

Computer Architecture :-

(7-8 marks)

(5-6 questions)

Syllabus:-

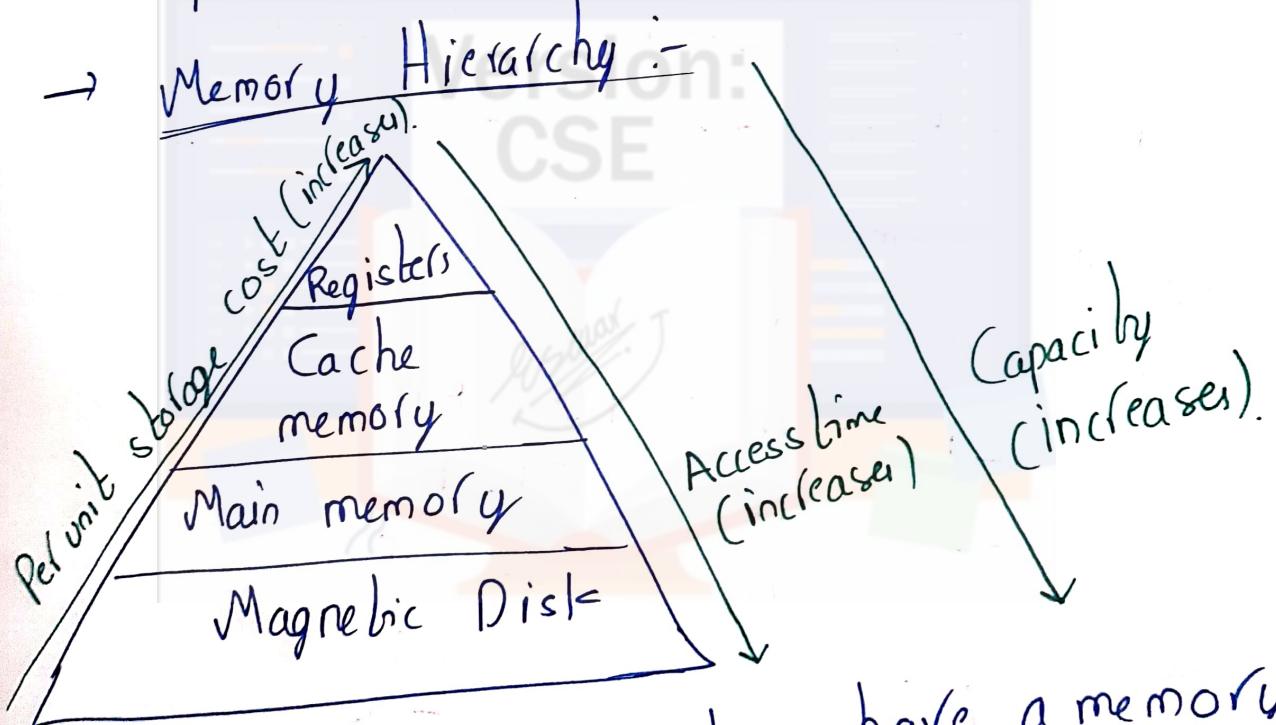
- Machine instructions & addressing modes ⁽⁴⁾
- ALU, data path & control unit ⁽⁵⁾
- Instruction pipelining ⁽³⁾
- Memory hierarchy ⁽¹⁾
- I/O interface. ⁽²⁾

(Morris Mano)

→ Computer Organization :-

Computer organization refers to the way in which the hardware components operate & the way they are connected together that realize the architectural specifications.

→ Memory Hierarchy :-



Our goal is to have a memory which has → Large capacity
→ Less per unit cost
→ Less access time
(Fast access).

but capacity $\propto \frac{1}{\text{access time}}$
hence both can't be low
that is the main reason why
we use multiple types of storage
devices.

Access Time / speed :- Register > Cache > MM > SM

Cost :- Register > Cache > MM > SM

Size :- Register < Cache < MM < SM,

ex:- my system
3 levels of cache.

L₁ cache :- 128 kB

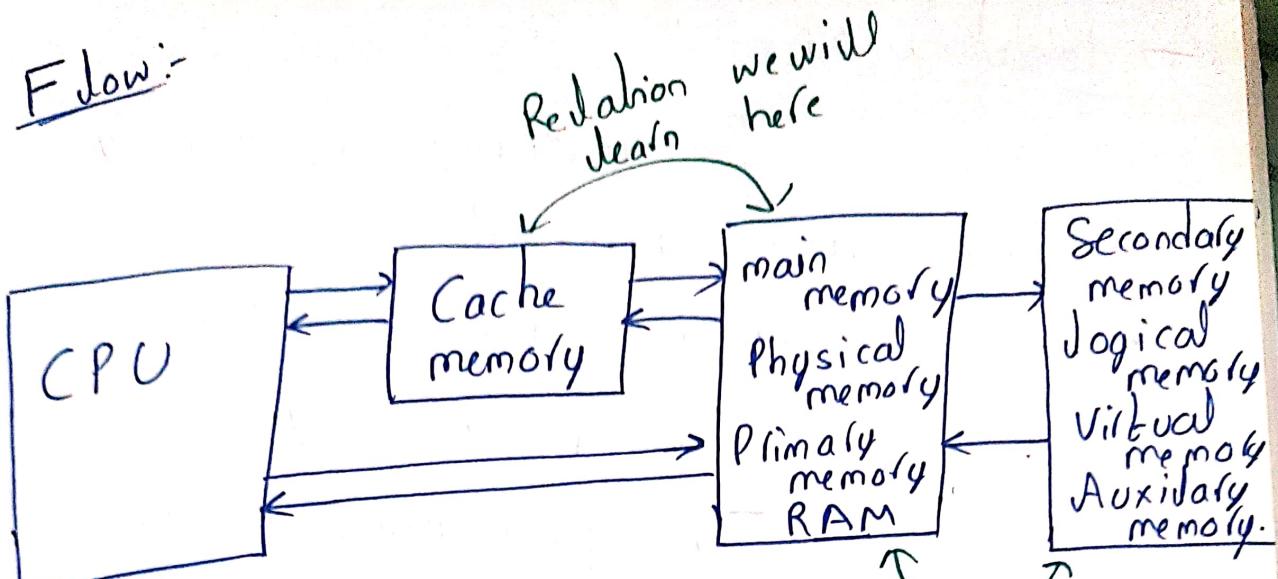
L₂ cache :- 512 kB

L₃ cache :- 3.0 MB

MM :- 8 GB

SM :- 2 TB

Flow:-



Logic:- if we are searching for a file/data
First we check in cache
if found its a (Cache hit)
else its a (Cache miss)
if not found then we search
in Main memory.
if Found copy it to cache
else with the help of main memory
search in secondary memory if found
bring to main memory & use it
else not found.

→ Locality of Reference:

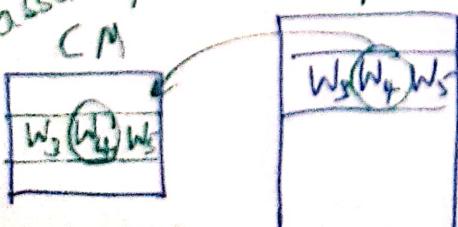
The references to memory at any given interval of time tend to be confined within a few localized areas in memory.

This phenomenon is known as property of Locality of reference.

Spatial Locality
(w.r.t space)

→ If a word is accessed now then the word adjacent to it will be accessed next.

whole block brought to cache assuming w_3 or w_5 will be accessed next.



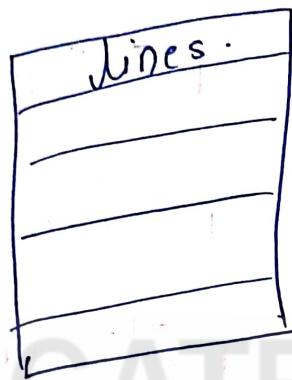
Temporal Locality
(w.r.t time).

→ It refers to the reuse of specific data or resources within a relatively small-time during. Or the most frequently used will be need soon.

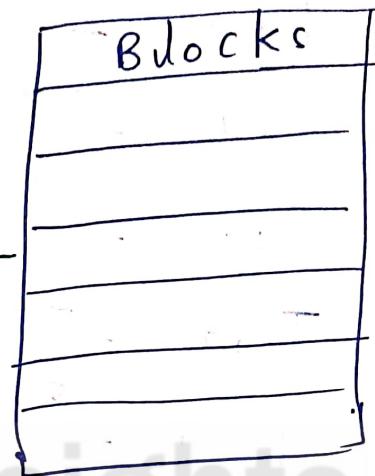
→ LRU (example).
(least recently used).

Note:

Cache Memory



Main Memory



Note size of lines - Blocks size.

Hence as we know MM size is greater than Cache memory.

Hence at any moment of time.

all blocks can't be in cache

Hence we using mapping requirement.
give slots as per

3 types.

i) Direct Mapping

ii) Associative Mapping.

iii) Set Associative Mapping.

Metric * * space
(very imp)

1 Kilo 2^{10}
1 Mega 2^{20}
1 Giga 2^{30}
1 Tera 2^{40}
1 Peta 2^{50}

Time
 1 sec
 $10^{-3} \text{ millisecond}$
 $10^{-6} \text{ micro second}$
 $\mu \text{ sec}$
 $10^{-9} \text{ nano second}$
 nsec.
 $10^{-12} \text{ picosecond}$
..

Basic:
(smallest unit we call a word)
(collection of words form block)
(collection of blocks form main mem).

| | Main memory | | | |
|----|-------------|----------|----------|----------|
| | w_0 | w_1 | w_2 | w_3 |
| B0 | w_0 | w_1 | w_2 | w_3 |
| B1 | w_4 | w_5 | w_6 | w_7 |
| B2 | w_8 | w_9 | w_{10} | w_{11} |
| B3 | w_{12} | w_{13} | w_{14} | w_{15} |
| B4 | w_{16} | w_{17} | w_{18} | w_{19} |
| B5 | w_{20} | w_{21} | w_{22} | w_{23} |
| B6 | w_{24} | w_{25} | w_{26} | w_{27} |
| B7 | w_{28} | w_{29} | w_{30} | w_{31} |

hence to find
 $w_{18} \Rightarrow$ write binary of 18

1 0 0 1 0

We have 8 blocks
means 3 bits
to differentiate

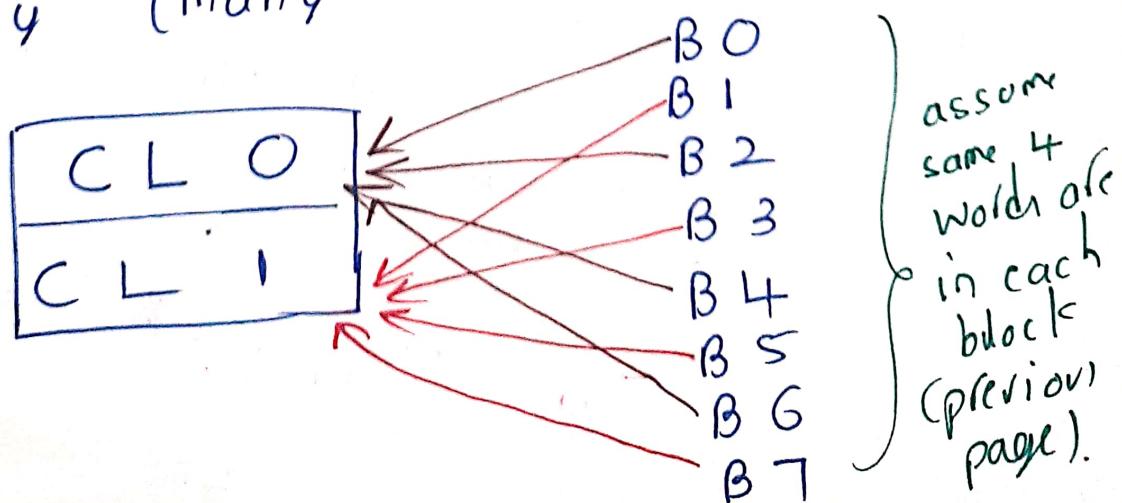
we have 4 word,
2 bits to
differentiate

i.e.
1 0 0 1 0
Block no. 4 - 2 Block offset

Hence in this we can get location.

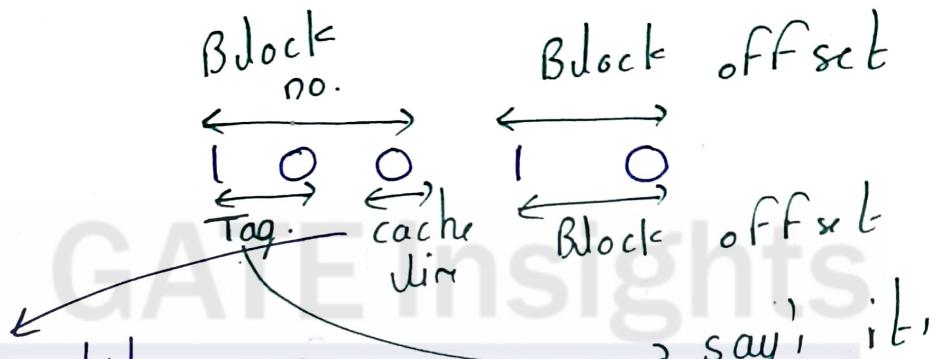
I) Direct Mapping:

In direct mapping scheme the main memory blocks are directly mapped onto a particular cache memory (many to one mapping).



to find exactly to which line

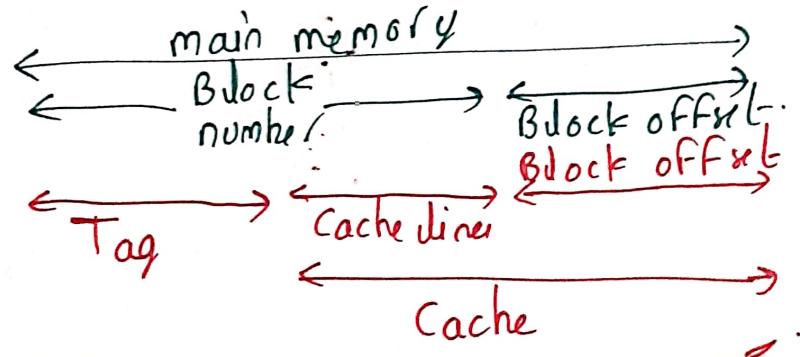
W_{18} goes
→ 10010



only 1 bit
coz we
have only 2 lines
1 bit is enough
0 means CL-0

CL-0 ⇒ B 0
B 2
B 4
B 6
2nd place (10)

1st time a bit confusing but
only remember the diagram it will be
easy.



For solving problems if you remember previous diagram is enough;

ex: given (direct mapping).

main memory size = 2^{32} bytes

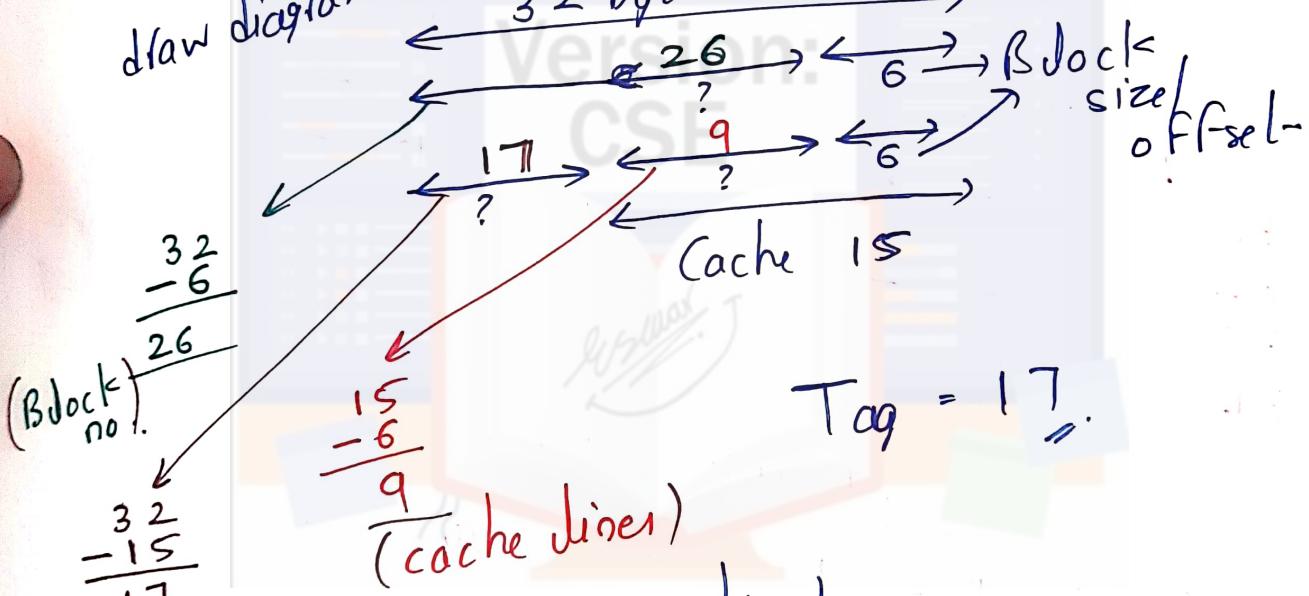
cache size 32 KB = 2^{15} bytes.

cache block size = 2^6 bytes.

(Offset) Tag = ?

(Date 2021)

draw diagram



(Simple questions).

Problem: with direct mapping:

Here we are assigning a particular cache line to a set of blocks.

Hence if that particular
blocks are called again & again
Even we have space in cache
memory as we are forced that
slot only for them hence this
is any issue.

II)

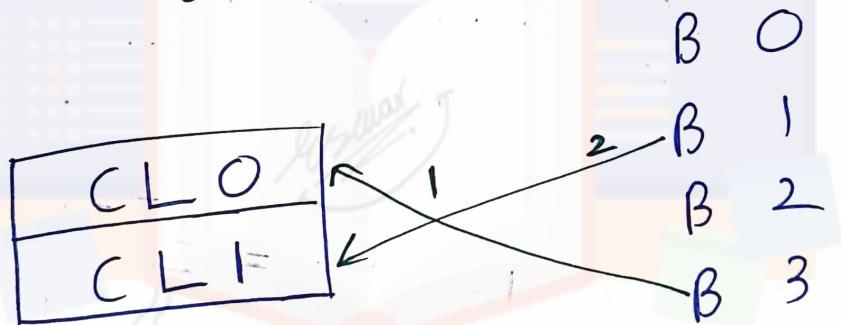
| ex: | CL 0 | B 0 |
|-----|------|-----|
| | CL 1 | B 1 |
| | CL 2 | B 2 |
| | CL 3 | B 3 |
| | : | : |
| | : | : |
| | : | : |
| | : | : |

Assume if a program used
B0, B4, B8, B12 in exact order
multiple times as they all
map to CL0 even CL1, 2, 3
are empty as this is direct
mapping we push them only B 16,,
to CL0

Hence due to this issue only we designed Associative mapping.

II) Associative Mapping:

here A block of main memory can be mapped freely to any available cache line. This makes fulling associative mapping more feasible than direct mapping (many to many mapping).



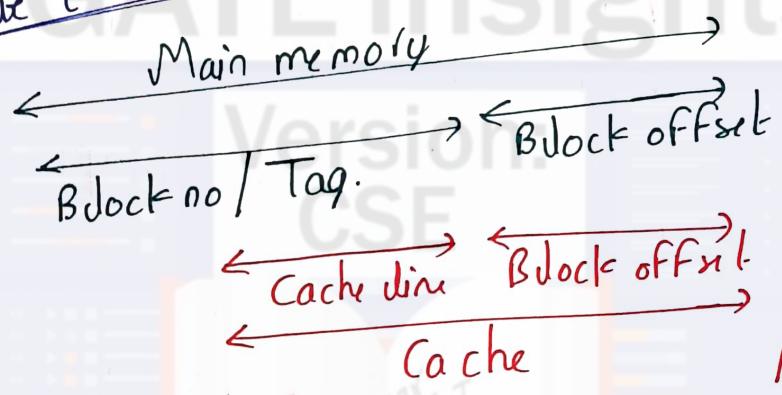
B₃, B₁ hence based on their respective order they get the slots. (no pattern). If Cache Full we Cache replacement algo's will learn soon =

hence as there is no pattern
here tag = Block no.

CL0 Tag = block no

CL1 Tag = block no //

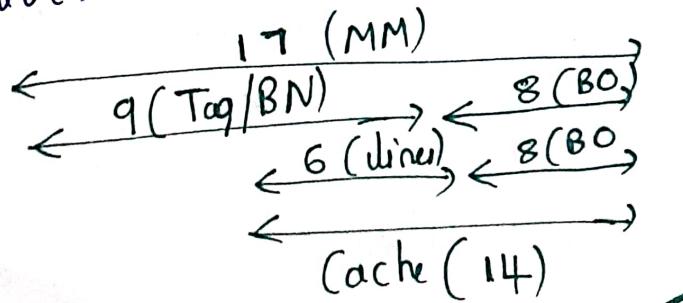
Update the diagram:



ex) MM size = 128 KB = $2^{7+10} = 2^{17}$

Cache size = 16 KB = $2^{4+10} = 2^{14}$

Blocksize = 256 B = 2^8



No. of tag bits = 9

$$\begin{aligned}\text{Tag directory size} &= 9 \times 2^6 \\ &= 9 * 64\text{.}\end{aligned}$$

As here in associative we gave full freedom which made tag bit = block no. hence which is not ideal.

III) Set Associative Mapping:

In K-way set associative mapping, cache lines are grouped into sets where each set contains K number of lines.

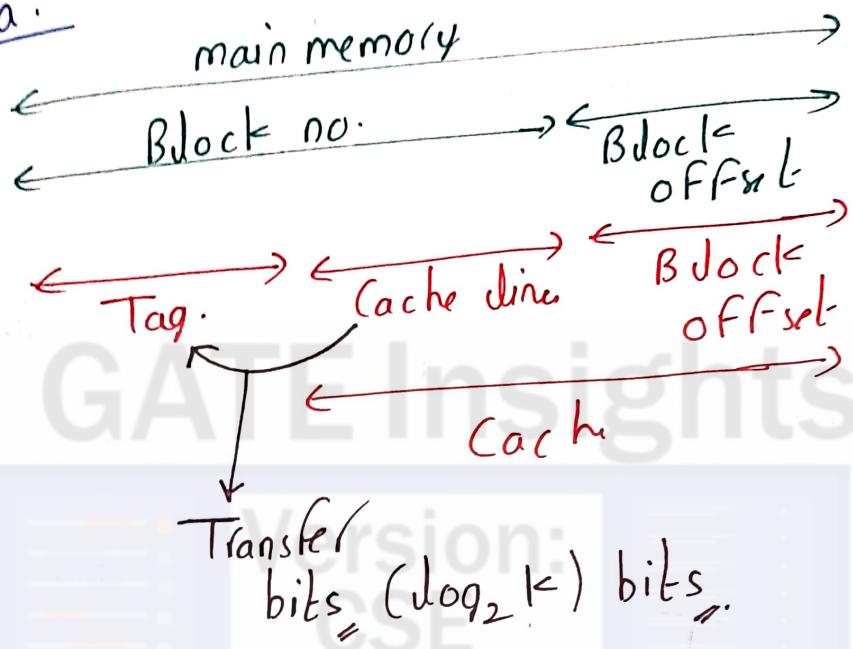
| | |
|-------|------|
| Set 0 | CL 0 |
| | CL 1 |
| Set 1 | CL 2 |
| | CL 3 |

alternate logic

B 15

Hence set associate is a combination of associative & direct.

formula:



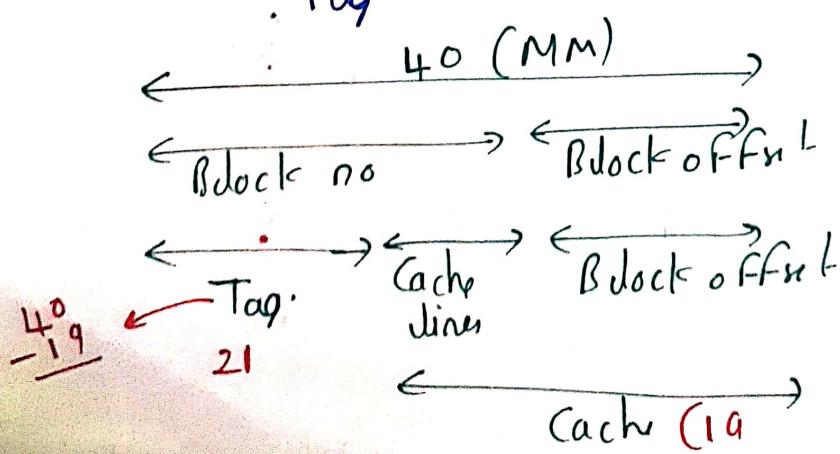
ex:

$$\text{Main memory} = 2^{40} \text{ bytes}$$

$$\text{Cache} = 512 \text{ KB} = 2^{19} \text{ bytes}$$

8 way set associative.

Tag size = ?

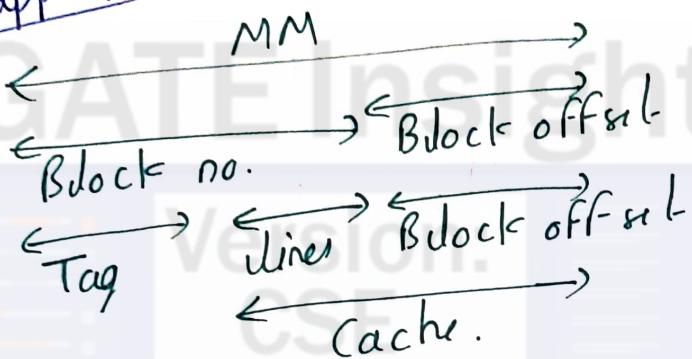


Grade
2016/1

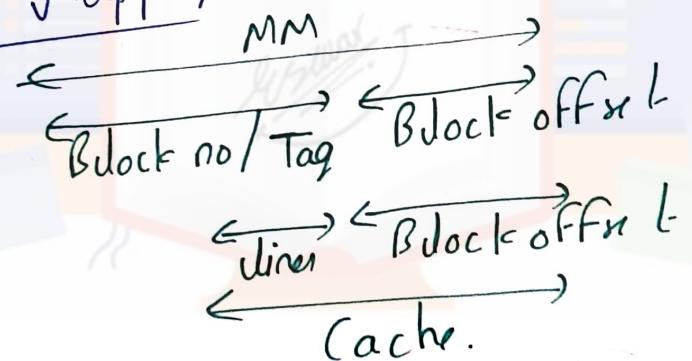
8 way
 ↙ 3
 $2^1 + 3 = 2^4$,
 24 bits of Tag.

Note:

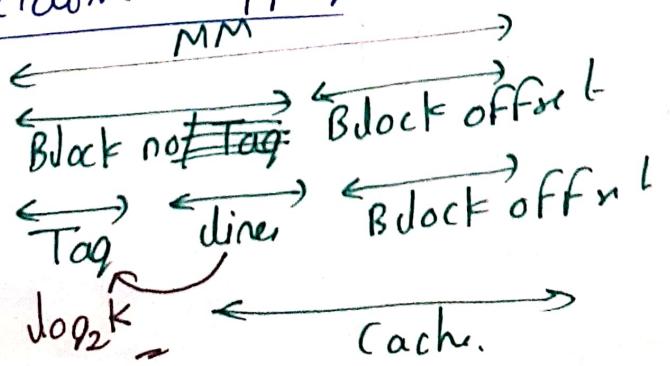
Direct Mapping:



Associative Mapping:



Set Associative mapping:



1.

→ Cache Replacement Policies:

As we know when cache is full we need to replace some to give place for a block. Hence to decide whom to remove/replace we use Cache Replacement Policies.

In direct mapped cache, the position of each block is predetermined hence no replacement policy needed.

In Fully associative & set associative caches there exist policies

- with repeat in OS.
- Type I
- i) FIFO
 - ii) LRU
 - iii) Optimal.

I) FIFO (First in First out)

The blocks which have entered first in the memory will be replaced first.

Block Sequence:

7, 0, 1, 2, 0, 3, 0, 4, 2, 3, 0, 3, 2, 1, 2, 0, 1, 7,
0, 1 & cache memory has 4 lines.

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 0 | 1 | 2 | 0 | 3 | 0 | 4 | 2 | 3 | 0 | 3 | 2 | 1 | 2 | 0 | 1 | 7 | 0 |
| 7 | 7 | 7 | 7 | 3 | 3 | 3 | 0 | 3 | 3 | 3 | 2 | 3 | 2 | 2 | 1 | 1 | 1 | 7 |
| 0 | 0 | 0 | 0 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 0 | 0 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| 2 | 1 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 |
| M | M | M | M | M | M | M | M | M | M | M | M | M | M | M | M | M | M | M |

10-miss Total / Faults.

Belady's Anomaly: (We assume if no. of cache lines increase faults will decrease)
But in FIFO we might have some cases in which faults increase if we increase cache lines i.e Belady's Anomaly.

e.g. (3 cache lines)

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 1 | 2 | 5 | 1 | 2 | 3 | 4 | 5 |
| 1 | 1 | 1 | 4 | 4 | 4 | 5 | 5 | 5 | | | |
| 2 | 2 | 2 | 2 | 1 | 1 | 1 | 3 | 3 | | | |
| 3 | 3 | 3 | 3 | 2 | 2 | 2 | 2 | 4 | | | |
| M | M | M | M | M | M | M | M | M | M | M | M |

9-faults

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 1 | 2 | 5 | 1 | 2 | 3 | 4 | 5 |
| 1 | 1 | 1 | 1 | 1 | 1 | 5 | 5 | 5 | 5 | 4 | 4 |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 5 |
| 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 2 | 2 | 2 | 2 |
| 4 | 4 | 4 | 4 | 4 | 4 | 4 | 4 | 3 | 3 | 3 | 3 |
| M | M | M | M | M | M | M | M | M | M | M | M |

10 Fault

3 lines | 4 lines } BeJady's
9 faults, | 10 faults } anomaly,

II) LRU (least recently used). (Past)

The page which was not used for the longest period of time in the past will get replaced first.

| | | | | | | | | | | | | | | | | | | |
|----------|----------|----------|----------|----------|----------|---|---|---|---|----------|---|---|---|---|---|---|---|----------|
| 7 | 0 | 1 | 2 | 0 | 3 | 0 | 4 | 2 | 3 | 0 | 3 | 2 | 1 | 2 | 0 | 1 | 7 | 0 |
| 7 | 7 | 7 | 7 | 3 | 3 | | | | | 3 | | | | | | | | |
| 0 | 0 | 0 | 0 | 0 | 0 | | | | | 0 | | | | | | | | |
| 1 | 1 | 1 | 4 | | | | | | | 1 | | | | | | | | |
| 2 | 2 | 2 | 2 | | | | | | | 2 | | | | | | | | 2 |
| <u>M</u> | <u>M</u> | <u>M</u> | <u>M</u> | <u>M</u> | <u>M</u> | | | | | <u>M</u> | | | | | | | | <u>M</u> |

8 faults.

III) Optimal algorithm: (Future work).

The page which will not be used for the longest period of time in future references will be replaced first.

Best performance but can't be used as we have no idea about future.

| | | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|----------|----------|---|---|---|---|----------|---|---|---|---|---|----------|---|
| 7 | 0 | 1 | 2 | 0 | 3 | 0 | 4 | 2 | 3 | 0 | 3 | 2 | 1 | 2 | 0 | 1 | 7 | 0 |
| 7 | 7 | 7 | 7 | 7 | 3 | 3 | | | | | 3 | | | | | | 7 | |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | | | | | 0 | | | | | | 0 | |
| 1 | 1 | 1 | 1 | 1 | 4 | | | | | | 1 | | | | | | 1 | |
| 2 | 2 | 2 | 2 | 2 | | | | | | | 2 | | | | | | 2 | |
| M | M | M | M | M | <u>M</u> | <u>M</u> | | | | | <u>M</u> | | | | | | <u>M</u> | |

8 - Faults

Type of miss:

Compulsory Miss: When CPU demands for any block for the first time then definitely a miss is going to occur as the block needs to be brought into cache. (brown underline).

Capacity Miss: Occur because blocks are being discarded from cache because cache cannot contain all blocks needed for program execution (red underline).

Conflict Miss :-

In the case of set associative or direct mapped block placement strategies, conflict misses occur when several blocks are mapped to the same set/block frame.

ex:- consider a 4-way set associative cache with 16 cache blocks. MM consists of 256 blocks. request order.

0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8, 48, 32, 73, 92, 155.

which of the following will not be in cache.

A) 3

B) 8

C) 129

D) 216

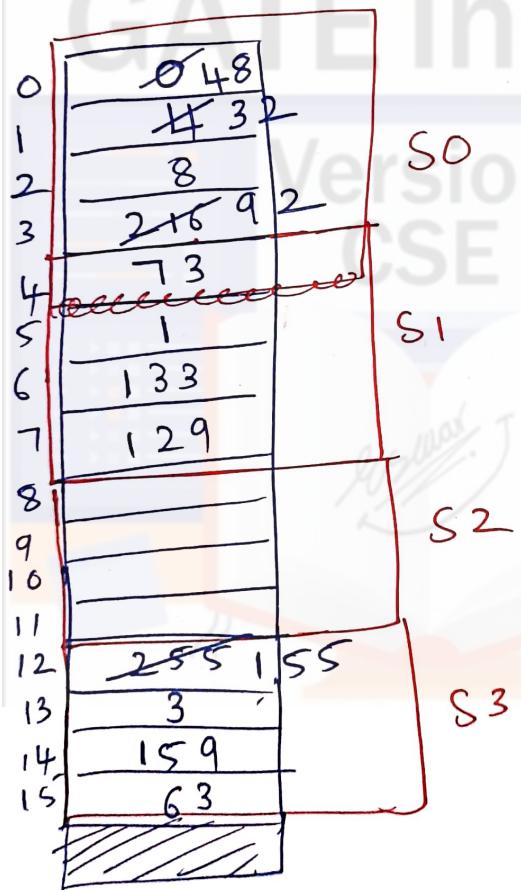
4 way set associative
4 lines in 1 set

16 cache blocks

256 blocks in memory

0, 255, 1, 4, 3, 8, 133, 159, 216, 129, 63, 8,
48, 32, 73, 92, 155.

same like hashing

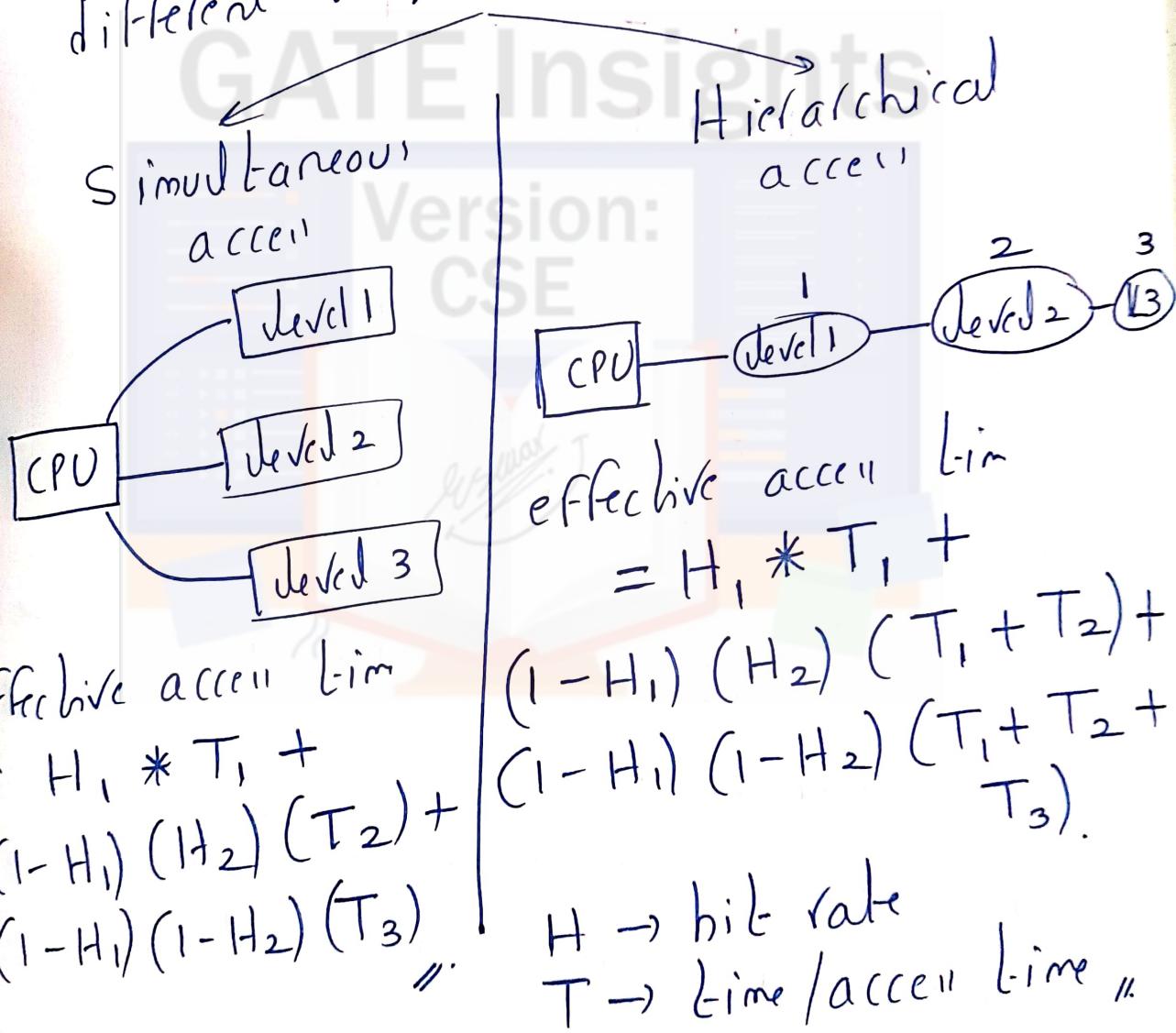


$$\begin{aligned}0 \bmod 4 &= 0 \\255 \bmod 4 &= 3 \\1 \bmod 4 &= 1 \\4 \bmod 4 &= 0 \\3 \bmod 4 &= 3 \\8 \bmod 4 &= 0 \\133 \bmod 4 &= 1 \\159 \bmod 4 &= 3 \\216 \bmod 4 &= 0 \\129 \bmod 4 &= 1 \\63 \bmod 4 &= 3 \\8 \bmod 4 &= 0 \\48 \bmod 4 &= 0 \\32 \bmod 4 &= 0 \\73 \bmod 4 &= 1 \\92 \bmod 4 &= 0 \\155 \bmod 4 &= 3\end{aligned}$$

→ Memory Organization:

→ Memory is organized at different levels.

→ CPU may try to access different levels of memory in different ways



ex:

(Grade 2015) { read request miss \Rightarrow 50 nsec }
 on cache. { read request hit \rightarrow 5 nsec } $\rightarrow T_1$
 80% request hit

$$\Rightarrow H_1 * (T_1) + (H_1 - 1) H_2 (T_1 + T_2)$$

$$\Rightarrow 0.8 * 5 + (0.2)(1)(50)$$

$$\Rightarrow 4 + 10 \text{ nsec}$$

$$\Rightarrow 14 \text{ nsec.}$$

→ Cache Coherence Problem :-

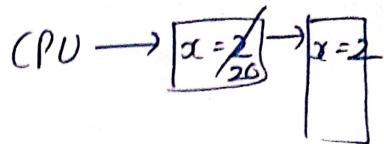
As we know the data which we store in cache is a copy of MM.
 Hence as multiple copies at different level of memory, the inconsistency may occur, this problem is known as Cache Coherence Problem.

+T₂
It can be solved using 2

techniques.

→ Write Through

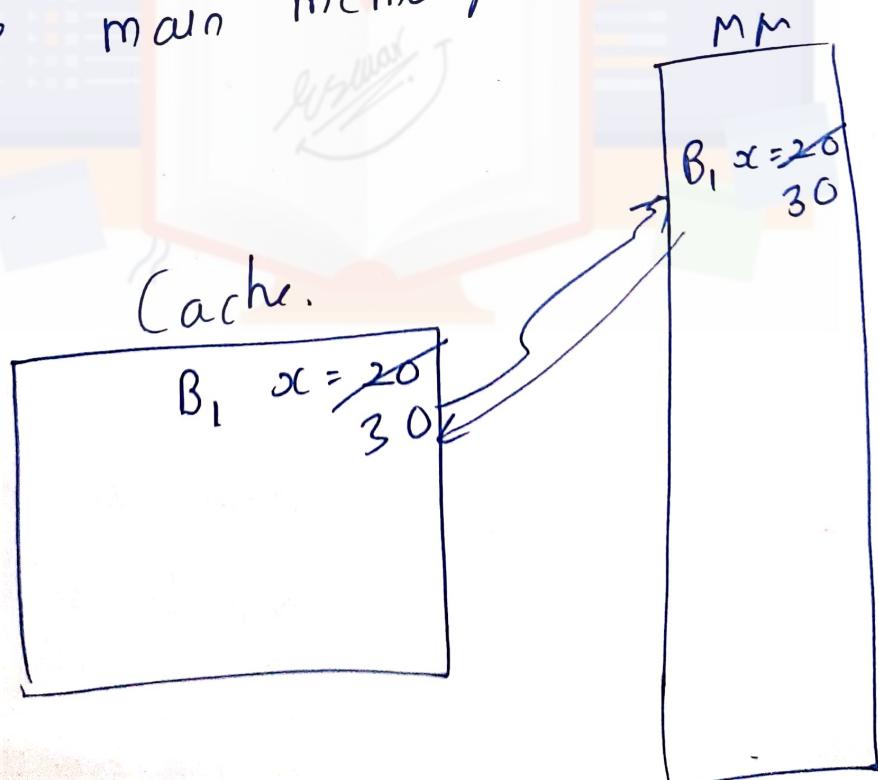
→ Write Back



I) Write Through:

Write through is used to maintain the consistency between the cache & main memory.

According to it if the cache copy is updated at the same time main memory is also updated.



Adv: It provides the highest level of consistency.

Disadv: It requires more number of memory access.

II) Write Back:

Write back is also used to maintain the consistency between the cache and main memory.

According to it all the changes performed on cache are

reflected back to the main memory in the end. (We maintain a bit modified bit $\leftarrow 1^{(\text{modified})} 0^{(\text{not modified})}$)

Adv:

→ less number of memory access & less write operations

Disadv:

→ Inconsistency may occur

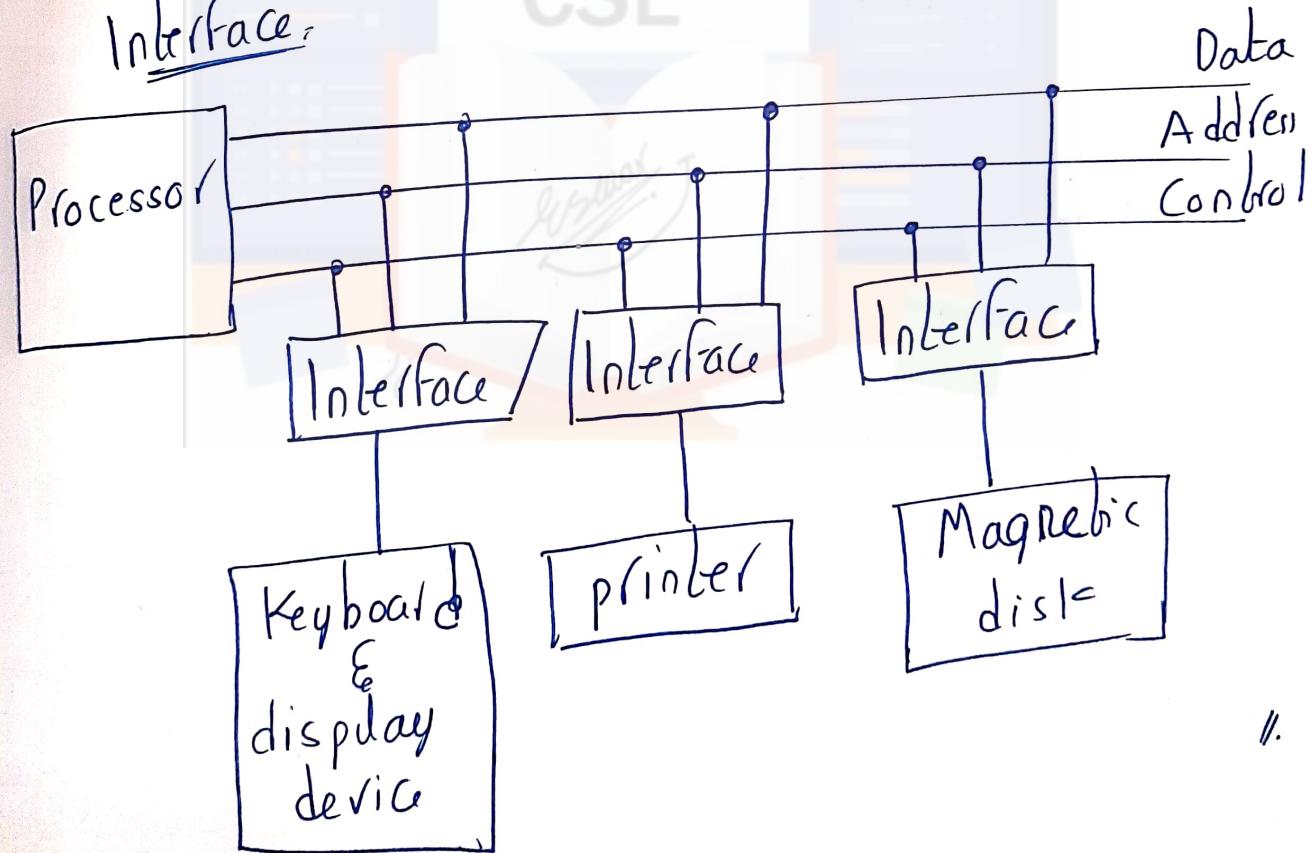
→ I/O management:

As we know the next most important component after CPU & memory is input & output (as without input & output the computer cannot serve any).

We cannot directly connect a I/O device to computer.

We use interface because

Interface:



Reason for having interface is:

Speed: The speed of CPU & I/O device will usually be different.

Format: The data code & Format of CPU & peripherals may be different
eg: ASCII, Unicode etc..

Physical orientation: Different devices have organizations like optical, magnetic, etc..

Signal conversion: peripherals are electro-magnetic & their manner of operation is different from the operation of CPU.

(Analog to digital etc..)

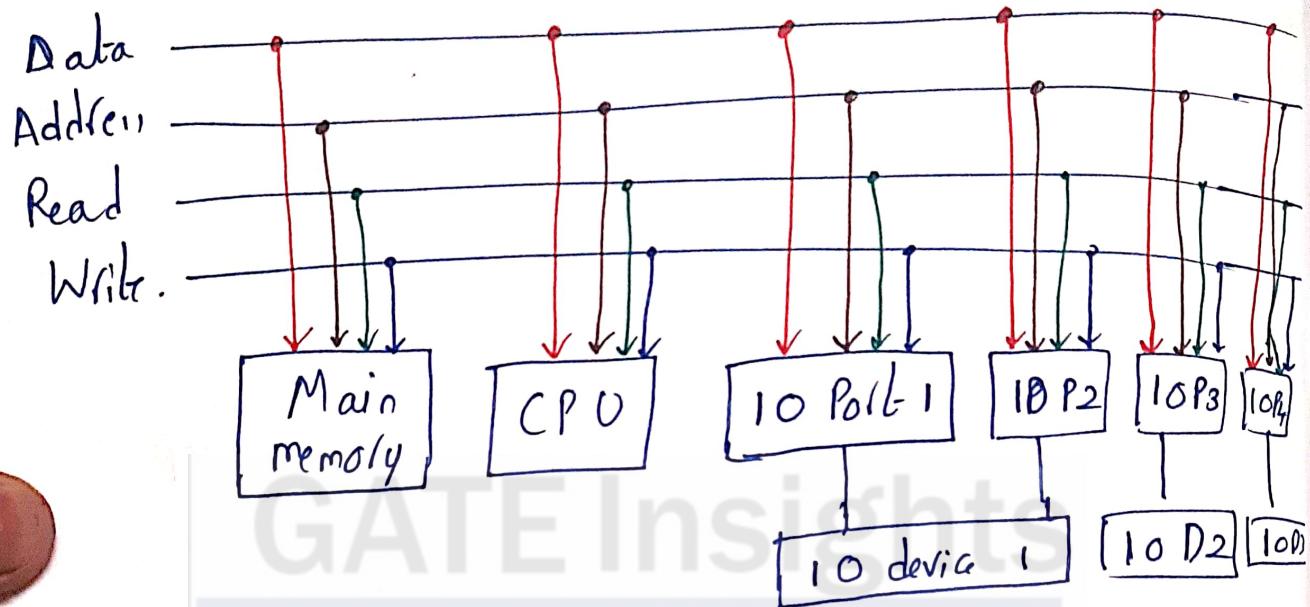
Address bus: Is Used to identify the correct i/o device among the number of i/o devices.

Control bus: After selecting a specific i/o device CPU sends a functional code on the control line. (read, write).

Data bus: In this final step depending on the operation either CPU will put data on the data line & device will store/display it & vice versa.

Based on the connection b/n CPU & i/o devices
3 ways:-

I) Memory Mapped I/O:

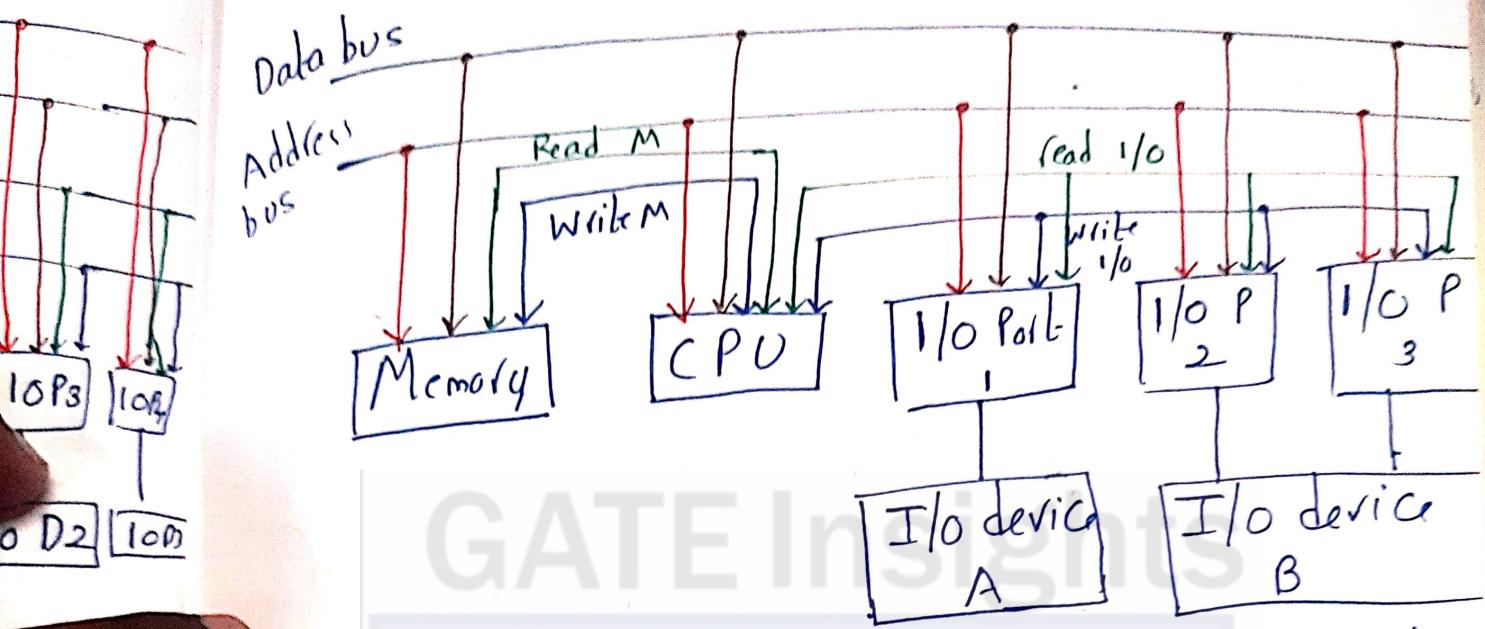


It uses the same address space
to address both memory & I/O
ex: 8085

Adv: In typical computer there are
more memory reference instruction than i/o
instruction, but in memory mapped i/o all
instruction that refer to memory are also
available for i/o

Disadv: Total address get divided.
some range is occupied by i/o while
some by memory

II) Isolated I/O:

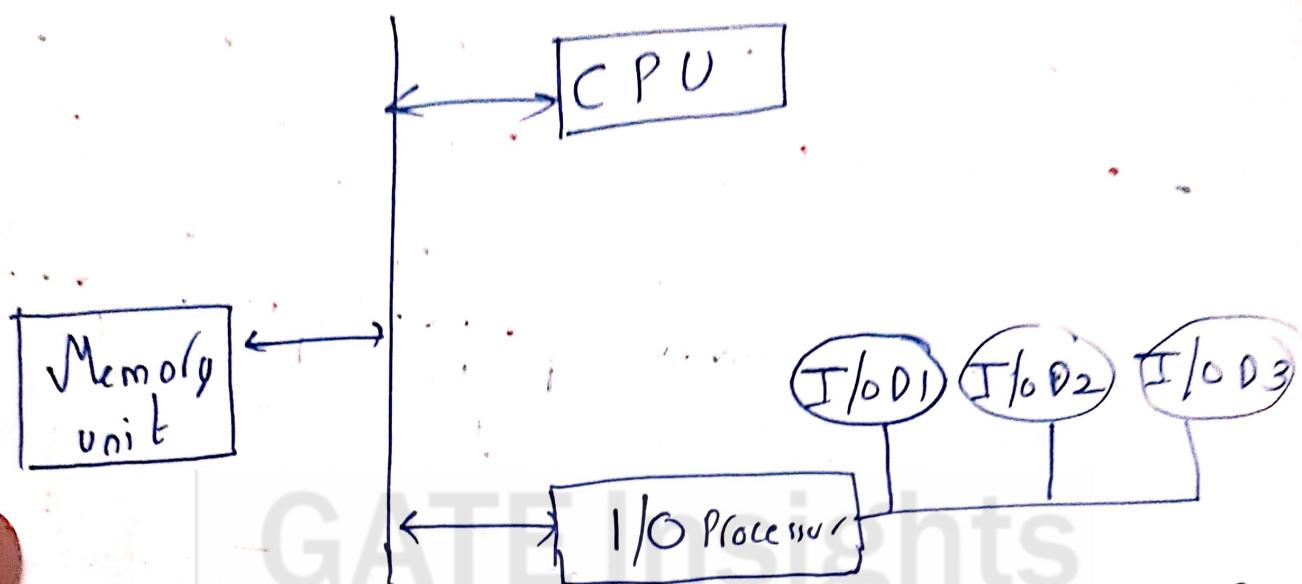


Here common bus to transfer data b/n the memory or I/O & CPU. The distinction between a memory & I/O transfer is made through separate read & write line.

Adv: Here memory is used efficiently as the address can be used two times
in 8086.

DisAdv: Need different control lines one for memory & other for I/O device.

III) I/O Processor:



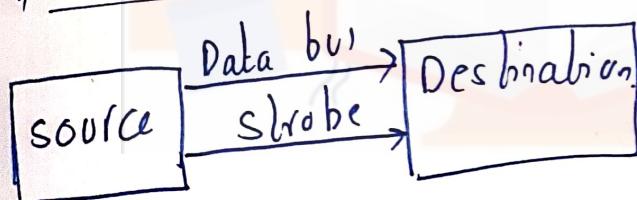
Computer has independent set of data, address & control buses, one for accessing memory & other for I/O. This is done in computer that provides a separate I/O processor other than CPU.

Synchronous v/s Asynchronous data transfer

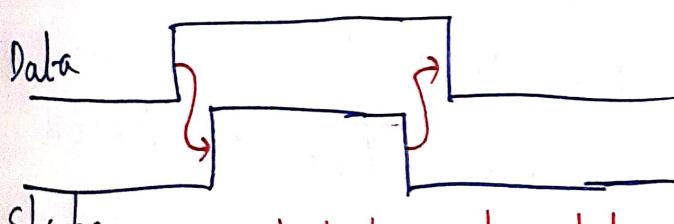
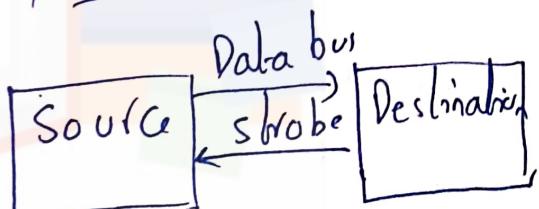
Synchronization is achieved by a device called master generator which generates a periodic train of clock pulse (The clock by CPU).

In Asynchronous data timing units of two devices are independent that is they are under different control.

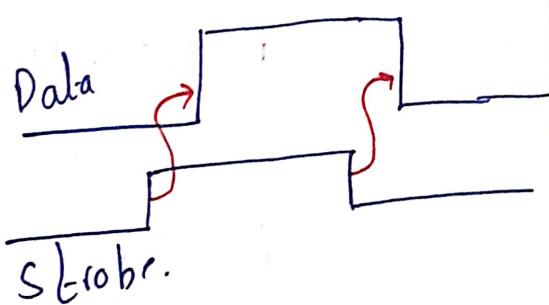
i) Source initiated



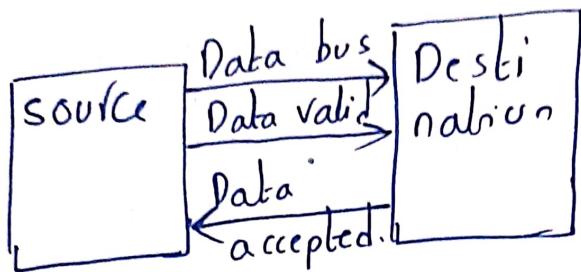
ii) Destination initiated



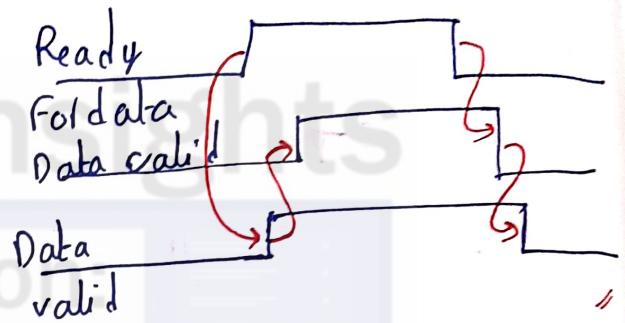
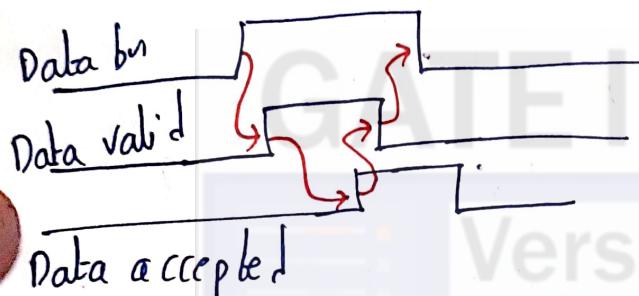
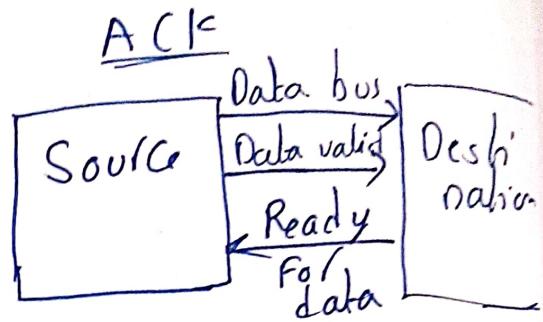
Strobe is a signal high when data is ready to copy.



iii) Source initiated with ACK



iv) Destination initiated with ACK



→ Mode of Data Transfer :-

here we discuss how data communication will take place b/n CPU

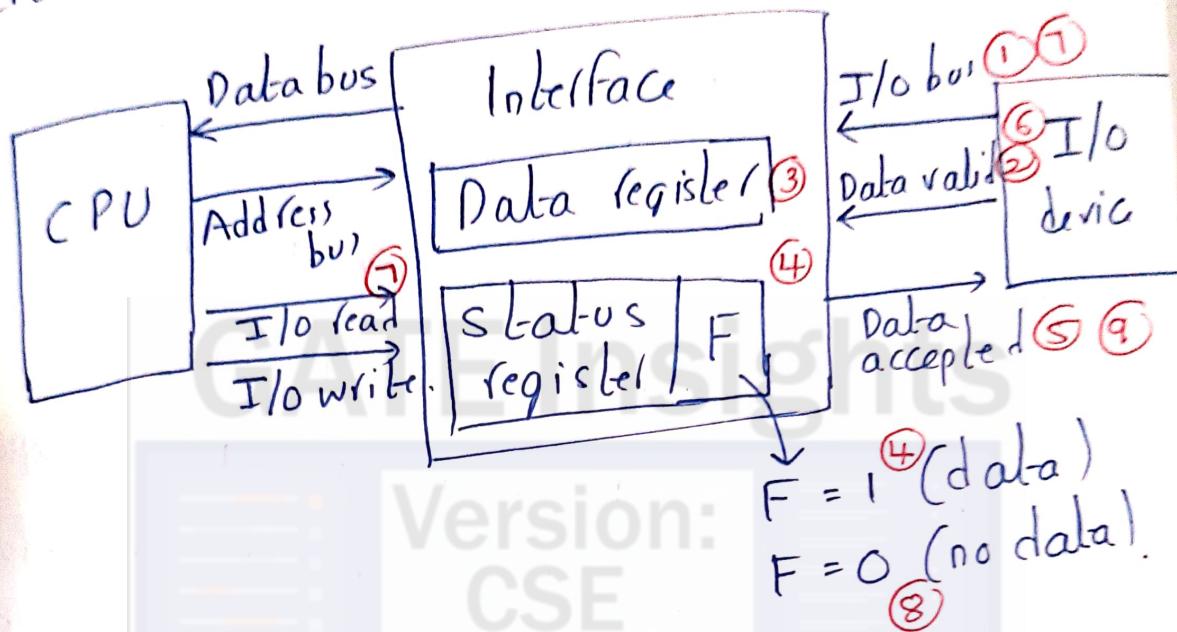
& I/O device.

3 methods

- i) Programmed I/O
- ii) Interrupt-initiated I/O
- iii) Direct Memory Transfer / Access

I) Programmed I/O:

In this I/O device cannot access the memory directly.



The flow is quite simple.
Note: Here CPU always checks Flag bit
whether it's 1 or 0
Hence many cycles of CPU are wasted

here.
Hence indirectly CPU is like slave;
Hence this method is not an ideal way.

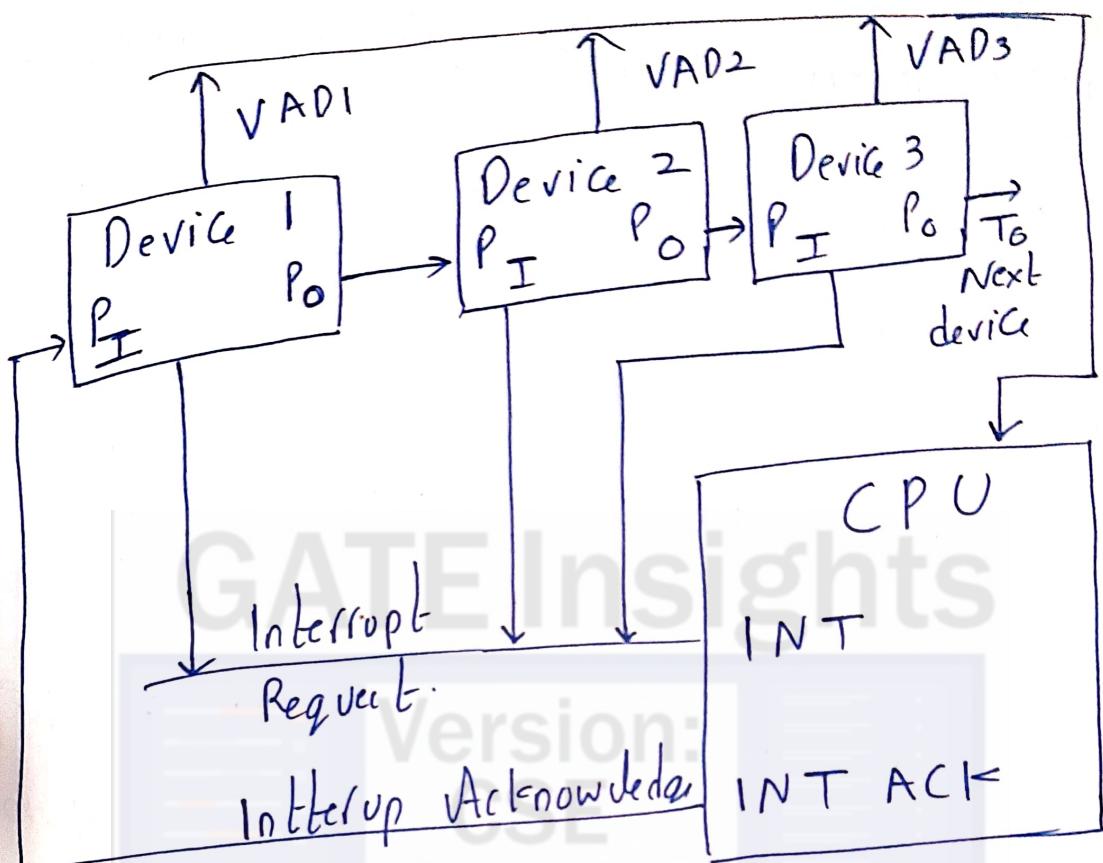
II) Interrupt initiated I/O :-

In this method I/O device interrupt CPU when it is ready for data transfer.

CPU keep executing instruction & after executing one instruction & before starting another instruction CPU wait & see if there is interrupt or not & if there is interrupt then taken a decision whether CPU should entertain this interrupt or continue with the execution.

Note: The one which we discuss is daisy chain

It has priority.



If an I/O device wants to communicate, it raises interrupt ($O_{II}(1, 2, 3)$). Once CPU completes its current instruction it will check & decide whether to continue or to work on I/O. It sends INT ACK to D_1 , if D_1 has request it will complete it and will move to $D_2 \dots$ (Priority $D_1 > D_2 > D_3$).

In this way it works.

Here interrupt are of 2 types.

i) Non vectored interrupt: Here there is a mutual understanding between CPU & device that where this routine is stored in memory (high priority device) (Well known friends/beil-friends?)

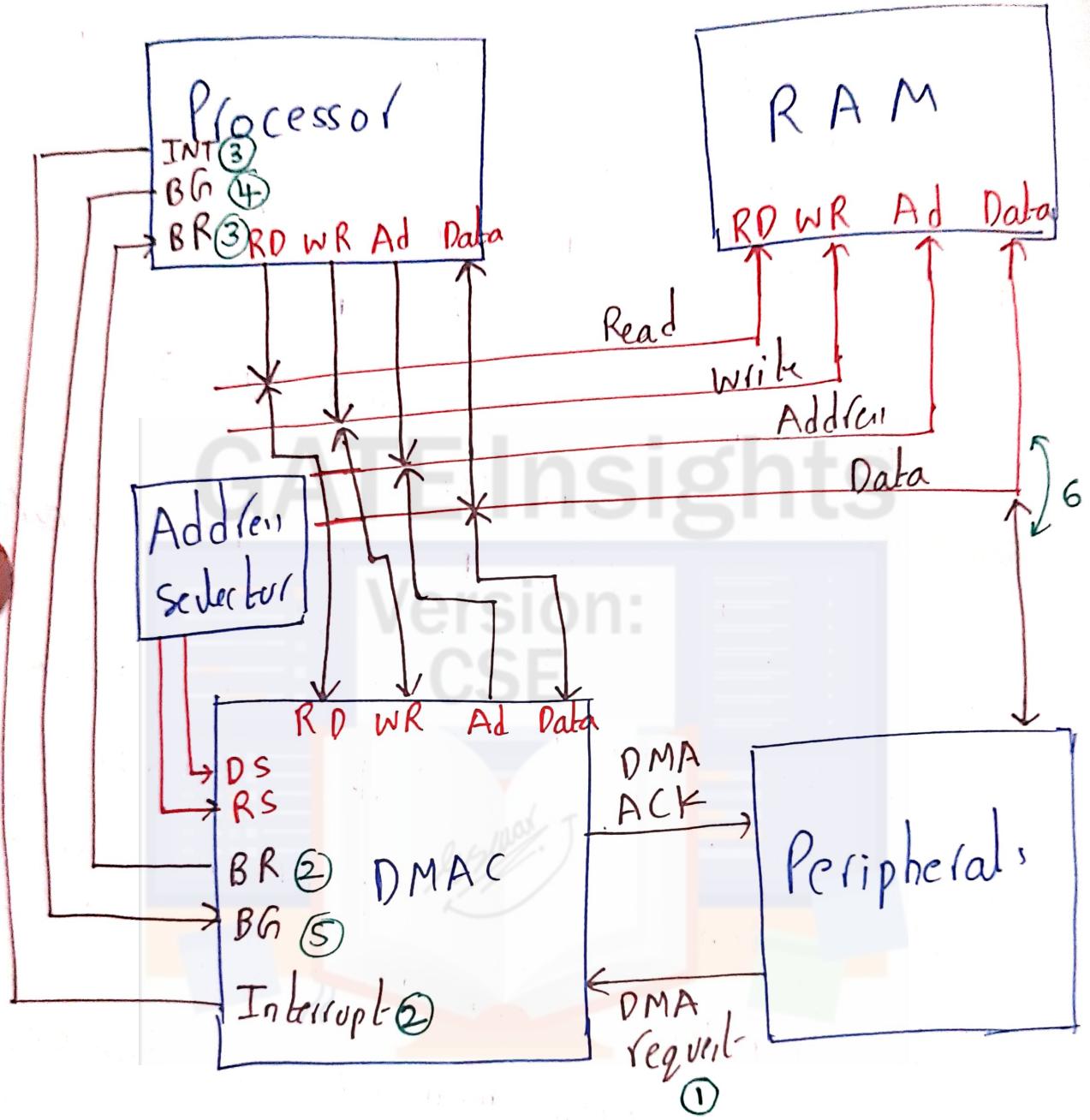
ii) Vectored interrupt: Some interrupt may be vectored where the interrupting device will also tell the address that where this routine is stored in the memory. (Address should be informed).

Adv: very simple, easy to use, easy to understand, relatively fast.

Disadv: here priority fixed & even in case of requirement- we cannot change it.

III) Direct Memory Access (DMA)

When we want to perform I/O operations then the actual source or destination is either I/O device or memory but CPU is placed in between just to manage & control the transfer. Hence DMA is the idea where we use a new device called DMA controller using which CPU allows DMA controller to take control of system buses & perform different data transfer either from device to memory or from memory to device.



BG - Bus Grant

BR - Bus Request

Model of transfer:

Burst mode: When entire I/O transfer is completed & the control comes back to CPU then it is called burst mode.

Cycle stealing mode:- When CPU executes an instruction then normally their could be following phases

IF : Instruction fetch

ID : Instruction Decode

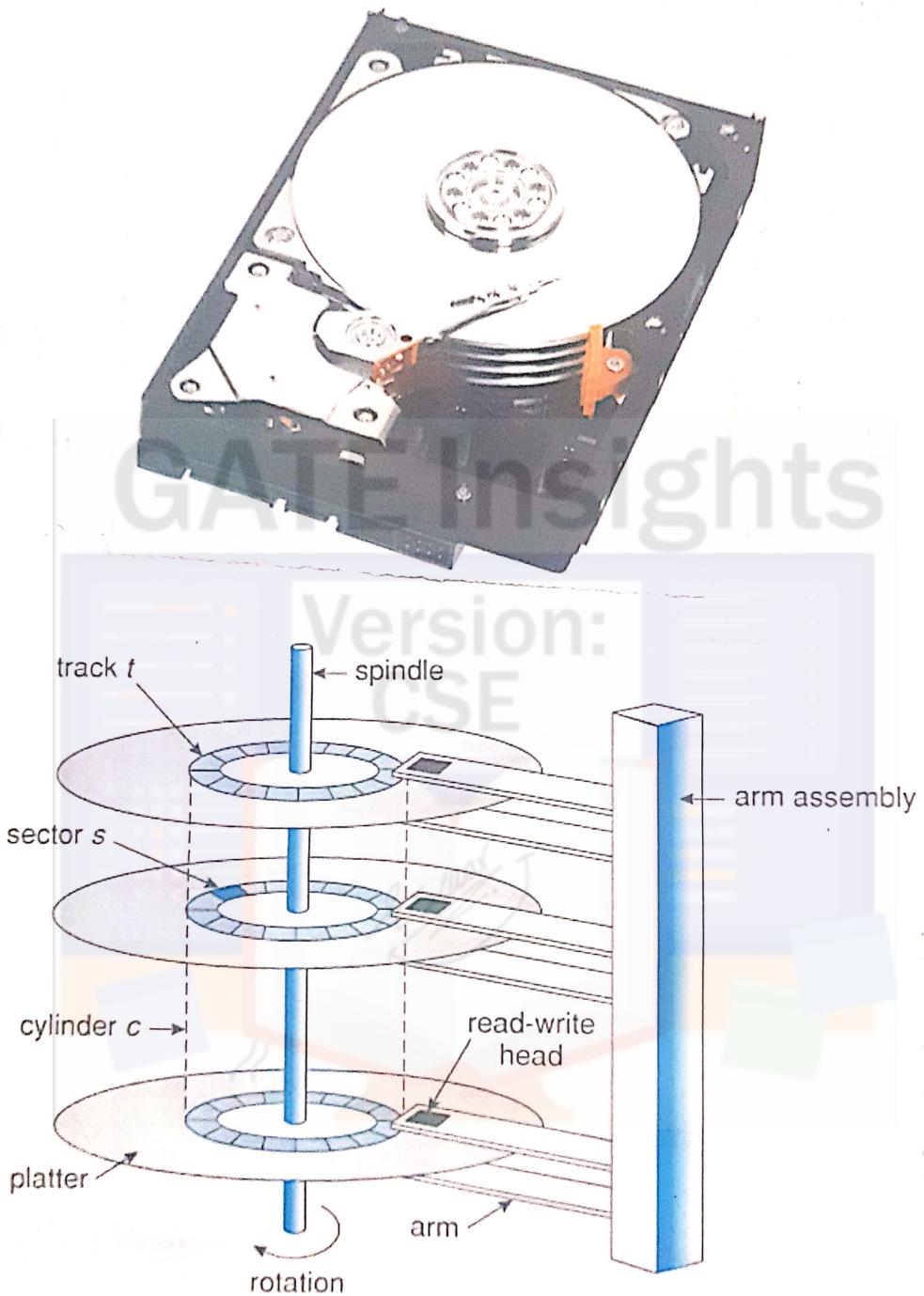
OF : Operand fetch

IX : Instruction Execute

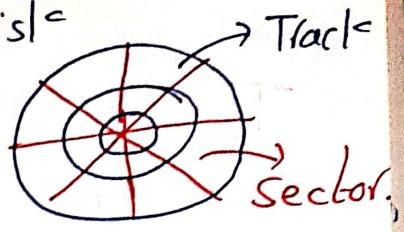
WB : Write back.

Here ID & IX phases CPU don't require system buses & if only in that time control is giving to DMA then it is called cycle stealing.

→ Mass storage structure:



- Spindle holds the disk (palettes) which are stacked one on other.
 - Spindle rotates at high speed. (counter clockwise)
 - Surface of palettes are covered by magnetic material.
 - Both surfaces are coated.
 - Data is written/read with the help of read/write head controlled by hand/arm assembly.
 - moving we can access different tracks.
 - These tracks are further divided into sectors.
- surface no → Track no → sector no



$$\text{Total Transfer Time} = \text{Seek Time} + \text{Rotational Latency} + \text{Transfer Time}$$

Seek Time: It is the time taken by Read/Write head during to reach the correct track (given in question).

Rotational Latency: It is the time taken by Read/Write head during the wait for the correct second.

In general it's a random value 0 to 360° rotation time.

So we take avg.

i.e half of complete time,

Transfer time: It is the time taken by read/write head either to read/write on a disk.

Total time = $\left(\frac{\text{file size}}{\text{track size}} \right) * \text{time taken by computer to complete one revolution}$

ex: disk speed = 360 rpm
track size = 512 bytes.

speed = ?
of data transfer.

(Date 1995)

$$\begin{array}{ccc} 360 & \xrightarrow{\quad} & 1 \text{ min} \\ \frac{360}{60} & \xrightarrow{\quad} & 1 \text{ sec} \\ 6 & \xrightarrow{\quad} & 1 \text{ sec} \\ 1 & \xrightarrow{\quad} & \frac{1}{6} \text{ sec.} \end{array}$$

$$\begin{aligned} \frac{1}{6} \text{ sec} &\rightarrow 512 \text{ B} \\ 1 \text{ sec} &\rightarrow 3072 \text{ B} \\ \therefore & 3072 \text{ B/sec.} \end{aligned}$$

ii) 15000 rpm

data transfer rate: 50×10^6 byte/sec

seek time = $2 * \text{Rotational latency}$.

combined transfer = $10 * \text{disk transfer time}$.

read/write 512 byte sectors

$$15000 \text{ r} \rightarrow 1 \text{ min}$$

$$\frac{1500}{6} \text{ r} \rightarrow 1 \text{ sec}$$

$$250 \text{ r} \rightarrow 1 \text{ sec}$$

$$1 \text{ r} \rightarrow \frac{1}{250} \text{ sec}$$

$$1 \text{ r} \rightarrow 4 \text{ msec}$$

$$RL = \frac{4}{2} = 2$$

$$ST = 2 * 2 = 4$$

disk Transfer time = ~~ST~~

$$50 \times 10^6 \text{ B} \rightarrow 1 \text{ s}$$

$$1 \text{ B} = \frac{1}{50 \times 10^6} \text{ sec}$$

$$\begin{aligned}
 512B &\rightarrow \frac{512}{50 \times 10^6} \text{ sec.} \\
 &\rightarrow \frac{512 \times 10^{-3}}{5 \times 10^6} \text{ msec} \\
 &\Rightarrow 0.01024.
 \end{aligned}$$

$$\begin{aligned}
 T.T.T &= ST + RL + TT + C.T \\
 &= 4 + 2 + 0.01024 + 10 \times \frac{0.01}{0.24}
 \end{aligned}$$

$$\frac{6.1148}{6 \text{ ms}}$$

→ Pipelining:
 If the system has only one processor then almost one instruction can be executed at a time.

If we really want to execute multiple instructions together or concurrently then we must have multiple processors

→ Pipelining is a phenomena or method using which we will be able to run more than one instruction at the same time, on a single processor.

Ex: Assume we have 2 instructions
& each instruction has 5 phases

(5 phases) {

- Instruction Fetch
- Instruction Decode
- Operand Fetch
- Instruction execute
- Instruction store / Write back.

each phase takes 1 cycle.

| Instruction | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|-------------|---|---|---|---|---|---|---|---|---|----|
| Fetch | ✓ | | | | ✓ | | | | | |
| Decode | | ✓ | | | | | ✓ | | | |
| Fetch | | | ✓ | | | | | | ✓ | |
| execute | | | | ✓ | | | | | | ✓ |
| Write | | | | | ✓ | | | | | ✓ |
| Clock | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |

For 2 instructions without pipeline
 it took 10 clocks
 \Rightarrow no. of instruction * no. of phase
 $\Rightarrow (2 * 5) = (10)$

with pipeline:

| Instruction | 1 | 2 | | | | |
|-------------|---|---|---|---|---|---|
| Fetch | ✓ | ✓ | | | | |
| Decode | | ✓ | ✓ | | | |
| Fetch | | | ✓ | ✓ | | |
| execute | | | | ✓ | ✓ | |
| write | | | | | ✓ | ✓ |
| clock | 1 | 2 | 3 | 4 | 5 | 6 |

for 2 instructions with pipeline it took 6 clocks
 \Rightarrow no. of phase + $(\frac{\text{no. of instruction}}{\text{instruction}} - 1)$
 $* \text{time of one clock}$

ex 10 instruction 4 phas
 without pipeline with pipeline

$$10 * 4 \rightarrow 40$$

$$4 + (9 * 1) \\ \rightarrow 13$$

$$\text{Speed up} = \frac{\text{Time without pipeline (T}_{wp})}{\text{Time with pipeline (T}_p)}$$

$$= \frac{40}{13} = 3.076$$

Max Speed up = No. of stage = 4.

$$\text{Efficiency} = \frac{\text{Speed up}}{\text{max speed}} * 100 = 76.92\%$$

1000 instruction 4 phas

$$T_{wp} \Rightarrow 1000 * 4 \\ \rightarrow 4000$$

$$T_p \Rightarrow 4 + 999 \\ \Rightarrow 1003$$

$$\text{Speed up} = \frac{4000}{1003} = 3.988 \quad \text{Max speed up} = 4$$

$$\text{efficiency} = \frac{3.988}{4} * 100 = 99.7\%$$

- As no. of instruction increase
efficiency also increases
Ideal state which we want
is 1 clock per instruction.

Two
P)
if we assume

$$m = \text{no. of stages}$$

$$n = \text{no. of instruc}$$

$$\text{Speed up} = \frac{m * n}{m + (n - 1)}$$

if n value is too big; m can
be neglected

$$\frac{m * n}{n} = m$$

Note that here we assumed each
phase takes 1 cycle but it really
is not.

ex: 5 instructions, 4 phases

| | F | D | E | WB | Total |
|----------------|-------|-----|-----------------------------------|----|--|
| I ₁ | 1 | 2 | 1 | 1 | $\Rightarrow 5$ |
| I ₂ | 1 | 2 | 2 | 1 | $\Rightarrow 6$ |
| I ₃ | 2 | 1 | 3 | 2 | $\Rightarrow 8$ |
| I ₄ | 1 | 3 | 2 | 1 | $\Rightarrow 7$ |
| I ₅ | 1 | 2 | 1 | 2 | $\Rightarrow 6$ |
| | | | | | $T_{WP} = \frac{32}{\cancel{\cancel{}}}$ |
| | D { X | D ✓ | X stall (can't execute) (wait) | | |

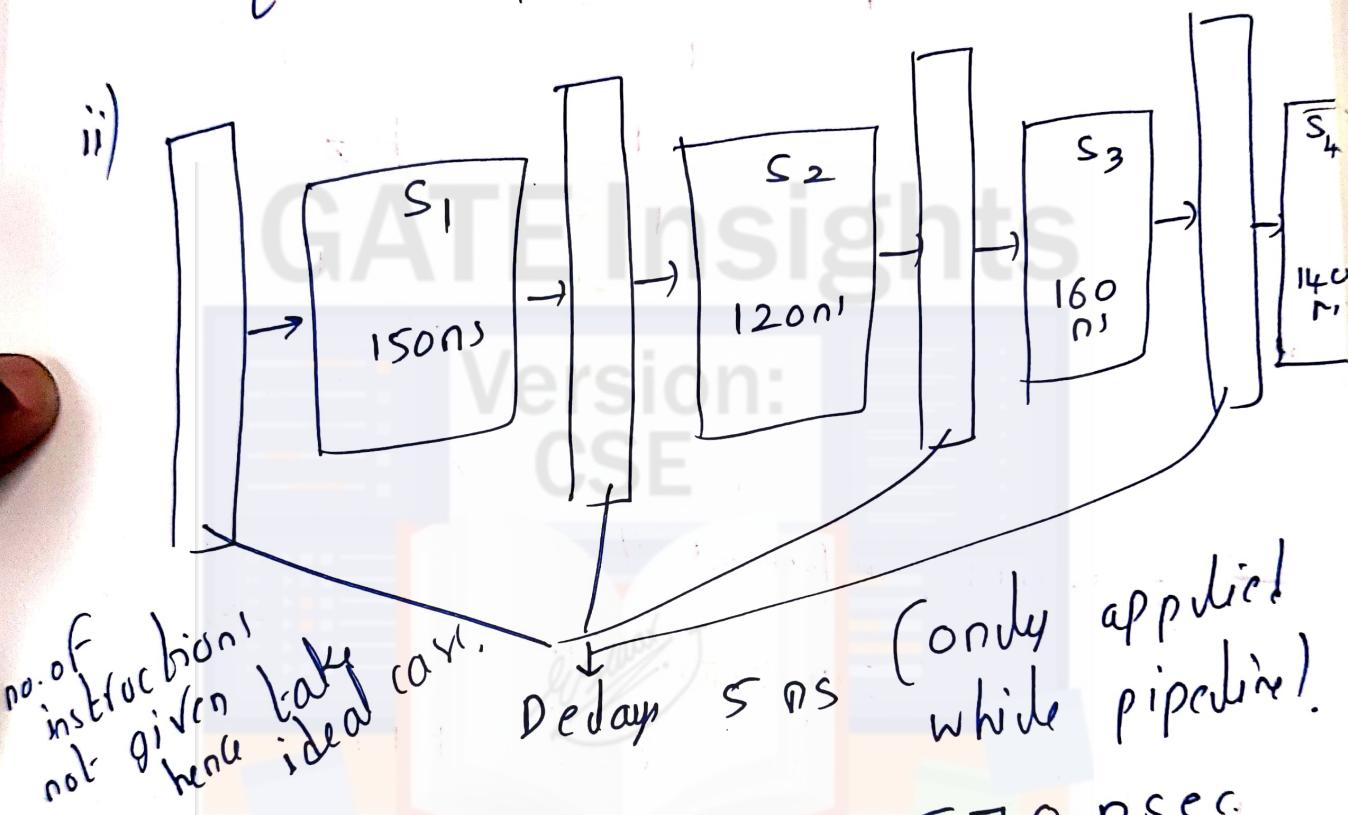
$T_P =$

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|----------------|---|---|---|---|----|---|---|----|---|----|----|----|
| I ₁ | F | D | D | E | WB | | | | | | | |
| I ₂ | | F | X | D | D | E | E | WB | | | | |
| I ₃ | | | F | F | X | D | X | E | E | E | WB | WB |
| I ₄ | | | | X | F | X | D | D | D | X | E | E |
| I ₅ | | | | | X | F | X | X | X | D | D | X |

| | 13 | 14 | 15 | 16 | 17 |
|----------------|----|----|----|----|----|
| I ₁ | | | | | |
| I ₂ | | | | | |
| I ₃ | | | | | |
| I ₄ | WB | | | | |
| I ₅ | E | WR | WB | | |

$$\text{Speed up} \rightarrow \frac{T_{WP}}{T_P} = \frac{32}{15} \\ \Rightarrow 2.17$$

$$\eta = \frac{2.17}{4} \times 100 = 53.33\%$$



$$T_{WP} = 150 + 120 + 160 + 140 = 570 \text{ nsec}$$

$$T_P = \frac{s_1}{155} \quad \frac{s_2}{125} \quad \frac{s_3}{165} \quad \frac{s_4}{145}$$

s added as delay is their

which is max 165 ns
hence it is the T_P

IF frequency asked

$$F = \frac{1}{T}$$

$$F \rightarrow 13 \text{ Hz}$$

$$F = \frac{1}{165 \text{ nsec}}$$

$$F = \frac{1}{165} \times 10^9$$

$$= \frac{1}{165} \text{ GHz}$$

$$\text{or } 6.06 \text{ MHz} \rightarrow 10^6$$

→ Hazards/dependency :-

Sometimes we cannot execute instruction with full efficiency in a pipeline because of certain dependency or hazards.

3 types

i) Structural hazard

ii) Control hazard

iii) Data hazard

hazard

hazard

hazard

I) Structural hazard:

Even by having m stages in pipeline processor and assuming that the combinational circuit of every stage is different many times there will be dependency because of external hardware parts from CPU like system bus, memory etc.. coz of this stall will occur which decreases speed up.

Structural hazard cannot be removed by programming. The only solution is replicate the hardware components (costly).

with pipeline

e.g. 4 phases

need memory access
Fetch & decode
Data read
ALU
Data write.

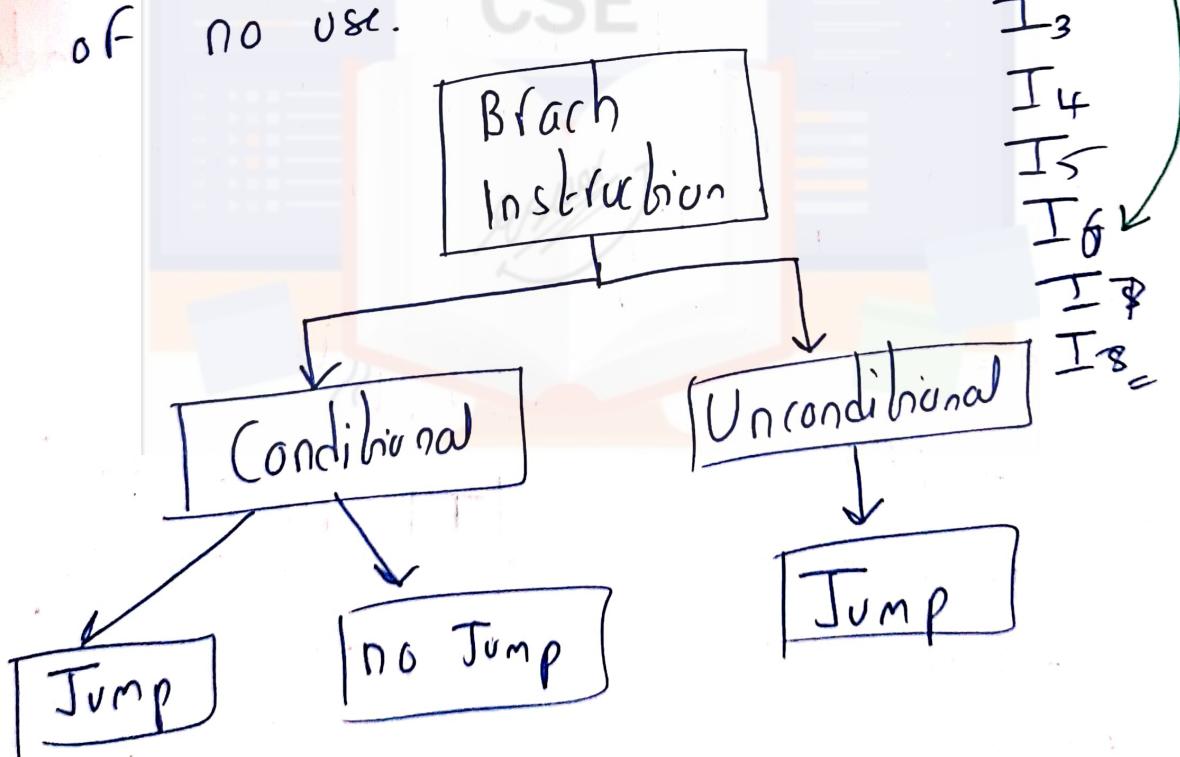
| IF | DR | ALU | DW | |
|----|----|-----|----|--------|
| | X | IF | X | DR ALU |
| 1 | 2 | 3 | 4 | 5 G T |
| | | | | |
| | | | | |

→ per 2 instructions

due to structural hazard //

II) Control Hazards:

IF some instruction I_1 is branch instruction & after executing the entire instruction we understood that there is jump & next instruction to be executed is I_{10} . This time we have partially loaded $I_2 I_3 I_4 \dots$ which is a problem of no use.



ex consider 20% of instructions are branch instruction

assume CPI = 1

now what is CPI = ?

if Branch instruction = 1 stall + 1
= $(0.8 \times 1) + (0.2 \times 2)$,
non cond branch only

$$\Rightarrow 0.8 + 0.4 \Rightarrow 1.2 \geq CPI$$

if BI = 2 stall
 $\Rightarrow (0.8 \times 1) + (0.2 \times 3) = 1.4$

if BI = 3 stall
 $\Rightarrow (0.8 \times 1) + (0.2 \times 4) = 1.6$

Best solution for this is code rearrangement

III) Data Hazards: (Gate 2010)

| ex | Instruction | meanings |
|------------------------|--|---|
| 3 PO \leftarrow M UL | R ₂ , R ₀ , R ₁ | R ₂ \leftarrow R ₀ * R ₁ |
| 6 PO \leftarrow D IV | R ₅ , R ₃ , R ₄ | R ₅ \leftarrow R ₃ / R ₄ |
| 1 PO \leftarrow ADD | R ₂ , R ₅ , R ₂ | R ₂ \leftarrow R ₅ + R ₂ |
| 1 PO \leftarrow SUB | R ₅ , R ₂ , R ₆ | R ₅ \leftarrow R ₂ - R ₆ |

5 stages

- 1 \leftarrow IF - Instruction fetch
- 1 \leftarrow ID - Instruction decode
- 1 \leftarrow OF - Operand fetch
- (Based on op) PO - Perform operation
- 1 \leftarrow WO - Write Operand.

(without forwarding) (from memory)

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| IF | ID | OF | PO | PO | PO | WO | | | | | | | | | |
| . | IF | ID | OF | X | X | PO | WO | | |
| . | IF | ID | OF | X | X | X | X | X | X | X | X | X | PO | PO | WO |
| . | IF | ID | OF | X | X | X | X | X | X | X | X | X | X | X | X |

17 18 19 20 21

OF PO WO

Note as their ii dependency

$$R_2 \leftarrow R_5 + R_2$$

Both result will be ready after 13

$$R_5 \leftarrow R_2 - R_6$$

Both result will be ready at 15

without Forwarding.

This approach is without Forwarding.

(with Forwarding) (Result taken a head)
→ (From ALU / CPU)

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| IF | ID | OF | PO | PO | PO | WO | | | | | | | | | |
| IF | ID | OF | X | X | Po | Po | Po | Po | Po | Po | WO | | | | |
| IF | ID | ← | — | — | — | OF | | | | | | | | | |
| IF | ID | ← | — | — | — | OF | | | | | | | | | |

Forwarding

as a simple logic

$I_2 \Rightarrow R_5 = R_3 / R_4$ (Result will be ready at 13 clock).

hence in 14 clock we can perform

I_3 by taking data ALU / CPU init'd
of memory (This is Forwarding logic).

Data hazard occur when instructions that exhibit data dependence, modify data in different stages of pipeline

Then all Bernstein condition

RAW (Read after Write) (Flow / True data dependency)

WAR (Write after read) (Anti Data dependency)

WAW (Write after write) (Output data dependency)

$$I : R_2 \leftarrow R_1 + R_3$$

$$J : R_4 \leftarrow R_2 + R_3$$

J is trying to read R_2 before I have written it (RAW).

$$I : R_2 \leftarrow R_1 + R_3$$

$$J : R_3 \leftarrow R_4 + R_5$$

J trying to write R_3 before I read it (WAR).

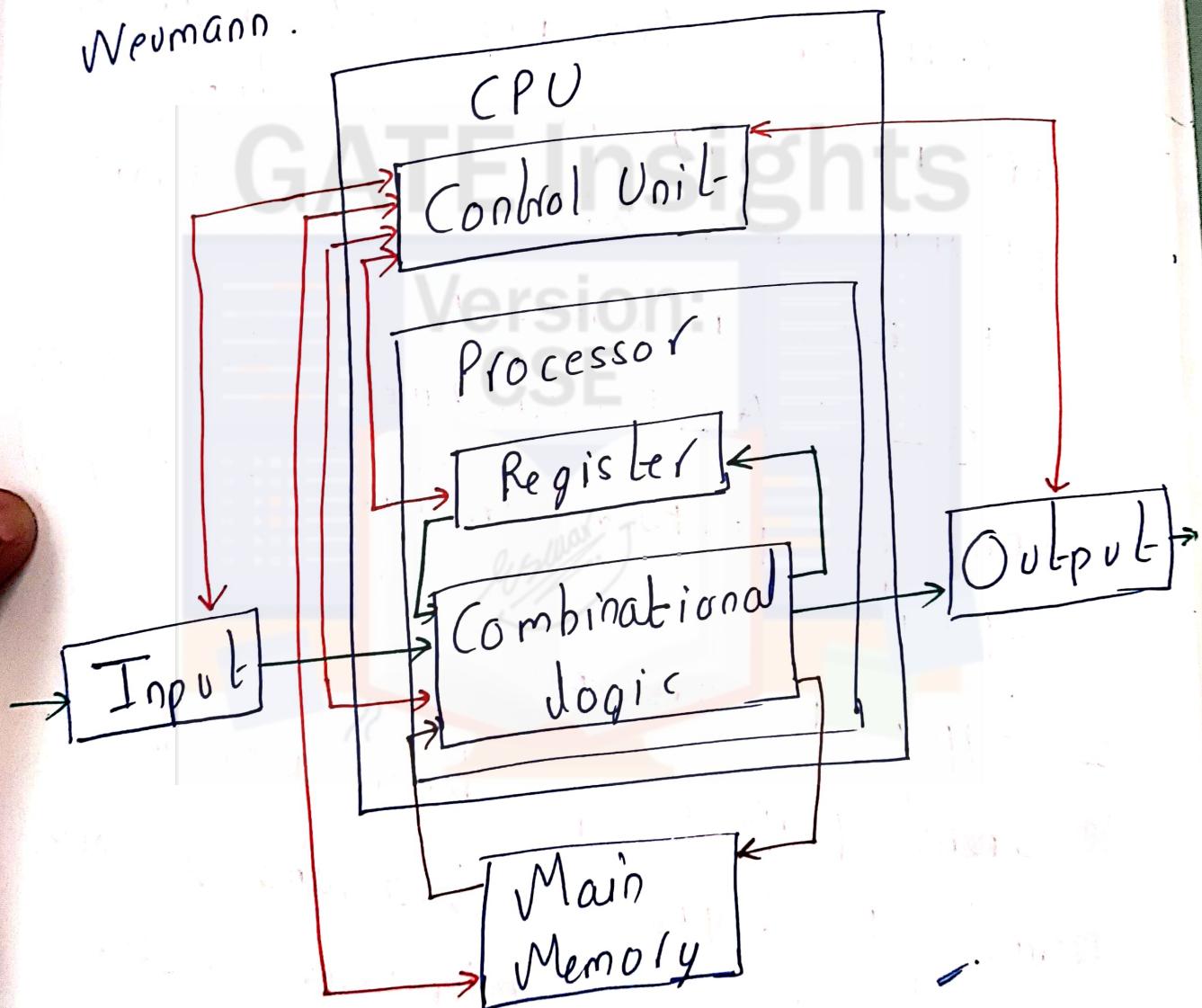
$$I : R_2 \leftarrow R_1 + R_3$$

$$J : R_2 \leftarrow R_4 + R_5$$

J trying to write R_2 before I (WAW)

→ Von neumann architecture :-

The Von Neuman architecture is a computer architecture based on a 1945 description by John von Neumann.



3 main components:

Control Unit: Which generates control signal (mask/generator) to control every other part of the CPU.

Registers Few important registers are in every processor like program counter (which consists address of next instruction).

ALU: ALU (Arithmetic logic unit) which can perform arithmetic & logical operations.

2 types of Computer

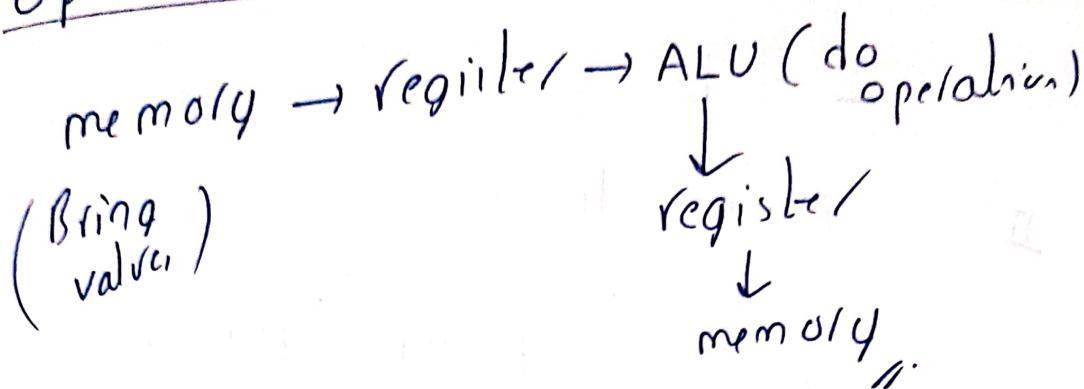
General purpose

Normal PC, Laptop, mobile etc.

Dedicated device

Washing machine, microwave, ~~fridge~~, ATM machine,

operation performing flow:



→ Instruction Format :-

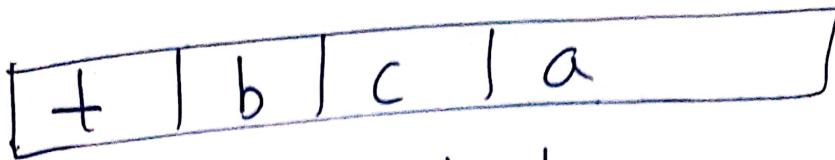
5 types
I) 4 Address Instruction
4 address & 1 op code/mode.

| Mode/Opcode | Op1end1 | Op1end2 | Result | Next instruction |
|----------------|---------|---------|--------|------------------|
| +,-,* / 2 bits | 40 bit | 40 bit | 4 bit | 4 bit |

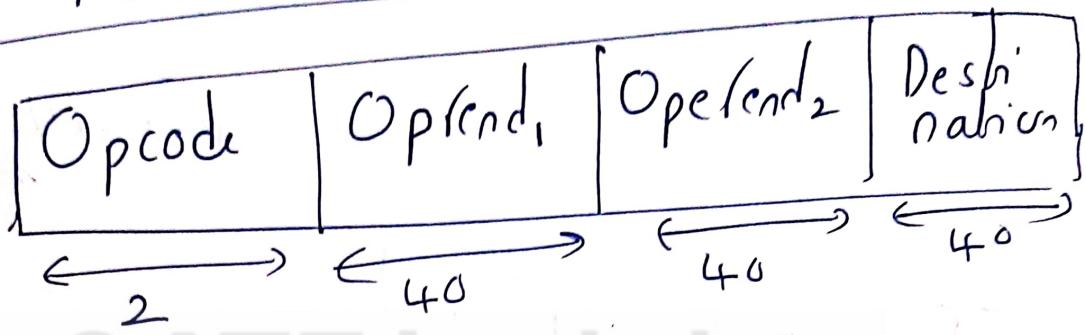
Adv: → Very simple, easy to use easy to understand.

Disadv: → Very lengthy & occupies a lot of space.

$$a = b + c$$



II) 3 Address Instruction:



$$x = (A + B) * (C + D)$$

$$R_1 \leftarrow M[A] + M[B]$$

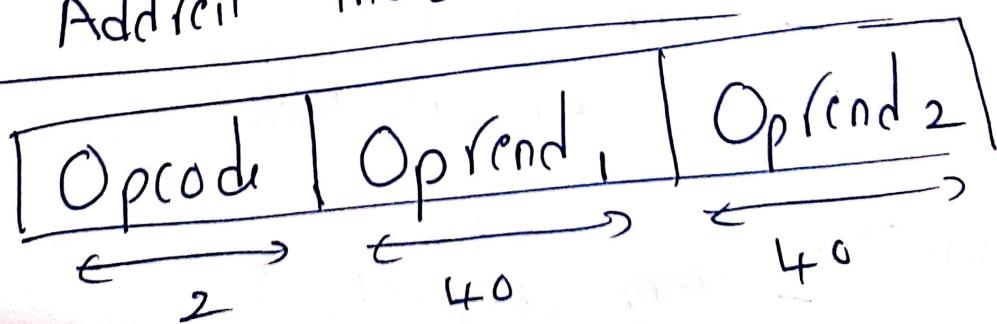
$$R_2 \leftarrow M[C] + M[D]$$

$$M[X] \leftarrow R_1 * R_2$$

Adv: Better than 4 add/c.,

Disadv: still lengthy.

III) 2 Address Instruction:



$$X = (A + B) * (C + D)$$

$$R_1 \leftarrow M[A]$$

$$R_1 \leftarrow R_1 + M[B]$$

$$R_2 \leftarrow M[C]$$

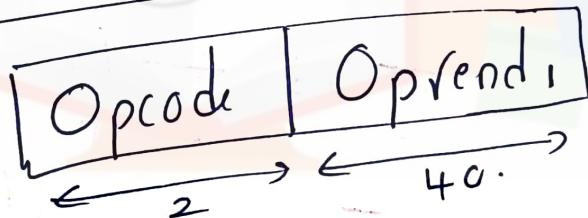
$$R_2 \leftarrow R_2 + M[D]$$

$$R_1 \leftarrow R_1 * R_2$$

$$M[X] \leftarrow R_1 //$$

Adv less length, more efficient, no. of registers low.
disadv - code optimization is relatively difficult.

IV) Address Instruction:



AC (Accumulator is used).

$$X = (A + B) * (C + D)$$

$$AC \leftarrow M[A]$$

$$AC \leftarrow AC + M[B]$$

$$M[T] \leftarrow AC$$

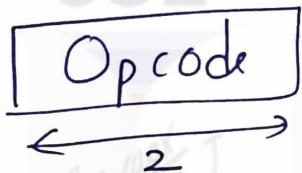
$$AC \leftarrow MC$$

$$AC \leftarrow AC + M[D]$$

$$AC \leftarrow AC * M[T]$$

$$M[X] \leftarrow AC$$

II) O Address Instruction:



stack is used here.

$$TOS \leftarrow A$$

$$TOS \leftarrow B$$

$$TOS \leftarrow (A + B)$$

$$TOS \leftarrow C$$

$$TOS \leftarrow D$$

$$TOS \leftarrow (C + D)$$

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

$$M[X] \leftarrow TOS$$

→ Addressing Mode:

It specifies the different ways possible in which reference to the operand can be made.

Effective address: It is the final address of the location where the operand is stored.

many types
I) Immediate mode addressing :-

It means the operand is itself part of instruction

MVI B, 30 H

~~location~~ 30 data move
to B,

Adv: Can be used for constants, where values are already known

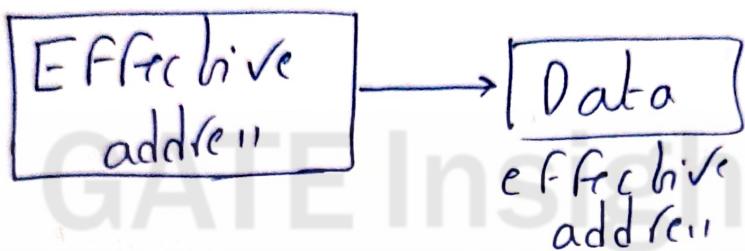
Fast & easy.

Disadv: Cannot access full space as its immediate it's size is limited.

II) Direct mode addressing:

It means instruction contains address of the memory where variable or data is present.

Instruction



Adv:

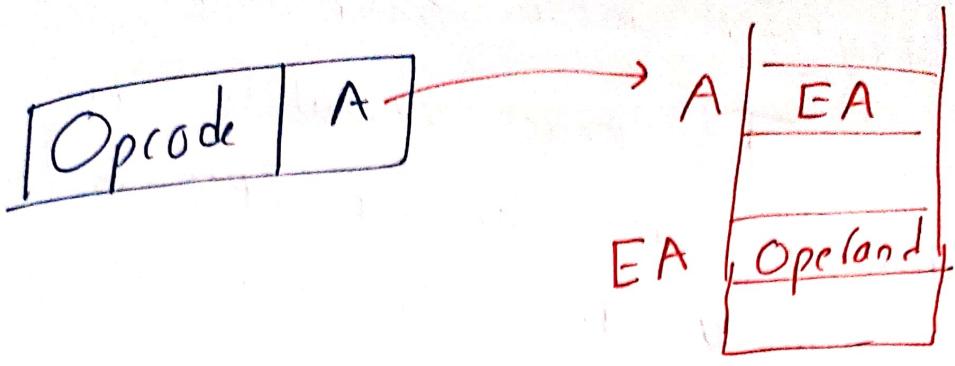
- Best for unknown variables
- No restricted range.

Dis Adv:

- relatively slow compared to immediate mode.

III) Indirect mode Addressing:

Here in the instruction we store the address where the effective address is stored using which we can access the actual data.



- Adv:
- No limitation on no. of variable/sizes
 - Implementation of pointer are feasible

- Disadv:
- Relatively slow.

IV) Implicit mode addressing:-

This mode is mostly used with 0-address instruction which means we only tell the operation system know location
e.g. AC (Accumulator)

e.g. increment, HLT.

V) Register mode addressing:

It means variable are stored in register of CPU instead of memory. In instruction we give the register no.

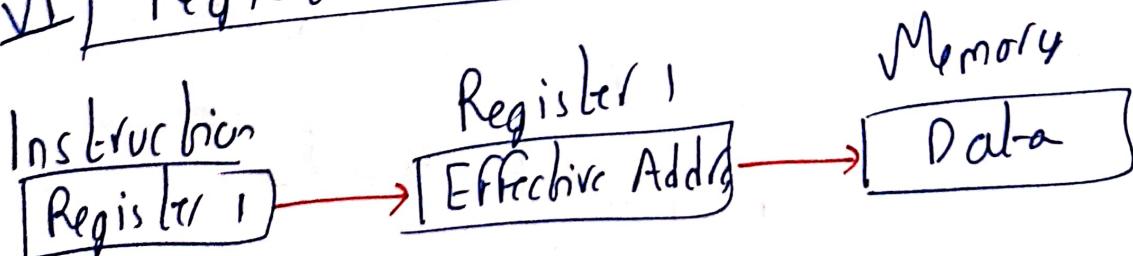


Adv:

- Low no. of bit (as min no. of register exist)
- extremely fast

Disadv: → few variables can only be stored.

VI) Register indirect mode:



III) Base register (OFF set) mode

Here we save the starting of the program address in a register called base register.

$$\text{effective address} = \text{base address} + \text{offset (instruction)}.$$

Now the advantage is even if we try to shift process in the memory, we only need to change the content of the base register.

III) Index addressing mode :- (large array access).

Base address of the array ~~(in)~~ in the instruction & the index which we want to access will be there in the register.

$$\text{location} = \text{Base} + \text{Index}.$$

Reduced Instruction Set Computer (RISC)

- few instructions
 - few address modes
 - Faster (hence hardwired)
 - all operations done within the registers of the CPU.
 - Memory access limited to load & store
 - Fixed length easily decoded instruction.
- ↳ Used commonly & powerful

Complex Instruction Set Computer (CISC)

- Large number of instructions typically 100 to 250 instructions.
- Some instructions that perform specialized tasks & are infrequently used.
- a large variety of addressing modes typically from 5 to 20 different modes.