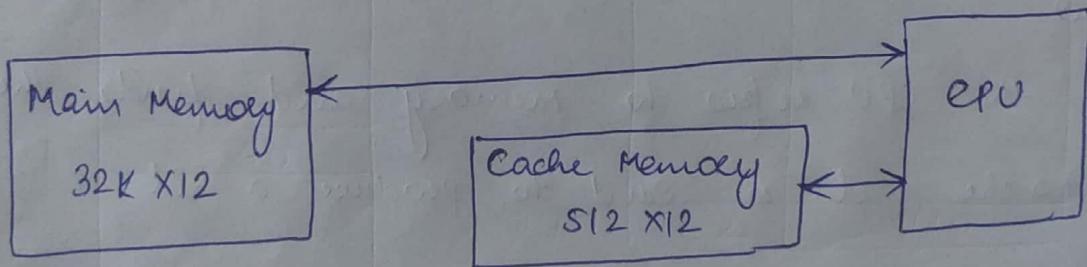
- The transformation of data from main memory to cache memory is referred as mapping process.
- There types of mapping.
  - 1) Associative Mapping
  - 2) Direct mapping
  - 3) Set - Associative Mapping

### Direct Mapping :-

To understand mapping procedure, consider an example

- The main memory can store 32 K words of 12 bits each.
- The cache is capable of storing 512 of these words at any given time.

$$\frac{32}{2^5} \text{ K} \times 2^{10} \rightarrow 2^{15}$$



- For every word stored in cache, there is a duplicate copy in main memory.
  - The CPU communicates with both memories. It sends 15 bit address to cache.
- Hit - CPU accepts data from cache.
- Miss - CPU read the word from main memory.

# I) Associative Mapping

CPU Address (15 bits)



Argument Register

Address	Data
01000	3450
02777	6710
22345	1234

$$\begin{aligned} 5 \times 3 &= 15 \rightarrow \text{Address} \\ 4 \times 3 &= 12 \rightarrow \text{bit data} \end{aligned}$$

- It stores both address and <sup>content of</sup> memory word.
- The address value of 15 bit is shown as five-digit octal number and its corresponding 12 bit word as four digit octal number.
- A CPU address of 15 bit is placed in the argument register and associate memory is search for matching process.

If the address is found, the 12 bit data is read and sent to CPU.

If no match occur, the main memory is accessed for the word.

The address data pair is then transferred to the associative cache memory.

## 2) Direct Mapping

- Associative mapping is expensive compared to RAM because of added logic associated with each cell.
- The CPU address of 15 bits is divided into two fields. The nine least significant bits constitute the index field and the remaining six bits form tag field.

Tag	Index
6 bits	9 bits

- The internal organization of the words in cache memory is as shown

Memory Address	Memory Data
00000	1220
01777	4560
02777	6710

Index Address	Tag	Data
000	00	1220
777	02	6710

b) Cache Memory

a) Main Memory

Fig :- Direct Mapping Cache Organization

- Each word in cache consist of data word and its 44 associated tag.
- When a new word is associated with cache, the tag bits are stored along the data bits.
- When the CPU generates memory request, the index field is used for the address to access the cache.
- The tag field of CPU address is compared with the tag in the word read from cache. If there is no match, there is a miss and required word is read from main memory.
- Disadvantage- Hit ratio drops, if two or more address have same index but different tags are accessed repeatedly.

### 3) Set Associative Mapping

It is an improvement over direct mapping organization in that each word of cache can store two or more words of memory under the same index address.

Each data word is stored together with its tag and the number of tag data items in one word of cache is said to form a set.

Index	Tag	Data	Tag	Data
000	01	3450	02	5670
777	02	6710	00	2340

Fig :- Two way Set - Associative mapping cache

- Each index address refers to two data words and their associated tags.
- Each tag require six bits and each data word has 12 bits, so the
 
$$\text{word length} = 2(6+12)$$

$$= 36 \text{ bits}$$
- Thus, An index address can accomodate 512 words.  
Thus size of cache memory is  $512 \times 36$
- It can accomodate 1024 words of main memory  
Since each word of cache contains two data words.
- In fig, the word stored at addresses 01000, & 02000 of Main Memory are stored in cache memory at index 000.

- When the CPU generates a memory request, the index value of the address is used to access the cache.
- The tag field of the CPU address is then compared with both tags in the cache to determine if a match occurs.
- The hit ratio will improve as the set size increases because more words with the same index but different tags can reside in a cache.

Note :- When a miss occurs in set associative cache, it is necessary to replace one of the tag data items with a new value.

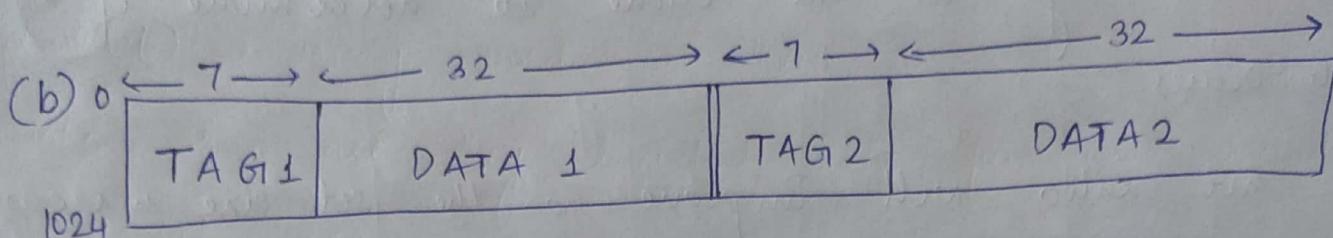
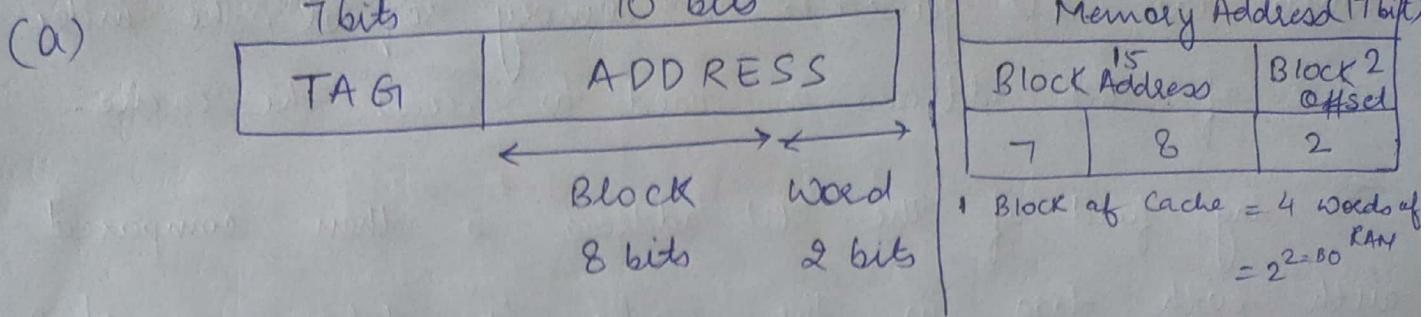
Q) A two way associate cache memory uses block of four words. The cache can accommodate a total of 2048 words from main memory. The main memory size is 128 K x 32. Formulate all pertinent information required to construct the cache memory. What is the size of cache memory.

Soln:-  $\therefore$  Size of Main Memory = 128 K

$$\therefore 128 \text{ K} = 2^{17}$$

2048 from Main Memory

i.e.  $2048 / 2 = 1024$  words of cache.



$$\begin{aligned}
 \text{Size of Cache Memory} &= 1024 \times 2 (7 + 32) \\
 &= 1024 \times 78
 \end{aligned}$$

Q) The access time of a cache memory is 100 ns and that of main memory is 1000 ns. It is estimated that 80% of the memory request are for read and 20% for write. The hit ratio for read access only is 0.9. A write through procedure is used.

- a) What is the average access time of the system considering only memory read cycles.
- b) What is the average access time of the system for both read & write requests.
- c) What is the hit ratio considering the write request also

Soln:-

$$\text{(a) Average Access time for Cache} = 0.9 \times 100 \\ \text{to read} \\ = 90 \text{ ms}$$

$$\text{Average Access time for cache & memory to write} = 0.1 \times (100 + 1000) \\ = 110 \text{ ms}$$

$$\text{Total time taken} = 90 + 110 \\ = 200 \text{ ms}$$

$$\text{(b) Average access time for read} = 0.8 \times 200 \\ = 160$$

$$\text{Average access time for write} = 0.2 \times 1000 \\ = 200$$

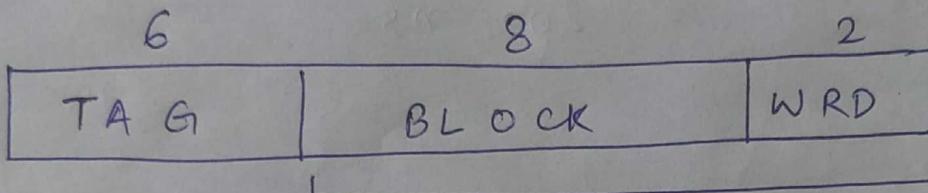
$$\text{Total time taken} = 160 + 200 \\ = 360 \text{ ms}$$

$$\text{(c) Hit ratio} = 0.8 \times 0.9 \\ = 0.72$$

Q A computer has memory unit of ~~64~~ 64 K X 16 and cache memory of 1K words. The cache uses direct mapping with block size of four words  
 (a) How many bits are there in the tag, index, block and words fields of the address format.  
 (b) How many words are there in each word of cache, how are they divided into functions. Include a valid bit.  
 (c) How many blocks can cache accommodate.

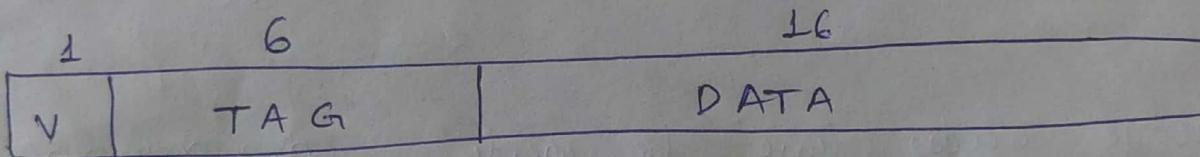
Soln:

(a)



Index = 10 Bit Cache Address

(b)



23 Bits in each word of cache

(c)

$$2^8 = 256 \text{ block of 4 words each.}$$

$$1K \cong \frac{1024}{4} = 256 \text{ Block}$$

## Virtual Memory.

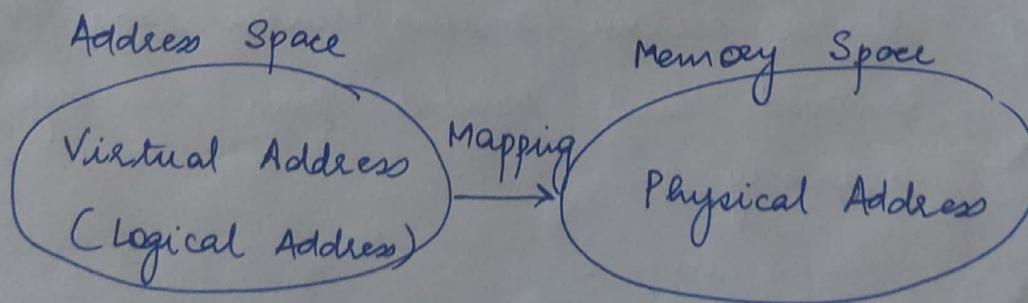
- It is a concept used in some larger computer systems that permit the user to construct programs as though a large memory space were available.
- A virtual memory system provides a mechanism for translating programs generated addresses into correct main memory locations.

→ Address Space

An address used by programmer will be called virtual address, and the set of such addresses the address space.

→ Memory Space

An address in main memory is called location or physical address. The set of such locations is called memory space.



Address generated by  
Programs

Actual Main memory  
address

## → Page Replacement

When a page fault occurs in a virtual memory system, it signifies that the page referenced by the CPU is not in the main memory.

If main memory is full, it would be necessary to remove a page from a memory block to make room for new page.

The most common replacement algorithms are:

- 1) FIFO (First In First Out)
- 2) LRU (Least Recently Used)
- 3) OPR (Optimal Page Replacement)

FIFO

- It selects for replacement the page that has been in memory for long time.
- Each time a page is loaded into memory, the identification number is pushed into FIFO stack.

Advantage - Easy to implement.

Disadvantage - Under certain circumstances pages are removed and loaded from memory too frequently.

LRU

- Implemented by associating a counter with every page
- Replace the page which has not been used for longest period of time.

**(OPT)**

- Optimal Replacement
- lowest page fault rate of all algorithms.

Consider three memory frames.

Reference String

First In First Out

**FIPO**

7 0 1 2 0 3 0 4 2 3 0 3 1 2 0

		1	1	1	X	0	0	Ø	3	3	3	3	2	2
f <sub>3</sub>														
f <sub>2</sub>	0	0	0	Ø	3	3	3	2	2	2	2	1	1	1
f <sub>1</sub>	7	7	7	2	2	2	2	4	4	4	0	0	0	0
*	*	*	*	*	Hit	*	*	*	*	*	*	Hit	*	*

Main  
Memory  
frame

$$\text{Page Hit} = 3, \quad \text{Page Miss} = 12, \quad \text{Hit Ratio} = \frac{\text{No. of Hits}}{\text{Total References}}$$

$$\text{Miss Ratio} = \frac{\text{No. of Miss}}{\text{Total Reference}} = \frac{3}{15} \times 100 = 20\%$$

$$= \frac{12}{15} \times 100 = 80\%$$

Optimal Page Replacement **OPR**

→ Replace the page which is not used in longest dimension of time in future.

		1	1	1	3	3	3	3	3	3	3	3	1	1
	0	0	0	0	0	Ø	4	4	4	0	0	0	0	0
	7	7	7	2	2	2	2	2	2	2	2	2	2	2
*	*	*	*	Hit	*	Hit	*	Hit	Hit	*	Hit	*	Hit	Hit

$$\text{Hit} = 7$$

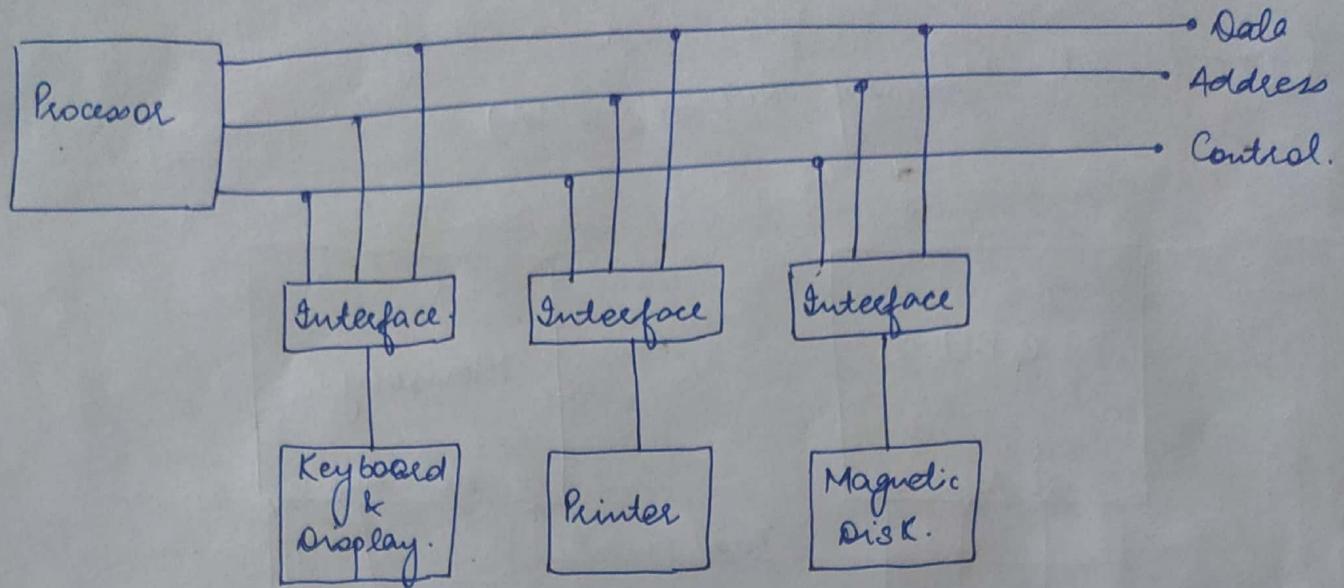
$$\text{Miss} = 8$$

$$\text{Hit Ratio} = \frac{7}{15} \times 100$$

$$\text{Miss Ratio} = \frac{8}{15} \times 100$$

Least Recently Used															LRU
7	0	1	2	0	3	0	4	2	3	0	3	1	2	0	Replace the page that is least recently used.
*	*	1	1	1	3	3	3	2	2	2	2	1	1	1	
0	0	0	0	0	0	0	0	3	3	3	3	3	3	0	
7	7	7	2	2	2	2	4	4	4	0	0	0	2	2	
*	*	*	*	Hit	*	Hit	*	*	*	*	Hit	*	*	*	

- If the page demanded, is not present in the main memory is refer as page fault.
- Then information is fetch from hard disk (virtual memory) and placed in main memory.
- If the frames in main memory are full, then replace the page (via page replacement algorithm)
- If the demanded page is already placed in main memory it is refer as hit.

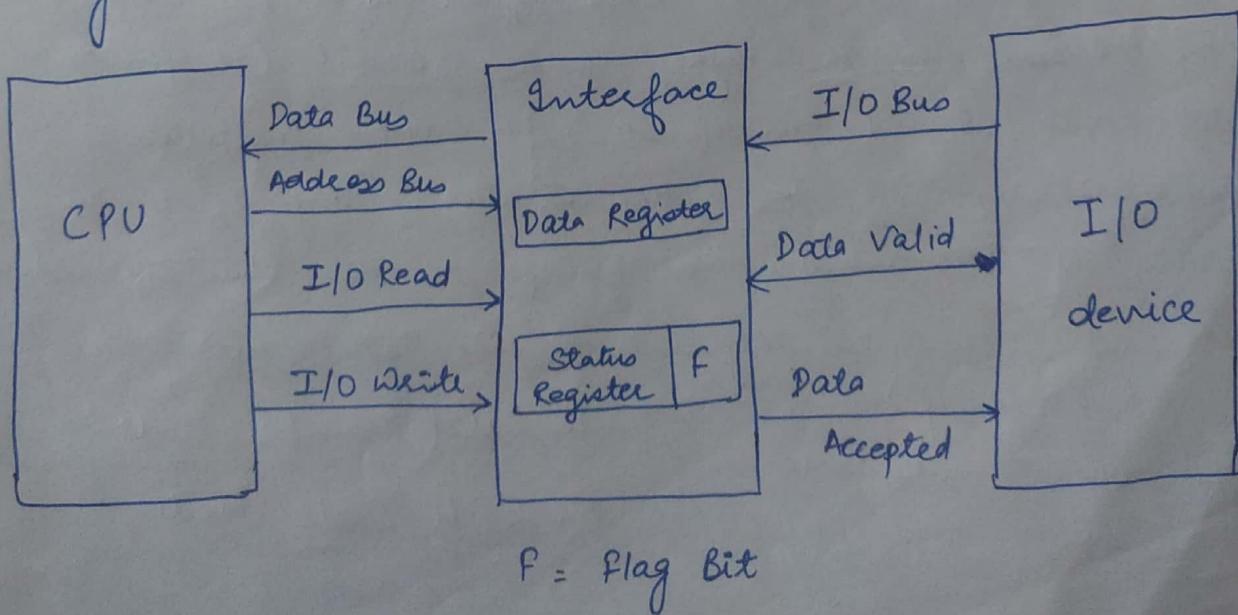


## Need of Interface

- |           |                                    |                 |
|-----------|------------------------------------|-----------------|
| 1) Speed  | Million of instructions<br>in secs | Slow            |
| 2) Signal |                                    | Electromagnetic |
| 3) Format | 32 byte, 64 byte                   |                 |
| 4)        |                                    |                 |

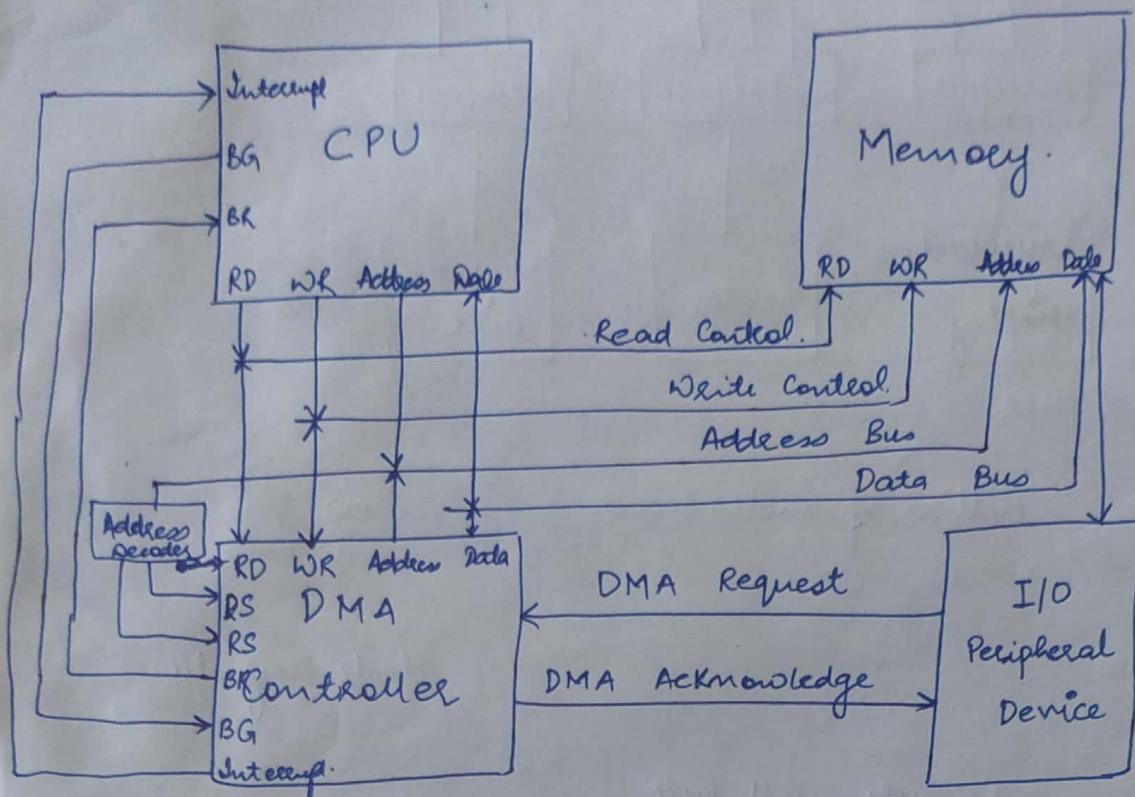
## Input - Output Transfer

① Programmed I/O



② Interrupt I/O Initiated

③ Direct Memory Access (DMA)



→ DMA transfers the data b/w I/O to Memory w/o using CPU

→ When I/O device is ready to send data, it sends DMA request (i.e ready to transfer the data)

then DMA controller sends BR (Bus Request to CPU), CPU stands RD, WR, Address, Data in high impedance circuit (i.e open circuit)

Now CPU sends BG (Bus Grant) to DMA controller.

Now bus is handed over to the DMA controller. Now it is the responsibility of DMA controller to transfer the data b/w I/O and Memory

Now DMA sends DMA Acknowledgment to the I/O device.

(i.e DMA is ready)

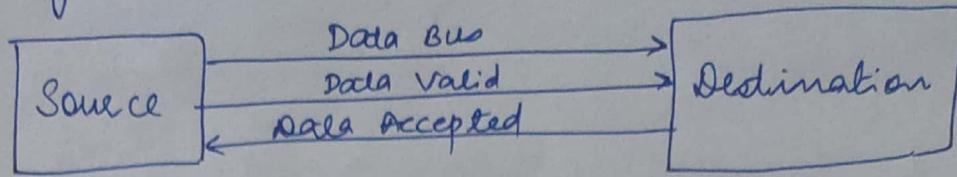
$BG = 1$  Buses with DMA  
 $= 0$  Bus with CPU

1. Burst Transfer Mode.  
Allows at a time (I/O)

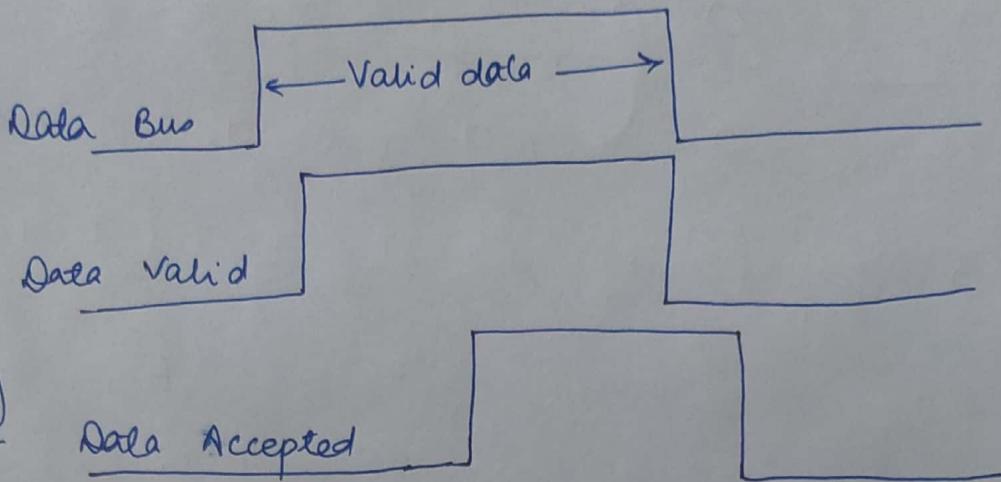
2. Cycle Steal Mode

## Handshaking.

Block Diagram



Timing Diagram



## Sequence

- ① Place data on bus
- ② Enable data valid .
- ③ Accept data from bus.
- ④ Enable data accepted .
- ⑤ Disable data valid
- ⑥ Invalidate data on bus
- ⑦ Disable data accepted .
- ⑧ Ready to accept data .