

INPUT - OUTPUT INTERFACE

→ The Input - Output interface allows transferring information between external input / output devices (i.e., peripherals) and processors.

Interface

Why Input - Output Interface

→ We require the input - output interface b/w many differences exist b/w each peripheral and the central computer while transferring data.

→ Some significant differences b/w the peripheral and CPU are:-

① The nature of the CPU is electronic, and that of the peripheral devices is electro mechanical and electromagnetic. So, there are many differences in the mode of operation of both CPU and peripheral devices.

② We have a synchronization mechanism because the data transfer rate is slower in peripheral devices than CPU.

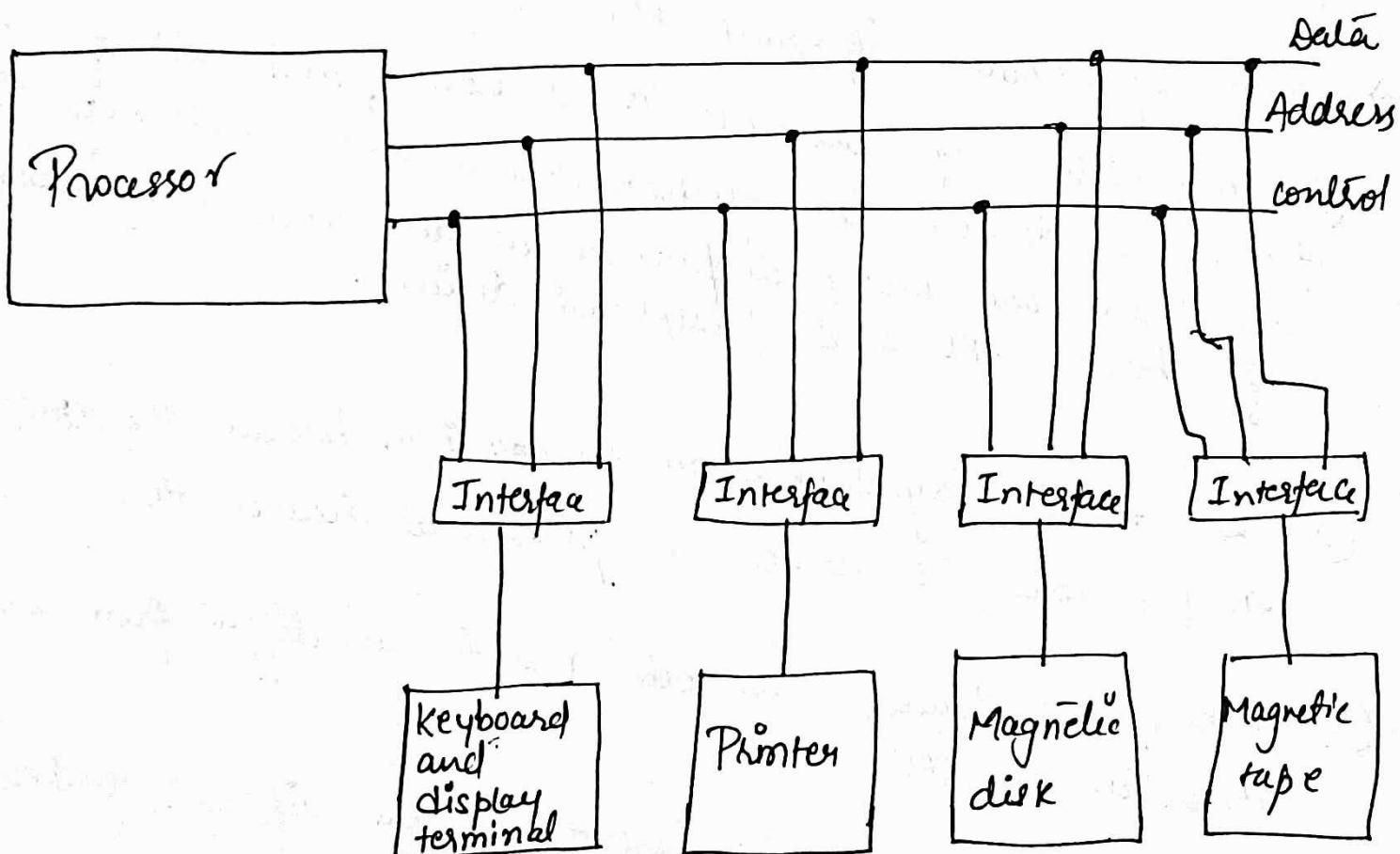
③ In peripheral devices, data code and formats differ from the format in the CPU.

④ The operating modes of peripheral devices are different, and each can be controlled don't to disturb the operations of any other peripheral devices connected to the CPU.

Note :-

- To resolve these differences, computer systems include special hardware components between the CPU and peripherals to supervise and synchronize all the input and output transfers. These components are called interface units because they interface b/w the processor bus and the peripheral devices.
- In addition, each device may have its own controller that supervises the operations of the particular mechanism in the peripheral.

I/O BUS AND INTERFACE MODULES



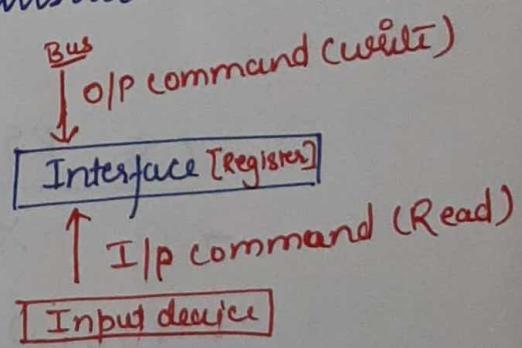
Connection of I/O bus to input output devices

- A typical communication line b/w the processor and several peripherals is shown in figure.
- The I/O bus consists of data lines, address lines and control lines.
- The magnetic disk, printer, and terminal are employed in practically any general-purpose computer.
- Each peripheral device has associated with it an interface unit.
- Each interface unit decodes the address and control received from the I/O bus, interprets them for the peripheral, and provide signals for the peripheral controller. It also synchronizes the data flow and supervises the transfer between the peripheral and processor.
- Each Peripheral has its own controller that operates the particular electromechanical devices.

I/O COMMANDS

I/O command commands are the instructions executed in the interface & its attached to peripheral devices.

- ① Control Command: used to activate the peripherals & to tell it what to do.
- ② Status Command: used to test various status conditions in the interface & the peripherals. like whether the peripherals are busy, or free or error.
- ③ Output data (I/O write): It causes the interface to respond by transferring data from bus into one of its registers. or cause the I/O module to take an item of data from the data bus & subsequently transmit that data item to the peripheral.

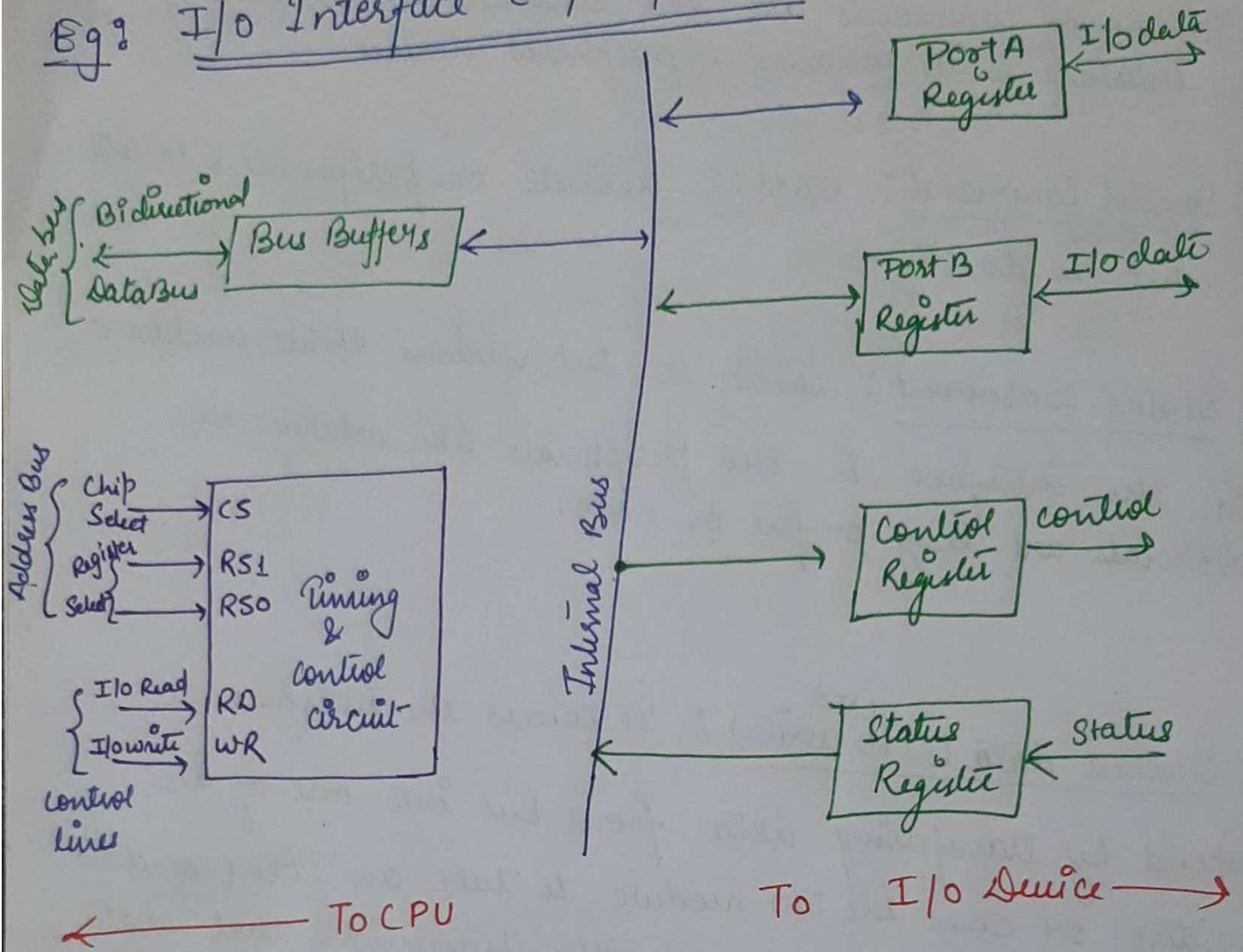


- ④ Input data (I/O read):

Interface receives an item of data from the peripheral & places it in its buffer register. or causes the I/O module to obtain an item of data from the peripheral & place it in an internal buffer.

(internal buffer register / data register)

Eg: I/O Interface (I/O Module)



CS	RSL	RSO	Register Selected
0	X	X	None: data bus in high-impedance
1	0	0	PORT A register
1	0	1	PORT B register
1	1	0	control register
1	1	1	status register

- I/O device transfers data or communicate with CPU with the help of data bus. And it consists of its own Bus Buffer.
- Here, CPU has another unit called Timing & control unit
Timing and control circuit consists of
① chip select (CS).

Address Bus :- Selects the interface unit through the chip
Select (CS) [enable] and the two registers (RS₁, RS₀)
select input. [among four register from I/O device part]

- I/O Data :- Transfer of data to the I/O device or from the I/O device happens either with the help of port A & port B
- Interface unit may work as a output device or Input device depends on the situation

Ex: I/O device means like magnetic tape. which

works as a I/P as well as O/P

In this case Both the port registers will use one for I/P & another for O/P.

If it uses only for I/P means if these ports are used as I/P data
If it uses only for O/P means if these ports are used as O/P data

Control Register :- It receives control information from the CPU by writing the appropriate bits Y₀ to the control register.

Status Register : It stores status conditions like error related I/O devices

I/O Mapping (I/O Addressing / I/O Interfacing Techniques)

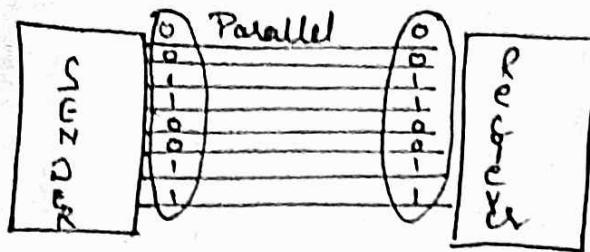
- When computers use one common-bus to transfer information between Memory or I/O & CPU. Then
- There are 2 ways that computer programs communicate with I/O devices. or 2 techniques for addressing the I/O device by the CPU.
 - 1) Isolated I/O (I/O-mapped I/O or Port mapped I/O)
 - 2) Memory Mapped I/O

Data Transfer
(Data Transmission)

Parallel
(Synchronous)

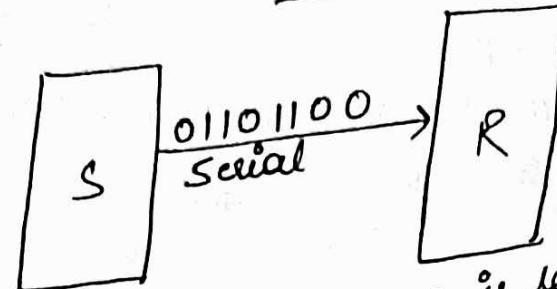
Serial
Synchronous Asynchronous

Parallel Transmission



- Each bit has its own path & Total message is transmitted at the same time.
- n-bit must be transmitted at the same time.
- n-bit must be transmitted through n-separate wires.
- Faster but requires many wires
- Used for short distance where speed is important
- Eg: CPU & memory (Bus)
CPU & Printer
- Parallel transmissions are always synchronous because sender and receiver clock is synchronized

Serial Transmission

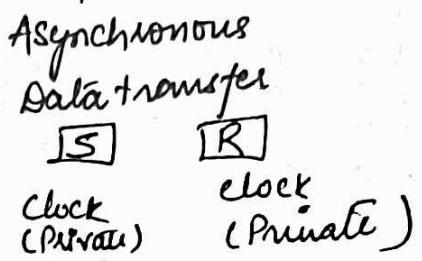
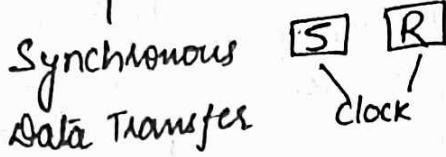


- Each bit in the message is sent in sequence one at a time.
- n-bit is transmitted through one wire.
- Slower but less expensive since requires only pair of wires.
- Used for long distance
Eg: CPU & I/O device
computer to computer

→ Serial transmission has 2 types

- ① Synchronous
- ② Asynchronous

Serial Transmission



Synchronous

- Two units share a common clock
- Data transfer b/w Sender & Receiver is synchronized with same clock pulse
- Used b/w devices their matches in speed.
- Bits are transmitted continuously to keep the clock frequency synchronise in both units.
- Synchronous means "at the same time" (due to common clock)
- Fast
- Costly

Asynchronous

- Two units are independent & each have it own private clock.
- Data transfer b/w Sender & Receiver is not synchronized with same clock pulse.
- Used b/w devices that do not matches in speed.
- Bits are sent only when it is available & line remains idle when there is no information to be transmitted.
- Asynchronous means "at the regular interval" (due to no common clock)
- Slow
- Economical

✓ Asynchronous Data transfer

→ It is used when speed of I/O devices don't match with processor & Timing characteristics of I/O devices is not predictable.

Definition

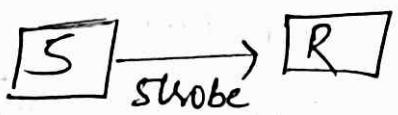
Asynchronous data transfer between two independent units requires control signals to be transmitted between communicating units to indicate the time at which data is being transmitted.



Asynchronous Data Transfer has 2 types of methods

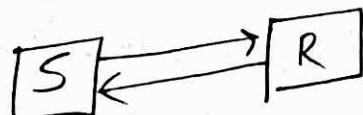
① Strobe control

- a strobe pulse is supplied by one of the units to indicate to the other unit when the transfer has to occur.
- After receiving the data at receiver end there is no acknowledgement is given back to sender that data is received to resolve this another method is introduced which is called handshaking method.



② Handshaking method:

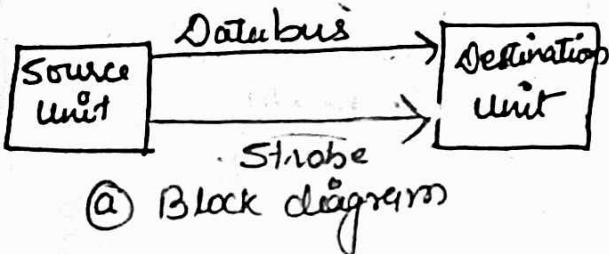
- A control signal is accompanied with each other data item being transferred transmitted to indicate the presence of data.
- The Receiving unit responds with another control signal to acknowledge the receipt of data.



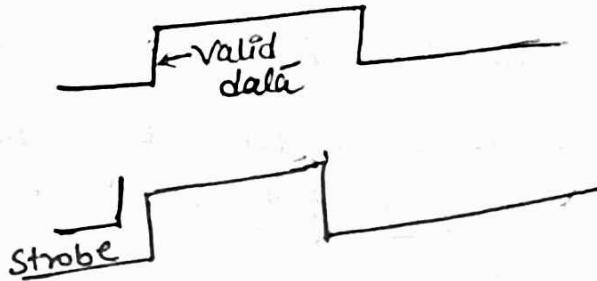
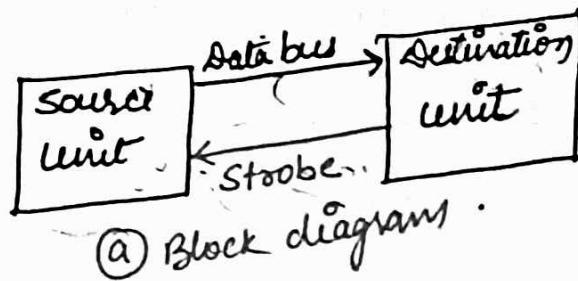
① "Strobe Control" Asynchronous Data Transfer

- it employs a single control line to time each transfer.
- the strobe may be activated by either source or destination unit

Source-initiated strobe for data transfer



Destination-initiated strobe for data transfer



(b) Timing Diagrams

Strobe is a single line that informs a destination unit when a valid data is available in the bus.

Here data is stored first into data bus then strobe pulse is sent. This will enable till the data is not send completely.

(b) Timing diagram

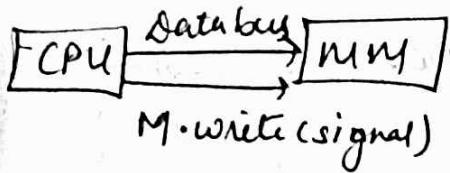
Destination sends a strobe pulse to source for informing that please provide for me

In that case source will store data into data bus & send to the destination.

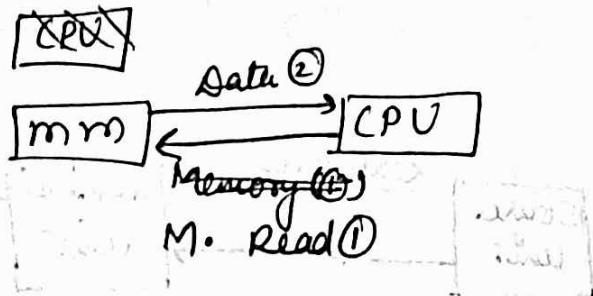
After that strobe pulse is disable.

Eg:-

memory write
control signal from
CPU to memory
unit.



Eg:- memory read control
signal from CPU to
main memory



Disadvantage of Strobe control

- Source unit that initiates the transfers has no way of knowing whether the destination unit has actually received the data items that was placed in the bus.
- Similarly, Destination unit that initiates the transfers has no way of knowing whether the source unit has actually placed the data on the bus or not.
i.e., No Acknowledgement.

Note

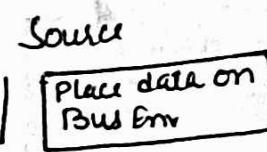
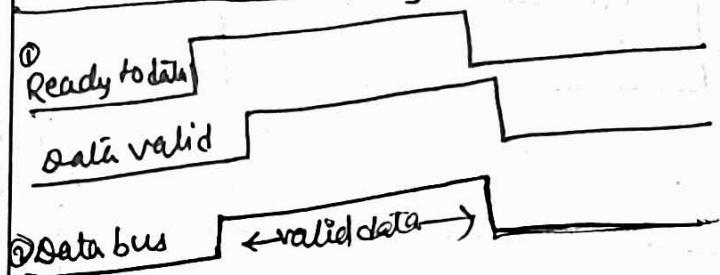
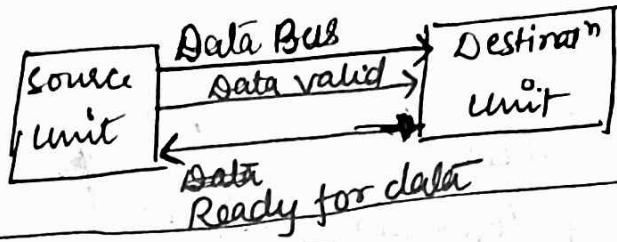
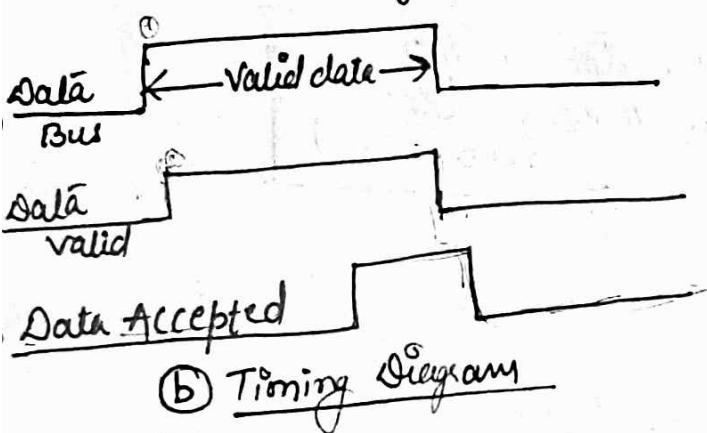
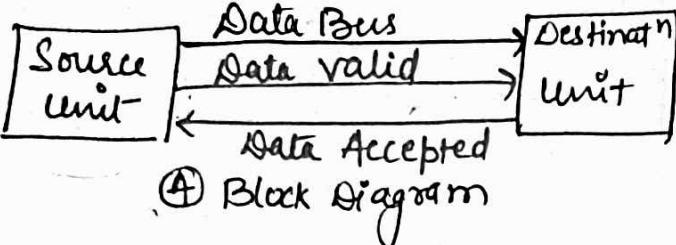


② "Handshaking Method"

for a synchronous data transfer

- Solves the problem strobe method by introducing a second control signal that provides a reply to the unit that initiate the transfer.
- Strobe control + Acknowledgement signal

Source-initiated Handshaking



Advantages

- It provide a high-degree of flexibility & reliability.
- If one unit is faulty, the data transfer will not be completed such an error can be detected by time out mechanism which produces an alarm if the data transfer is not completed within a predetermined time.

Source-initiated Handshaking

Source

Destination

Place data on Bus
Enable data valid

Accept data from Bus
Enable data accepted.

Disable data valid
Invalidate data on Bus

Disable data accepted.
Ready to accept data.
(Initial state)

Destination-initiated Handshaking

Source

Destination

Place data on Bus
Enable data valid

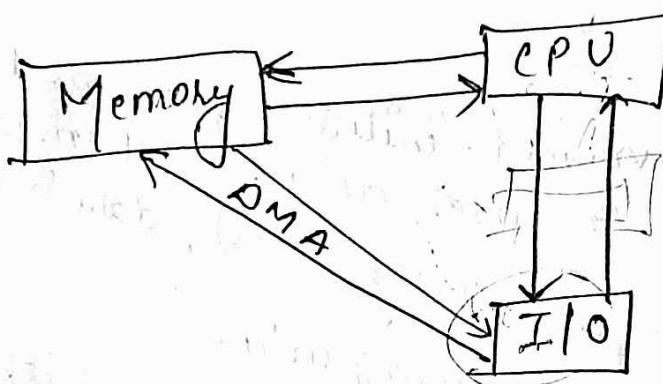
Ready to accept data
Enable ready for data

Disable data valid
Invalidate data on Bus
(Initial state)

Accept data from Bus.
Disable ready for data

Modes of transfer

- Binary information received from an external device is usually stored in memory for later processing.
- Information transferred from the central processing computer into an external device originates in the memory unit.
- The CPU merely executes the I/O instructions and may accept the data temporarily, but the ultimate source or destination is the memory unit.



→ Some nodes use the CPU as an intermediate path; others transfer the data directly to and from the memory unit.

→ Data transfer to and from peripherals may be handled in one of three possible modes:

① Programmed I/O

② Interrupt-initiated I/O

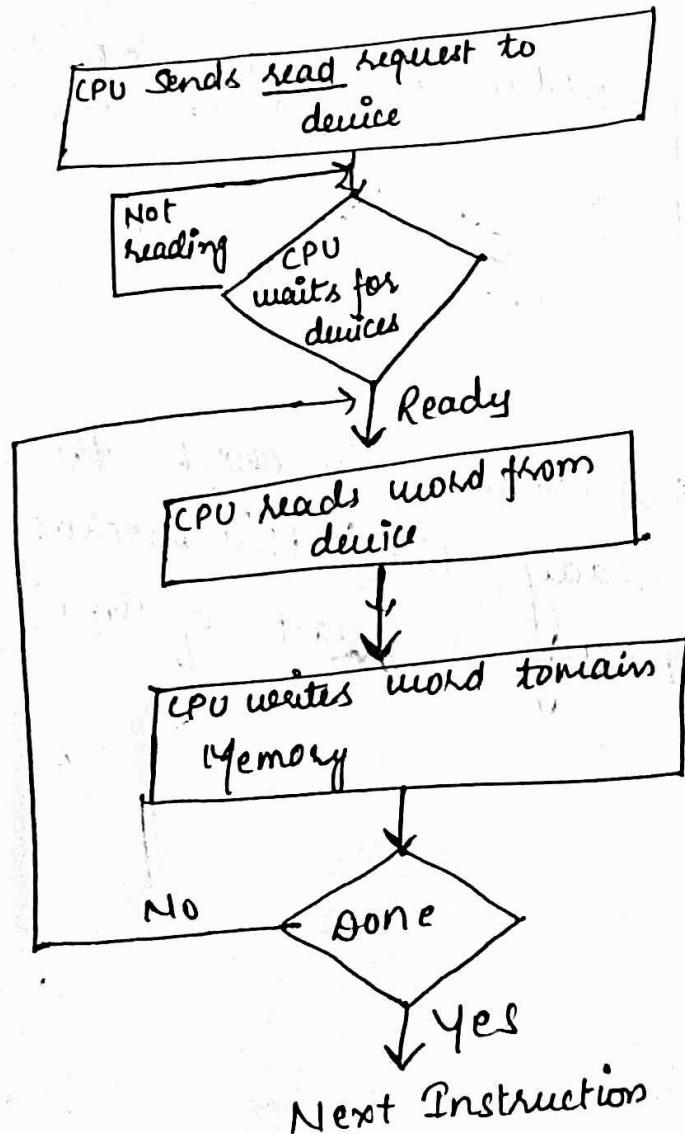
③ Direct memory access ~~acc.~~ (DMA)

Note: In programmed-initiated I/O and interrupt-initiated I/O, CPU is used as an intermediate path while in DMA (direct memory access), data is exchanged directly from the memory unit.

In DMA there is a channel between Input-output devices and Memory through which data is transferred.

Programmed I/O

- Programmed I/O operations are the result of I/O instructions written in the computer program.
- Each data transfer is initiated by an instruction in the program.
- Usually, the transfer is to and from a CPU register and peripheral.
- Transferring data under program control requires constant monitoring of the peripheral by the CPU.



PROGRAMMED I/O FOR INPUT OUTPUT DATA TRANSFER

→ In programmed I/O, CPU makes a request & then CPU stays in ~~in program~~ loop (Polling) until the I/O device indicates that it is ready for data transfer.
[the I/O devices takes no further actions to alert the CPU (i.e., it does not interrupt the CPU)]

Disadvantage:

Time consuming process since it keeps CPU busy needlessly.

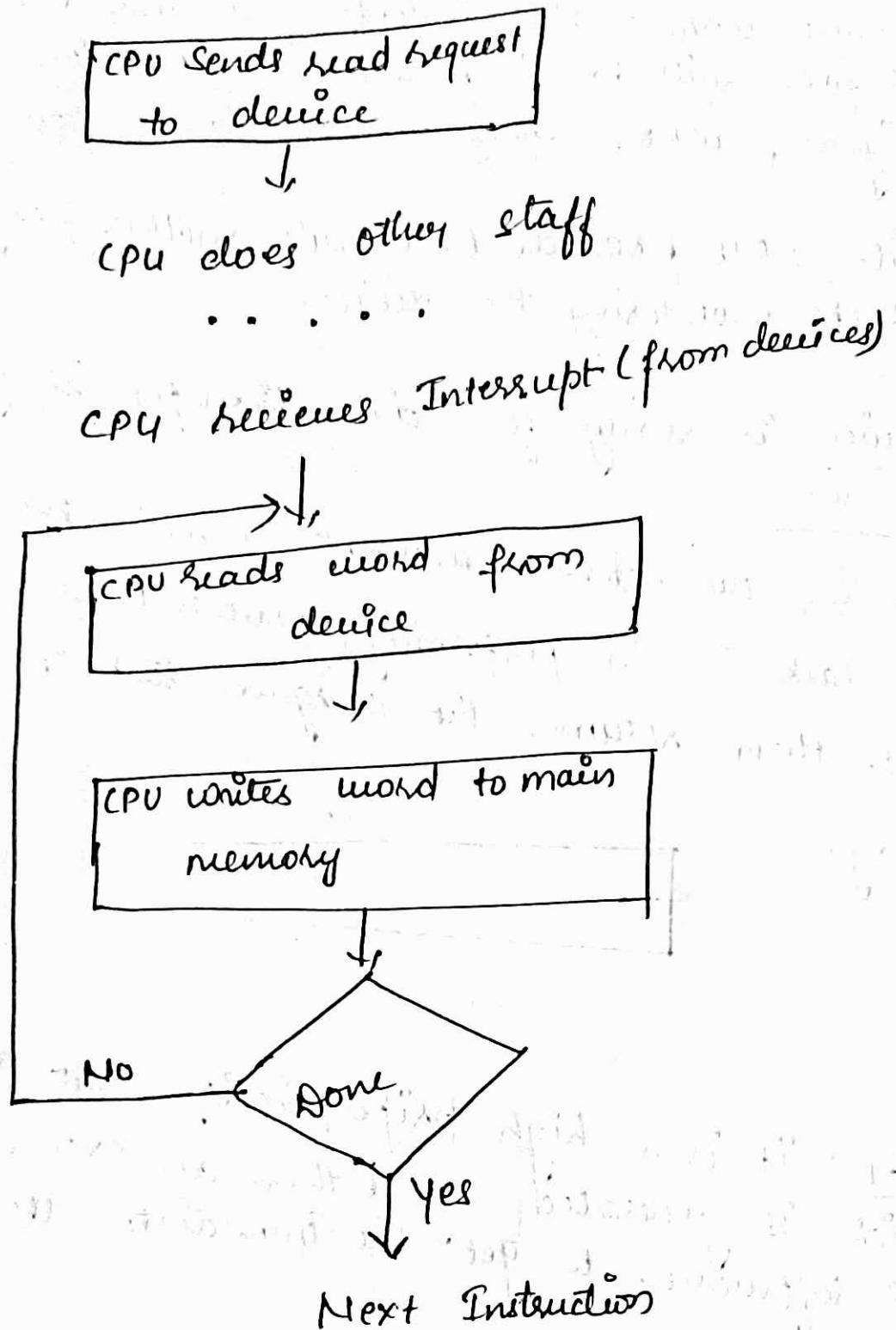
* To avoid this problem interrupt facility can be used

Polling:- When CPU continuously check the devices to see if it is ready, then this method is called polling [which is totally a waste of time]

Interrupt-Initiated I/O

- In Interrupt-initiated I/O, instead of continuous monitoring of CPU, interface will be informed to issue an interrupt request signal, when data are available from the device.
- Meanwhile, CPU proceeds to execute another program & interface keeps monitoring the device.
- When device is ready for data transfer it generates interrupt request.
- Upon detecting the external interrupt signal, the CPU stops the task it is performing, processes the I/O data transfer & then resumes the original task it was performing.

Interrupt: it is a high priority ~~int~~ signal which is generated by either an external device or some software to get the immediate attention of the CPU.



Interrupt Initiated I/O For Input Output

Data transfer

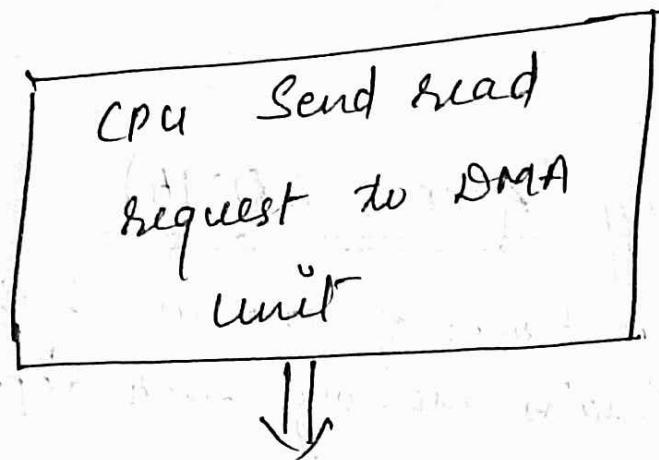
Advantages

Using interrupts the wait period is ideally eliminated.

Direct Memory Access (DMA)

- To transfer large blocks of data at high-speed between external device & main memory DMA approach is used often.
- DMA allows data transfer directly between I/O device & main memory with minimal interventions of CPU.
- DMA means CPU grants I/O Interface authority to read from or write to memory without its involvement.
- DMA itself controls data transfer between main memory & I/O device.
- CPU is only involved in beginning beginning & end of the transfer & interrupt only after entire block has been transferred.
 - * CPU asks the DMA controller to transfer data between a device & main memory & then CPU proceeds to execute other tasks.
 - * The DMA controller issues a request to the right I/O device, waits & manages data transfer between the device & main memory.

* When data transfer is finished the DMA controller
interrupts the CPU.



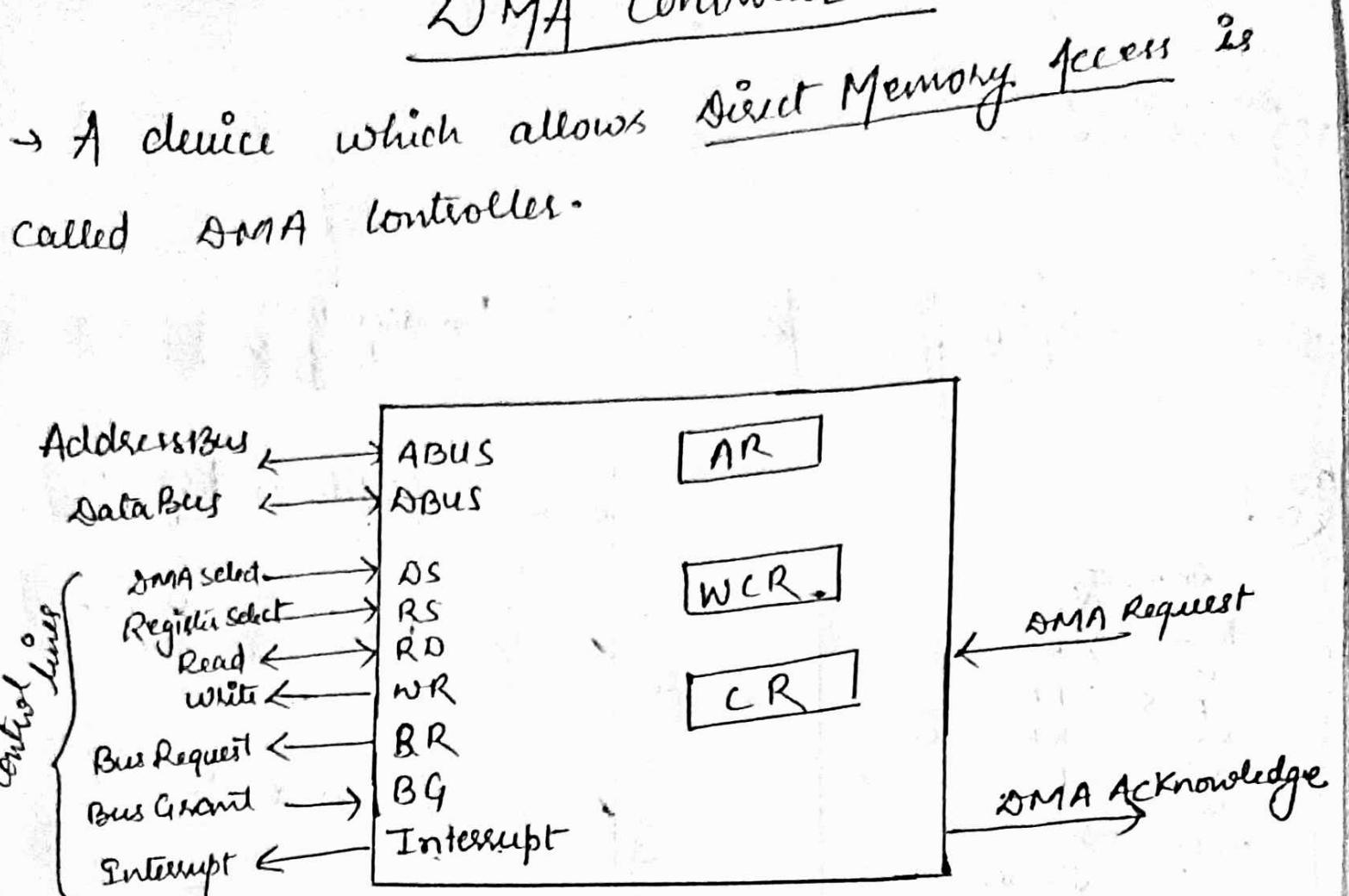
CPU does other stuff

CPU receives DMA
interrupt (from I/O device)

Next instruction

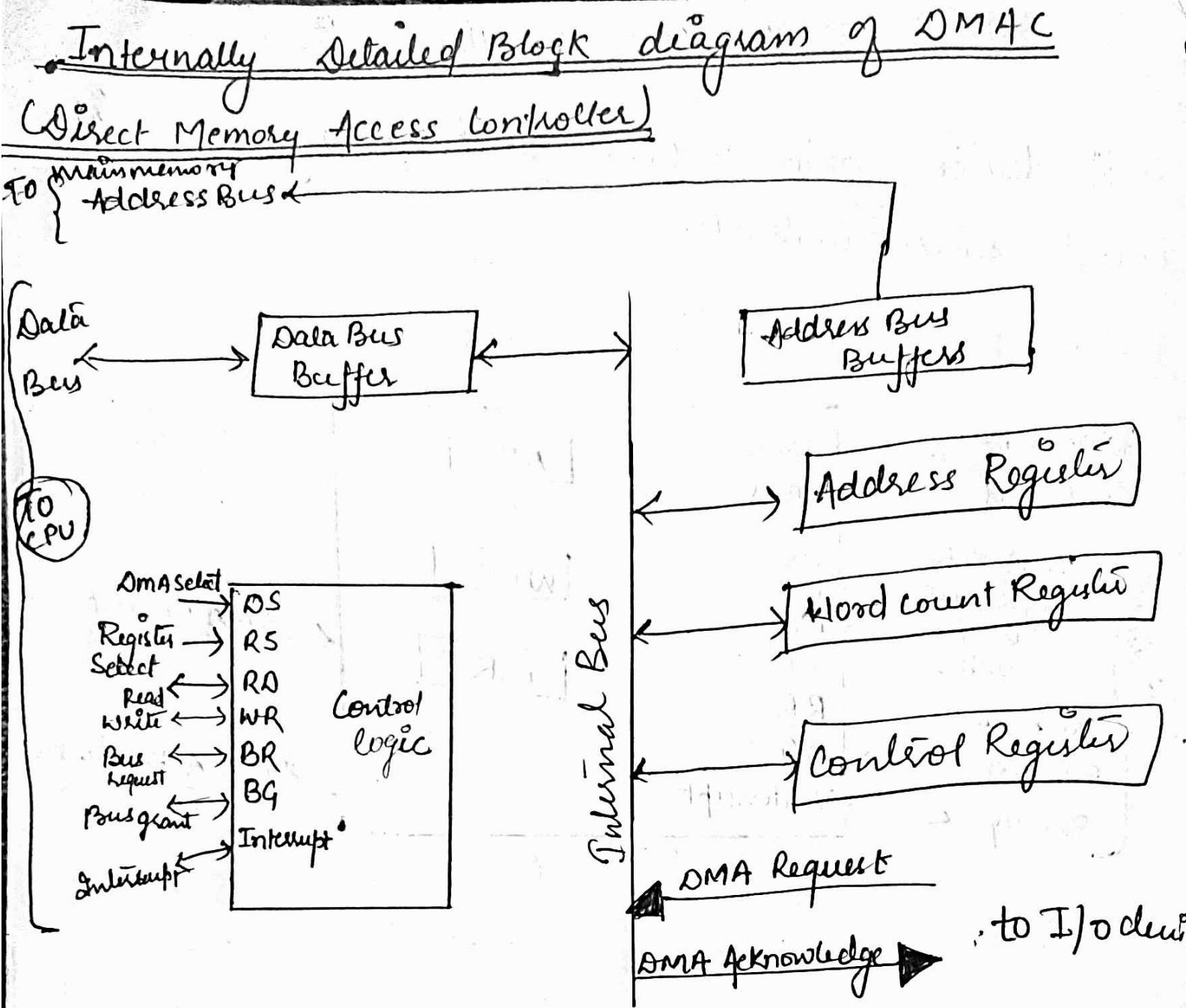
DMA Controller

28



There are 3 register

- ① Address Register (AR) : It holds the address of memory on which memory would be accessed either for Read or write operations. It communicates directly with memory through Address Bus. Here after each transfer (AR) value get increment by 1 (+1).
- ② Word count Register (WCR) : It holds the no of words to be transferred. After each transfer word count value is get decrease by 1 (-1), when it reaches to zero means data transfer is finished completely.
- ③ Control Register (CR) : - It specifies the mode of transfer [whether read or write in memory]



Block Diagram of DMA controller

DMA controller

- The DMA controller needs the usual circuits of an interface to communicate with the CPU & I/O devices.
- The DMA controller has three registers:
(i) Address Register (ii) Word count Register, (iii) Control Register

(i) Address Register: It contains an address to specify the desired location in memory. It is incremented by 1 after each word is transferred to memory.

(ii) Word count Register:

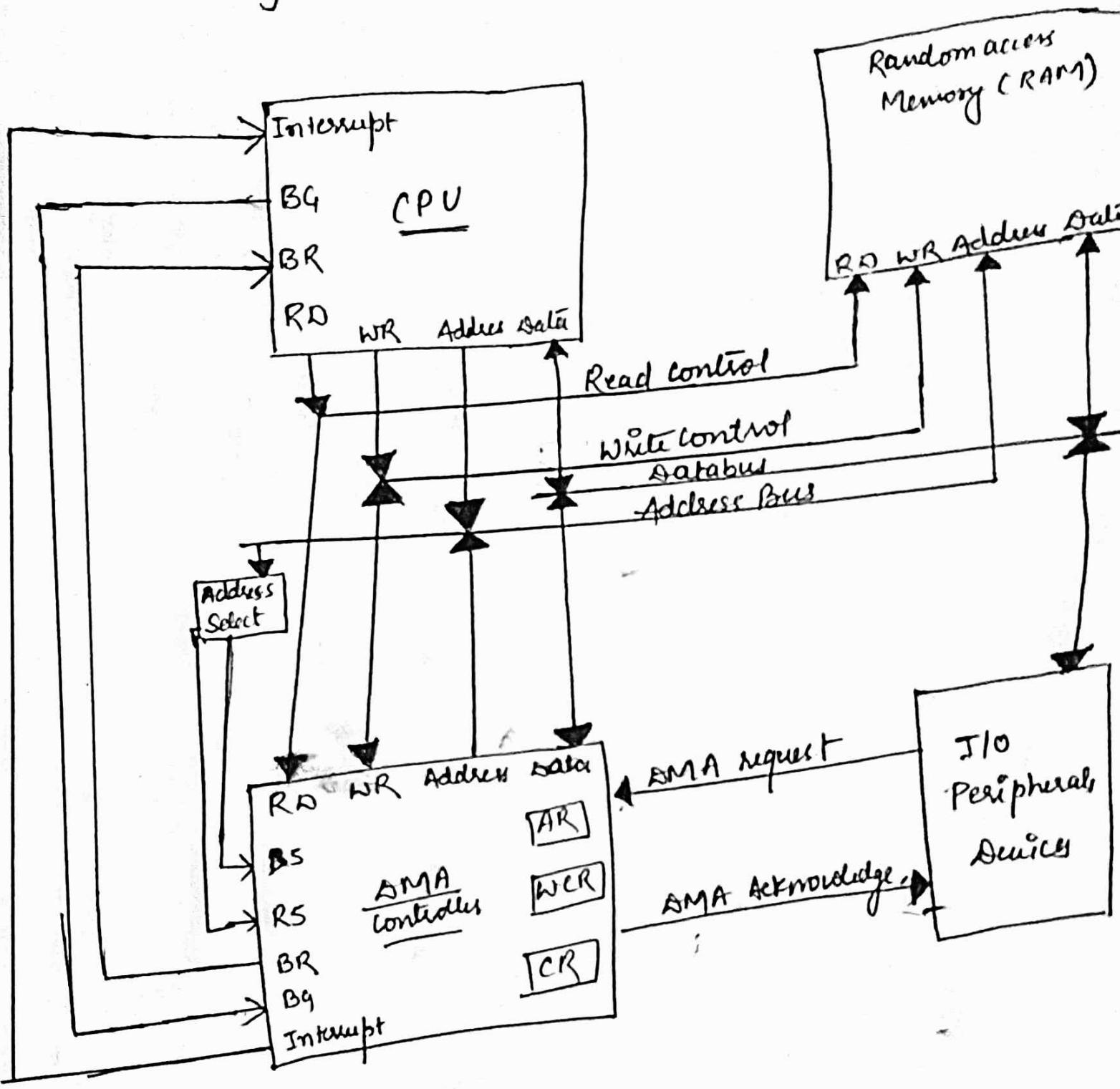
It holds the number of words to be transferred. It is decremented by 1 after each word transferred & internally tested for zero.

(iii) Control Register:

It specifies the mode of transfer, i.e., Read/ write the memory.

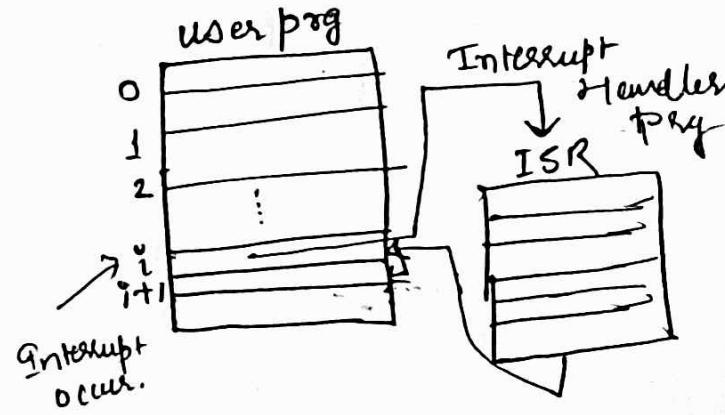
- The DMA controller communicates with the CPU via databus & control lines.
- The registers in DMA are selected by the CPU through the address bus by enabling DS(DMA select) and RS(Register Select) inputs.
The RD(Read) & WR(Write) inputs are bidirectional & either of the one is activated at a time (RD or WR)
- When BG(Bus Grant) IP is 0, the CPU can communicate with the DMA registers through the data bus to read from or write to the DMA registers.
- When BG(Bus Grant) IP is 1, the CPU can relinquish the buses & the DMA can communicate directly with memory by specifying an address in the address bus & activating RD or WR control.
- DMA communicates with the External input/output devices through the request & acknowledge lines by using Handshaking procedure.

DMA Transfer: (How DMA transfer data with I/O)



Interrupts

- Normal execution of programs may be pre-empted if some device requires urgent servicing to deal with the situation immediately, the current execution program must be interrupted.
- Interrupt is a high-priority instruction generated by H/W device on SW program instruction to get immediate attention of CPU
- or
- An interrupt is a request from a H/W or SW program for a service. the make immediate service by the processor.



Interrupt Handler program: Execute the interrupt which is occur.

It will find out which code we execute for this interrupt for executing this code their is an ISR (Interrupt Service routine)

→ This ISR has some hand address, & this will find out by the interrupt handler. By checking the address of this ISR in I.V.T (Interrupt vector Table)

which is stored in a memory somewhere.

Ex: for explanation purpose

- ① printer out off paper (all signals goes to CPU as a interrupt)
- ② Divide by zero, I/O device request for data transfer,

Power failure, overflow.

These all will generate an interrupt which will have to handle immediately.

numeracy zero, divide by zero,
Power failure, overflow.
These all will generate an interrupt which will have to handle
immediately.

Diagram Description

→ As interrupt generated, the processor provides the requested
service by executing an Interrupt service routine (ISR)

→ The state of processor [content of PC, general register, & PSW
(Process Status word)] is first

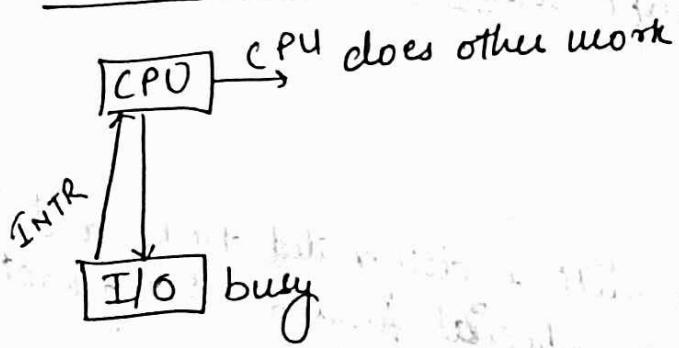
~~PSW~~ PSW cont.

Saved in memory before servicing the interrupt.

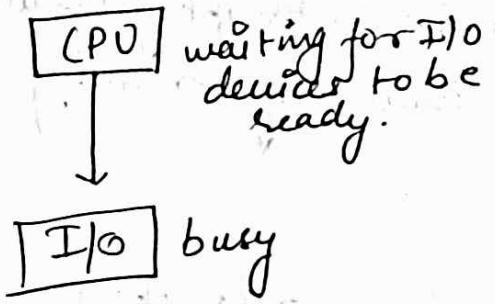
→ When ISR is completed the state of the processor is restored so that interrupted program may continue.

PSW: contain some status flag bit which we will be use to store control information

With interrupt



Without interrupt

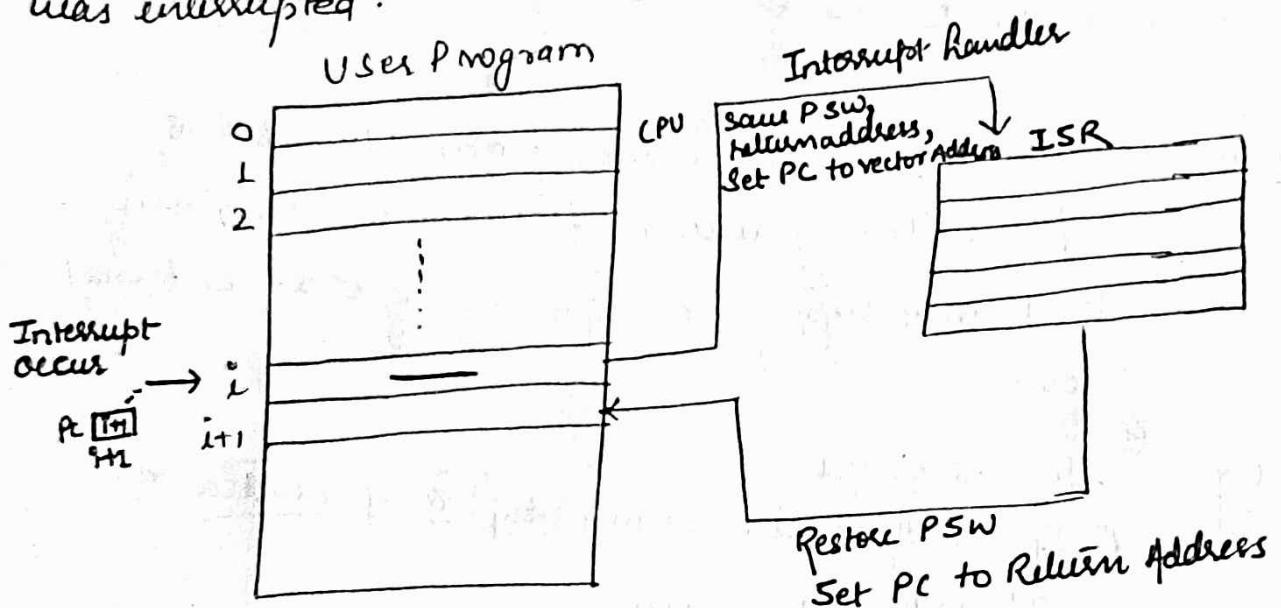


why? Interrupt is used

It is used to improve the performance of CPU.

INTERRUPT HANDLING MECHANISM

- ① Interrupt is generated by h/w device or s/w programs.
- ② CPU suspends current program execution by storing the return address from the Program Counter (PC), general registers values & values in PSW (Program Status word) into a memory stack by push operation.
- ③ CPU executes the Interrupt Service Routine (ISR)
- ④ When ISR finishes, CPU restores (pop) the previous program state & continue the original program that was interrupted.

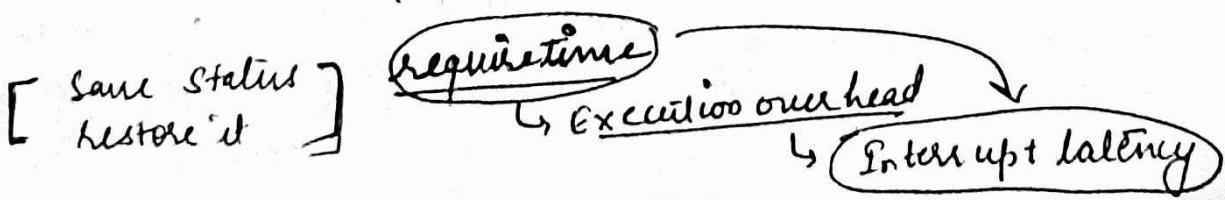


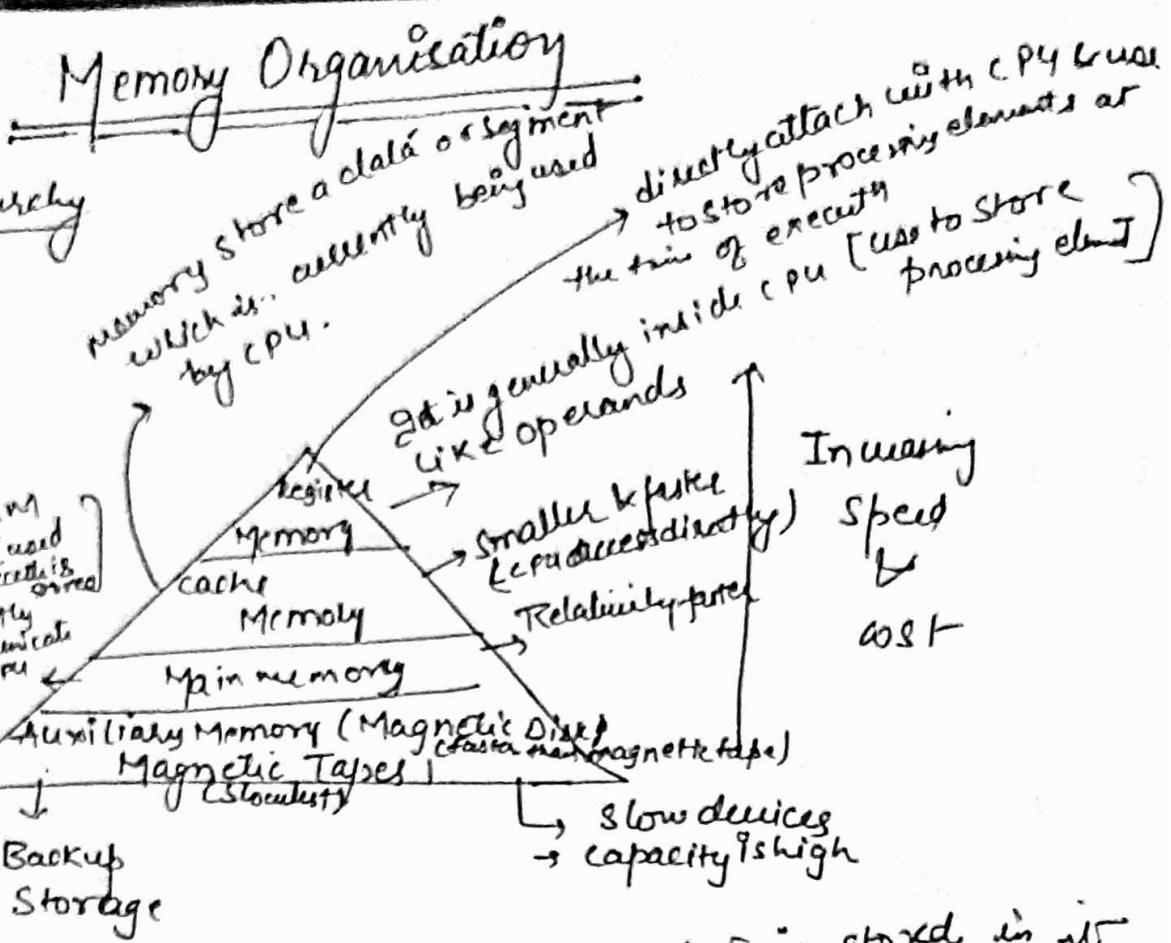
Advantages :-

Efficiency of CPU is improved.

Disadvantages :-

Overhead required to service the interrupt.





Memory:- Essential component because all data is stored in it but it is not possible to store all program or data in a same memory.

- If we use only auxiliary or secondary memory it makes the system slower and not possible to access by CPU same time and very expensive as well as it take time.
- Therefore, economically way is to use more than one memory in memory hierarchy. in digitized computer.
- Memory hierarchy contains all the memory storage devices

Auxiliary Memory (Magnetic Disk & Magnetic Tapes)

- Use for Back up storage
- Slow devices with high capacity
- In which Magnetic tapes are slowest as compared with Magnetic disk.

Main Memory

- Relatively faster than Auxiliary Memory
- It store data & program which is to be executed
- It directly communicate with CPU.

Cache memory

- Smaller & faster
- CPU access data directly
- memory store a data or segment which is currently being used by CPU.

Register Memory

- It is generally inside CPU like operands
- Use to store processing element like operands
- Directly attach with CPU & use to store processing elements at the time of execution.

In Memory hierarchy when you move

① Top to bottom

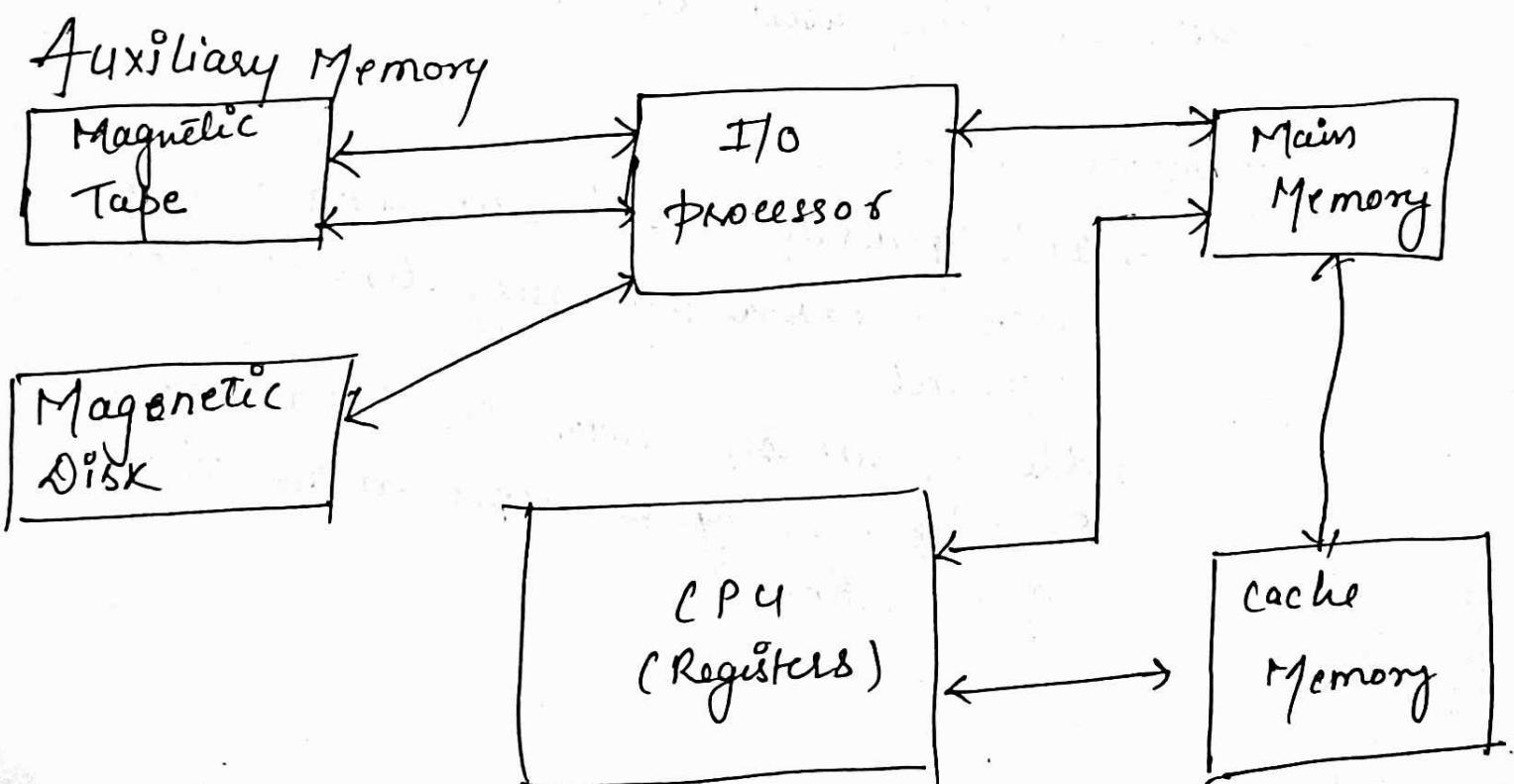
Memory Size Increases &
access Time also Increases

Hence it is slower to access.

② Bottom to Up

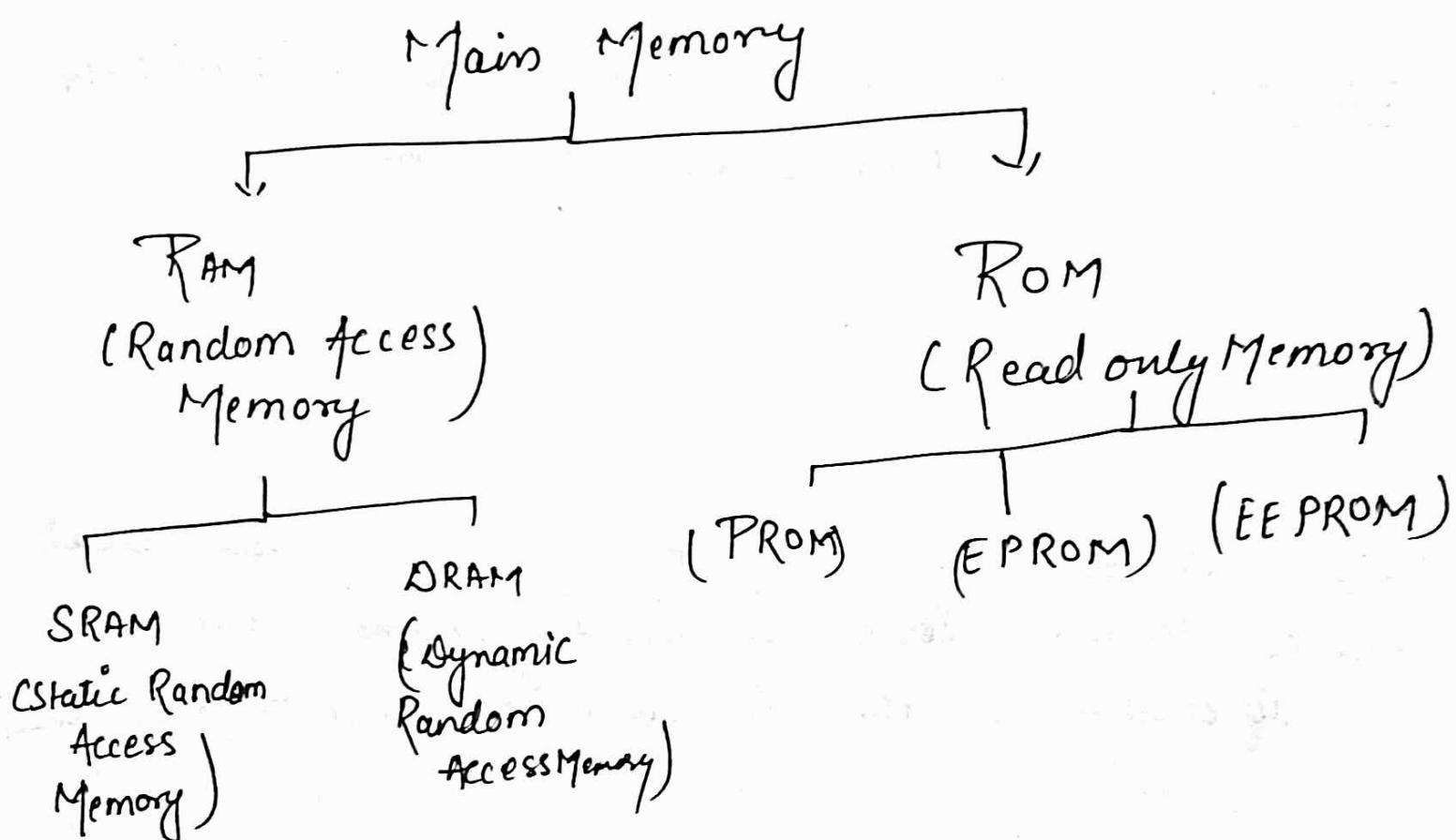
Memory Speed Increases and cost
is also increases

Hence it is faster to access.

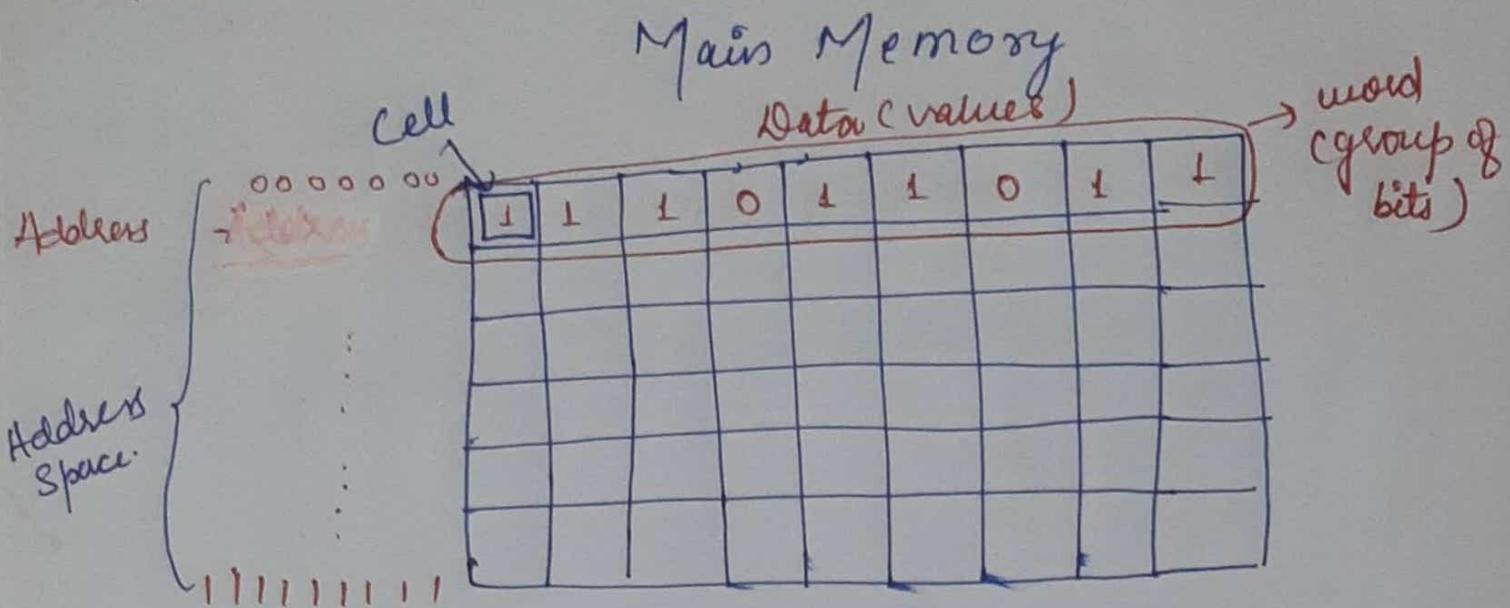


Main Memory

- Central storage unit in computer System
- Fast memory used to store programmes & data during execution.
- Principle technology: Semiconductor IC's.



Internal Structure of Memory (RAM/ROM)



→ MM contains large no of cells. (storage cells)

cells: → Each cell is capable of storing one bit of information (I/O).

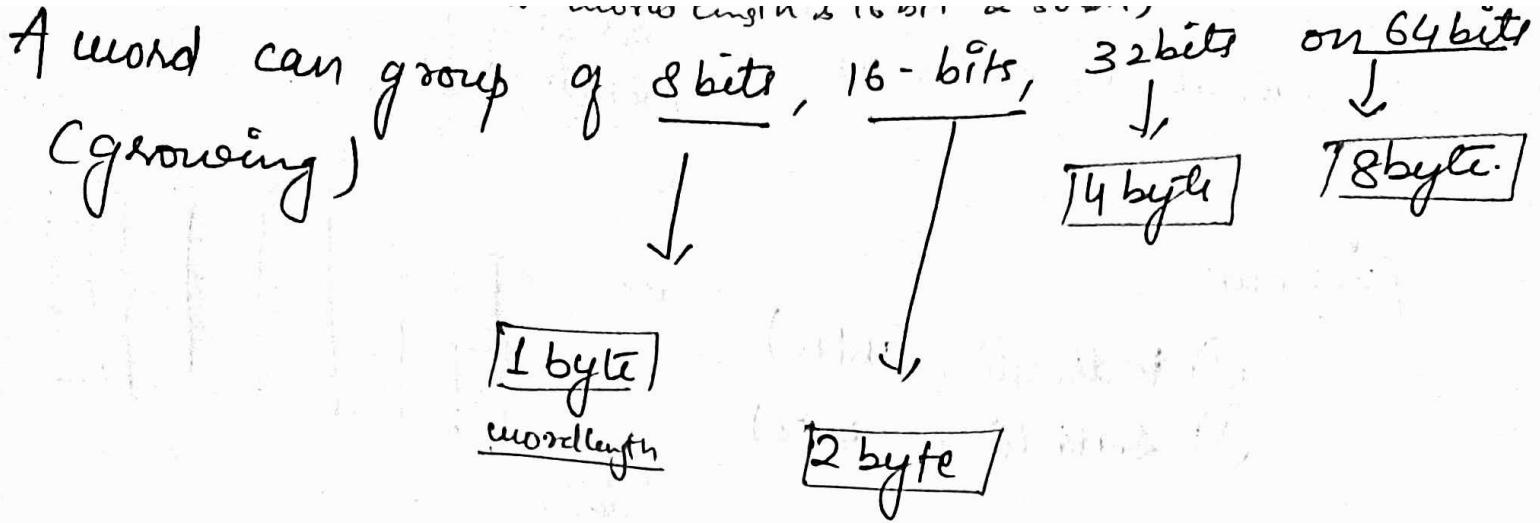
→ cells are processed in groups of fixed size words.

i.e., data is transferred to & from memory in groups of bits called words

address: For easy access different addresses are associated with each word location in memory.

and we also used some control cmd that we wanted to store or retrieve from memory.

word length: No of bits in each word is often referred to
as word length of the computer .. it would be 8 bit,



Data is transmitted from the memory or to memory
 is always in the form of word.

Address space

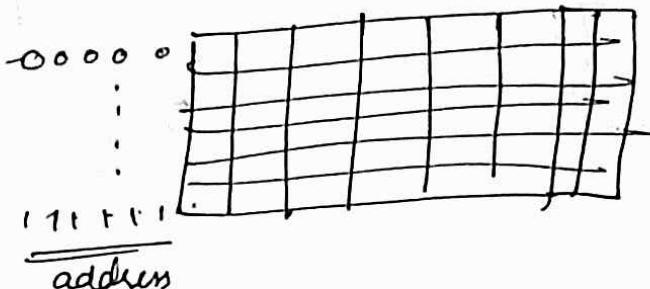
Total no of uniquely identifiable locatⁿ in
 memory.

Ques: Suppose a size of memory is 64KB & a word size of 1 byte
(in which one word is of 1 byte)

[It can also be written as $\frac{64\text{KB} \times 8\text{bit}}{\text{memory size}}$]

find out -

- ① Address line (bits)
- ② Data lines (bits)



Sol

① Address line

Given $64\text{KB} \times 8\text{bit}$

For specifying the address how many bits require

You can write

$$64 \times 8$$

~~64~~ ~~8~~

$$64 = 2^6$$

$$\frac{64\text{ KB}}{\downarrow}$$

$$2^6 \times 2^{10} = 2^{16}$$

Address bit = 16 bit

Data line

from $64\text{ kb} \times \boxed{8\text{ bit}}$

we know

one word is of 8 bits

Therefore, data line also become 8 bits

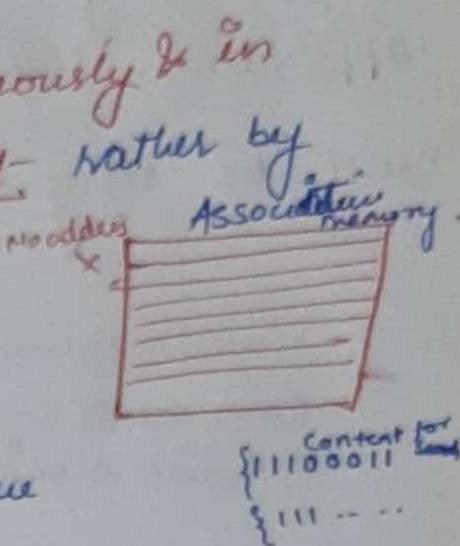
Associ

$1\text{ KB} = 1024\text{ bytes}$

$\Leftarrow \underline{2^{10}\text{ B}}$

Associative Memory

- A memory unit accessed by content is called an Associative Memory or Content Addressable Memory (CAM).
- Associative Memory is accessed simultaneously & in parallel on the basis of data content, rather by specific address or location.
- When a word is written in an associative memory, no address is given.
- This memory is capable of finding an empty unused location to store the word.
- When a word is to be "read" from an associative memory, the content of word is specified or part of word is specified.
- The memory locates all words which match the specified content and marks them for reading.
- Because of its organization, the associative memory is uniquely suited to do parallel searches by data association.
- An Associative Memory is more expensive than Random Access Memory - because each cell must have storage.

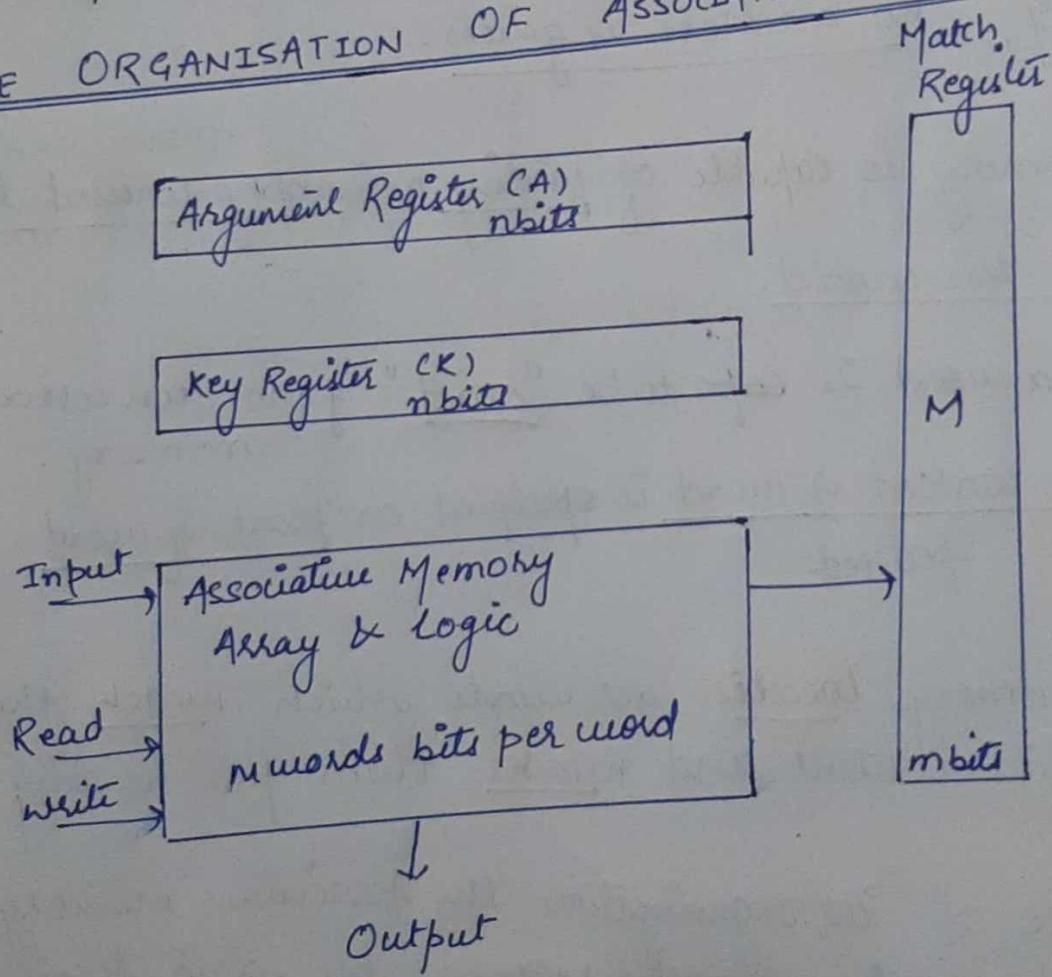


capacity capability as well as logic circuits for matching its content with an external argument

→ That's why, Associative Memory is useful in application. Where Search Time is very critical & must be very short

→ Eg:- Neural Networks because we stores data in patterns. Google search engine, because we provide content on which it will provide result.

HARDWARE ORGANISATION OF ASSOCIATIVE MEMORY



BLOCK DIAGRAM OF ASSOCIATIVE MEMORY

ARGUMENT: It contains words to be searched.

Register (A): It has n bits (one for each bit of word).

KEY REGISTER: It provides a mask for choosing a particular field/key in argument (K) or word.

- It specifies which part of the argument word needs to be compared with words in memory.

- It all bits in key Register are 1's, the entire word should be compared. otherwise, only the bits having 1's in their corresponding position are compared.

Associative Memory Array } → It contains the word that are to be compared. otherwise, only the bits having 1's in their corresponding position are compared. with the argument word in parallel.

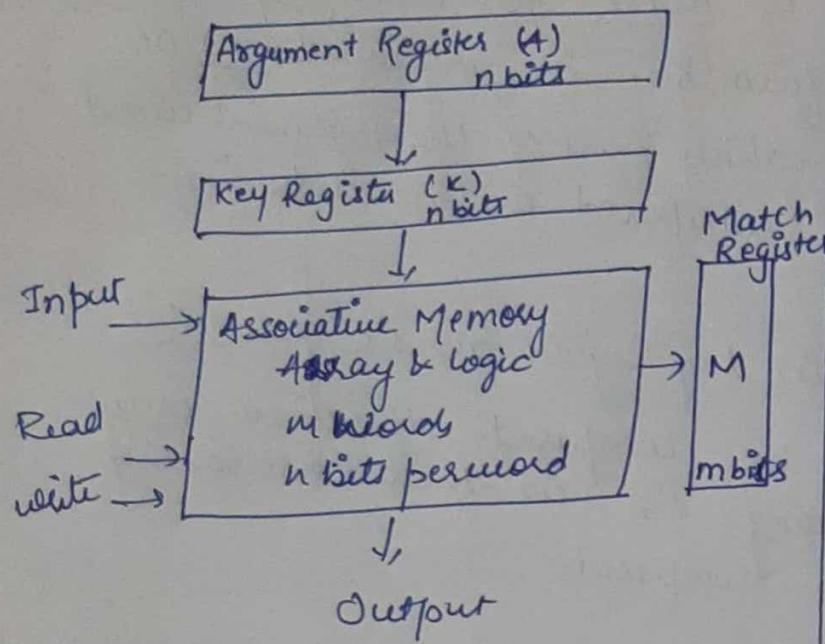
- It consists of m words with n bits per word.

Match logic (m): → It has m bits, one bit corresponding to each word is the memory array.

- After the matching process, the bits corresponding to matching words in match register are set to 1.

Reading is accomplished by sequential access in memory for those words whose match bits are set (or 1).

Explanation of Hardware organisation of Associative Memory with Examples



A	101	LLL LLL 00	M
K	111	000 000 00	1
Word 1	100	1111 00	0
Word 2	101	0000 01	1
Word 3	101	1111 00	1
Word 4	110	0010 10	0
Word 5	111	0001 11	0

Below the table, there is a legend:

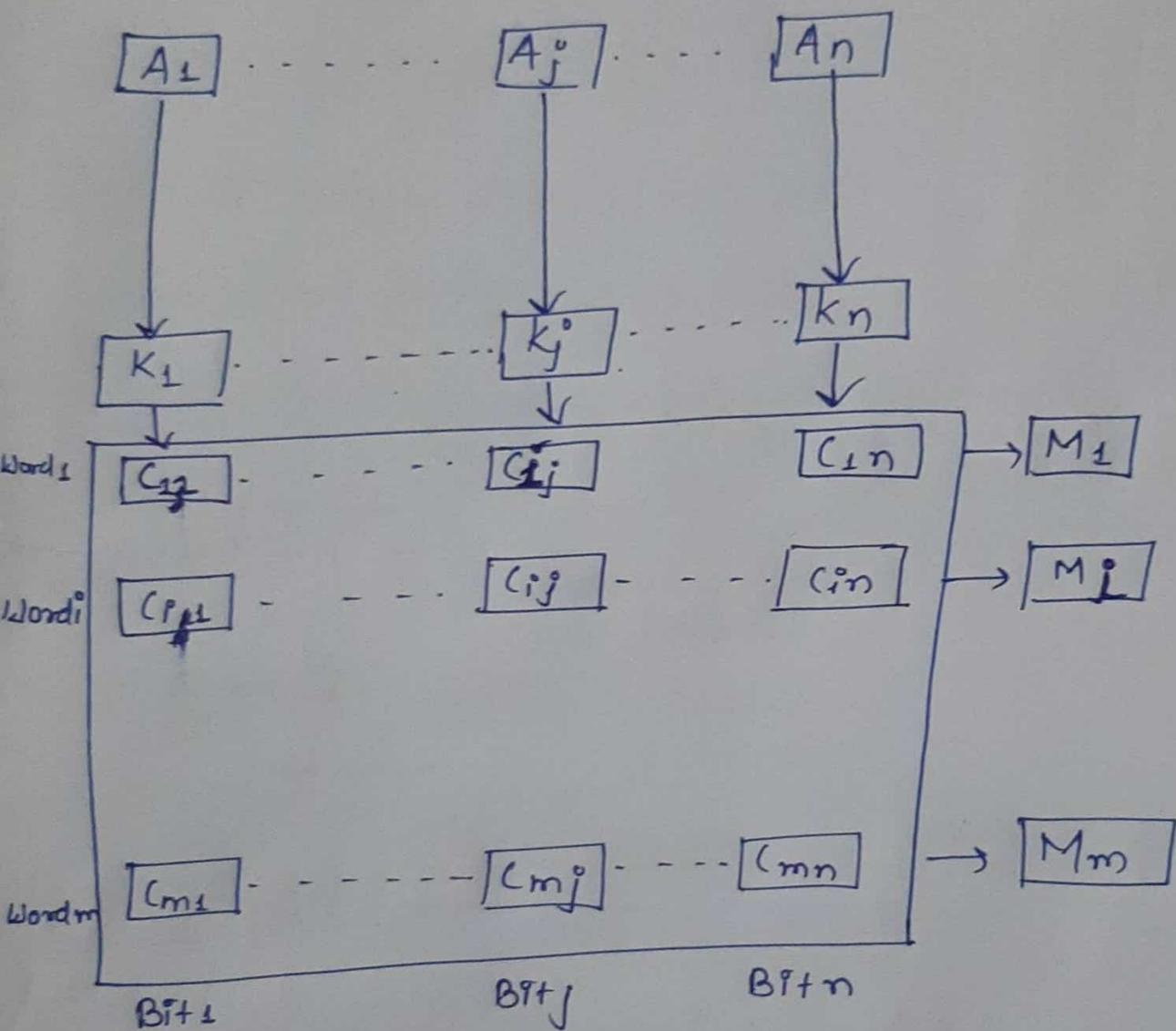
- 0 → no match
- 1 → match

Associative memory of m words, n cells per word

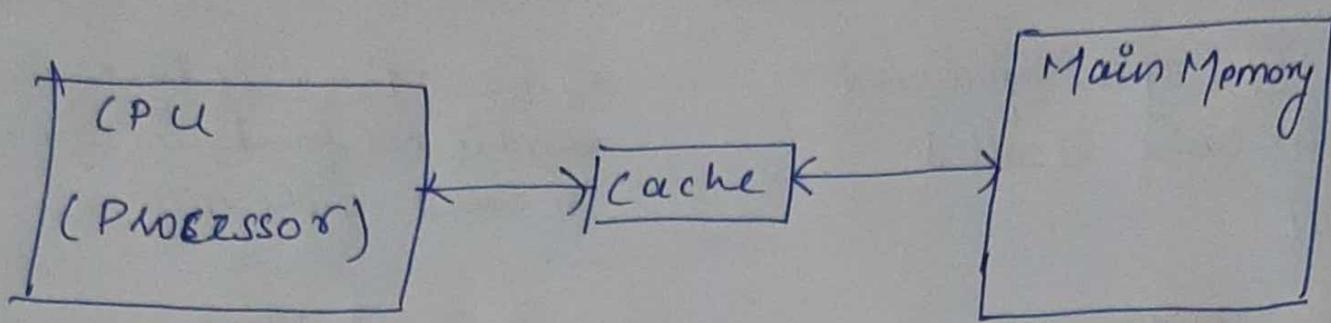
The cell in array is represented by c_{ij}^o : a cell for bit j in word i .

where i = word number

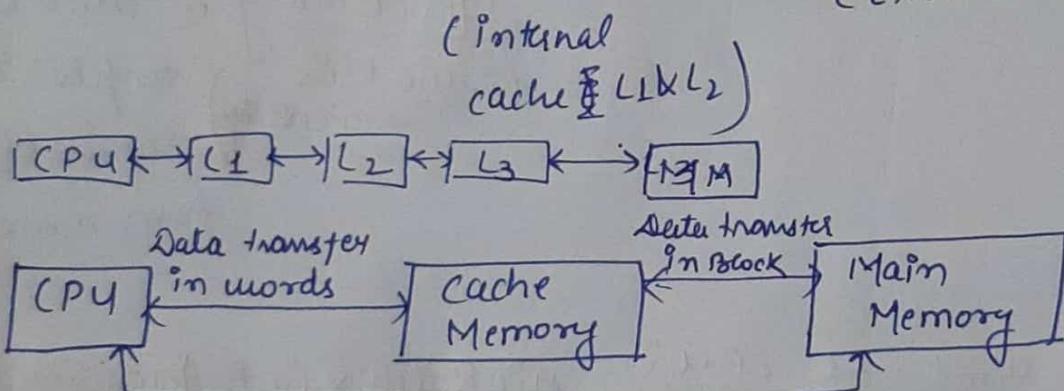
j = bit position in the word.



~~Cache~~ Cache Memory



- Small sized fast memory
- Placed between Main Memory & CPU
- high-Speed volatile Memory & CPU
- Contains most frequently accessed instructions & data
- Located inside the CPU chip or motherboard
(External cache - L3)



Cache Performance

→ It is measured in terms of Hit Ratio

Cache Hit: If the required word is found in cache is called cache hit

Cache Miss: If the required word is not found in cache is called cache miss

$$\text{Hit Ratio} = \frac{\text{Hits}}{\text{Hits} + \text{Miss}}$$

$$= \frac{\text{no. of hits}}{\text{Total no of CPU references}}$$

$$\text{Miss Ratio} = \frac{\text{Miss}}{\text{Hits} + \text{Miss}}$$

$$= \frac{\text{Total no. of miss}}{\text{Total no of CPU references}}$$

Cache access Time : Time required to access word from the cache.

Miss Penalty : Time required to fetch the required block from memory.

$$\text{Average Access Time of CPU} = \frac{\text{Hit Ratio} \times \text{Cache Access Time} + (1 - \text{Hit Ratio}) \times \text{Miss Penalty}}{(1 - \text{Hit Ratio})}$$

$$= \text{Hit Ratio} \times \text{Cache Access Time} + (1 - \text{Hit Ratio}) \times \text{Miss Penalty}$$

$$= h \times T_c + (1-h) + T_m$$

$\therefore T_m = \text{Memory Access Time or Miss Penalty}$

$$\boxed{\text{Miss Penalty} = \text{Cache access time} + \text{Main Memory Access Time}}$$

Ques: The access time of a cache memory is 100ns & that of main memory 1000ns. It is estimated that 80 percent of the memory requests are for read & the remaining 20 percent for write. The hit ratio for read accesses only is 0.9. A write through procedure is used.

- procedure is used.

a) what is the average access time of the system considering only memory read cycle.

b) what is the average access taking into consideration the write cycles.

Solution

or

Misp Penalty = Extra time needed to bring the data from one (block)
level to another
(memory to cache)

a) Total read access time = $0.9 \times 100 + (1 - 0.9) \times (100 + 1000)$

$$= 0.9 \times 100 + 0.1 \times 1100$$

$$= 90 + 110$$

$$= 200 \text{ ns.}$$

Total write access time = $1 \times \max(100, 1000)$

\uparrow \uparrow
cm num

$$= 1000 \text{ ns}$$

Bcz it is write through cache,
largest of the
(cm & num)
access time
is considered

⑤ Average Access Time for both read & write
 $= 80\%$ are Read request + 20% are write request.

$$= 0.8 \times 200 + 0.2 \times 1000$$

$$= 160 + 200$$

$$= 360 \text{ ns}$$

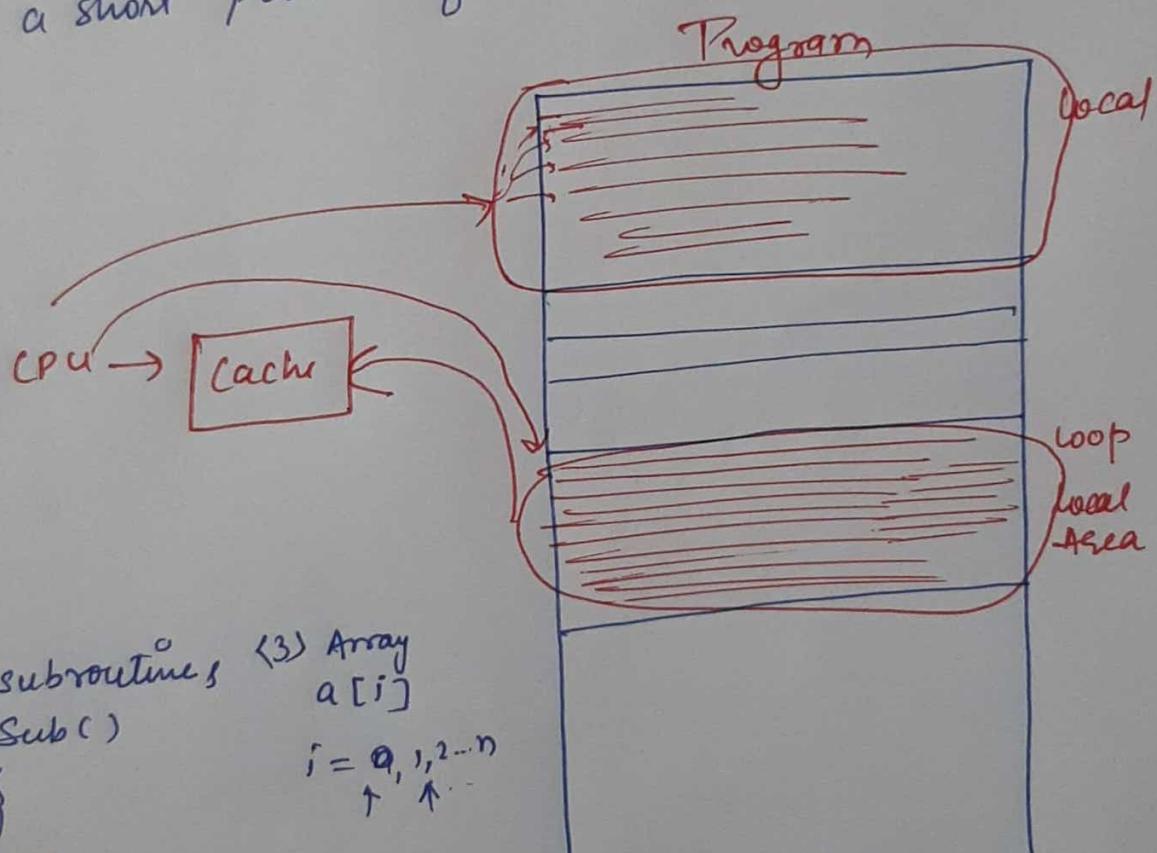
⑥ Hit Ratio = $0.8 \times \underbrace{0.9}_{\text{Hit Ratio}} + 0.2 \times \underbrace{0}_{\text{Hit ratio for write access is always 0}}$

for Read access.

$$= 0.72$$

Locality of Reference in Cache Memory

- Locality of reference is also known as Principle of locality
- Locality of reference is a term for the phenomenon, in which the same values or related storage locations are frequently accessed, depending on the memory access pattern.
- OR
- Locality of reference is the tendency of a processor to access the set of memory address locations repetitively over a short period of time.



(4) Table-lookup procedure

MM

2 Types of Locality of Reference

① Temporal Locality ② Spatial Locality

Temporal Locality :- It says that if a program accesses one memory address, there is a good chance that it will access the same address

→ It refers to reuse of specific data/resources within a short period of time

Eg: loop
for loop (1 to 50)
{
 }
 }

② Spatial Locality : says that if a program access one memory address, (Data Locality) there is a good chance that it will also access other nearby addresses.

→ It refers to use of data elements within relatively close storage location.

Eg: ① Nearly "Every program" exhibits spatial locality bcoz instructions are usually executed in sequences

② In terms of data: Array $a[i], i=0, 1, \dots, n-1$

③ Records: Employee Id/Name/Address/Phone



Levels of Cache Memory

Level 1 (L₁) Cache :

- L₁ is fastest cache & it usually comes within the processor.
- Typically it ranges in size from 8KB to 64KB & was the high-speed SRAM (Static RAM) instead of the slower & cheaper DRAM (Dynamic RAM) used for main memory.
- It is referred to as Internal cache or Primary cache.

Level 2 (L₂) cache :

- L₂ cache is larger but slower in speed than L₁ cache.
- stores recently accessed information.
- Also known as Secondary cache, it is designed to reduce the time needed to access data in cases where data has already been accessed previously.
- L₂ cache comes between L₁ and RAM & is bigger than Primary cache.
- typically ranges from 64KB to 4MB



Level 3 (L₃) Cache

- L₃ cache memory is an exchanged form of memory present on the motherboard of the computer.
- L₃ cache is a memory cache that is built into the motherboard.
- It is used to feed the L₂ cache, and is typically faster than system's main memory but slower than L₂ cache.
- It has more than 3MB of storage unit. in it.
- Not all processors have this L₃ cache, b'coze it is used to enhance the performance of L₁ & L₂ cache.

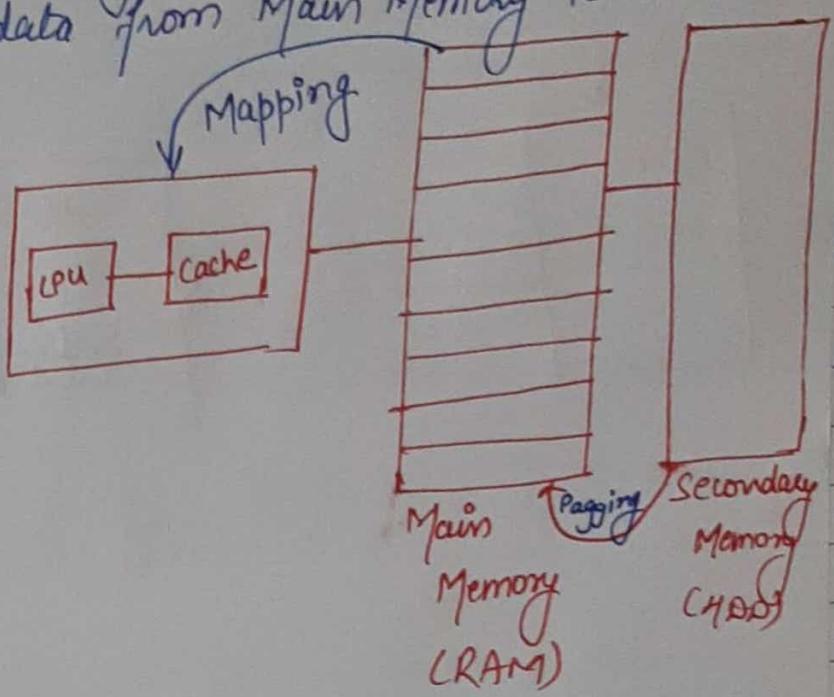
Cache Mapping

Cache Mapping is a technique by which content of Main Memory is brought into the Cache Memory.

→ It is a transformation of data from Main Memory to Cache Memory.

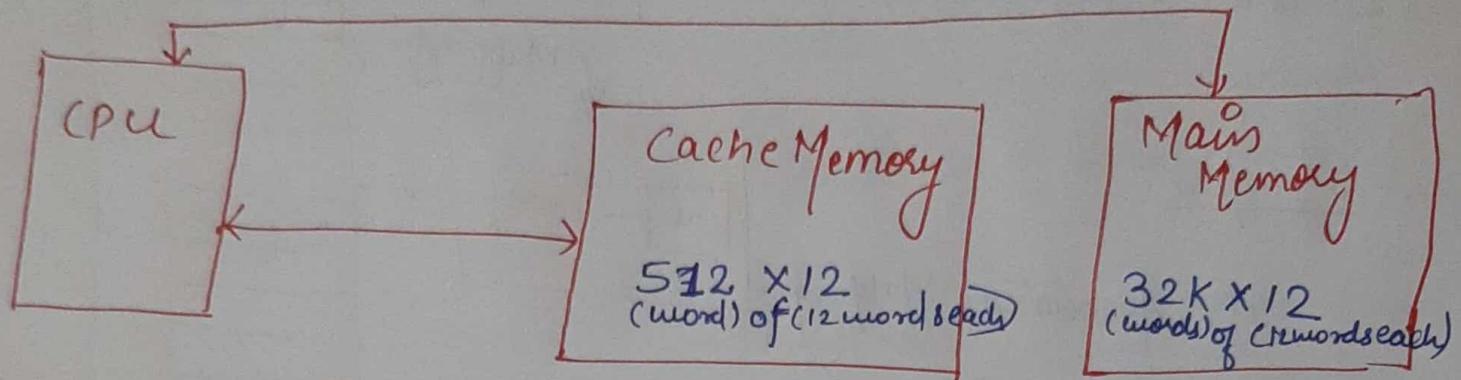
→ 3 Types :-

- ① Associative Mapping
- ② Direct Mapping
- ③ Set-associative Mapping



* content of Transferring (loading) data (content) of secondary memory to Main Memory is called Paging.

Example! - Common Example for understanding all these mapping.



32 K [words]

$$2^5 \times 2^{10} = 2^{15} \text{ words}$$

CPU address = 15 bits

data = 12 bits