

MEMORY ORGANIZATION

Unit - 4

Topics will be covered :

1. Basic concepts, memory hierarchy and its need.
2. Parameters used to measure the performance.
3. Types of memory components, semiconductor RAM memories.
4. Memory module design.
5. ROM
6. Cache memories
7. Mapping functions
8. Replacement algorithms.
9. Memory interleaving
10. Memory performance considerations.
11. Virtual memory organization

- Both program and data are present in the memory.
- The execution speed of programs is highly dependent on the speed with which instructions and data can be transferred between the processor and the memory.
- The memory would be fast, large, and inexpensive. However it is impossible to meet all ~~these~~ three requirements simultaneously.
- Increased speed and size are achieved at increased cost.

Basic Concept

→ Maximum memory size of a computer can be determined by the addressing scheme.

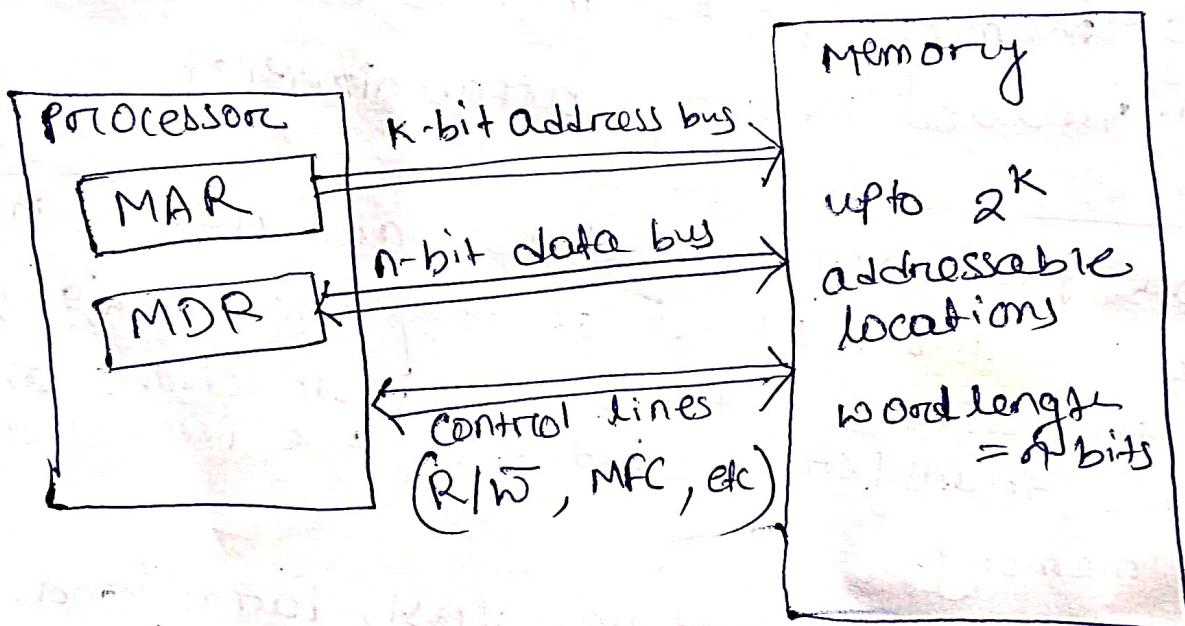
Ex :- A 16-bit computer generates 16 bit addresses is capable of addressing upto $2^{16} = 64K$ memory location.

A 32-bit Computer $\rightarrow 2^{32} = 4G$ memory location.

A 40-bit Computer $\rightarrow 2^{40} = 1T$ memory location.

→ The no. of locations represents size of address space of the computer.

→



Ex :- A 32-bit Computer

Address bits = 32 bits.

High-order 30 bits determine which word to be accessed.

Low-order 2 bits of the address specify which byte location is involved.

If MAR is k bits long and MDR is n bits long, the memory unit may contain up to 2^k addressable locations. During a memory cycle, n bits of data are transferred between the memory and the processor.

- The control lines Read / write (R/W) and memory function complete (MFC) for coordinating data transfers.
 - The processor reads data from memory by loading the address of the required memory location into the MAR register and setting R/W line to 1. The memory responds by placing the data from the addressed location onto the data lines and confirms this action by asserting the MFC signal. After receipt of the MFC signal, the processor loads the data on the data lines into the MDR register.
 - The processor writes data into a memory location by loading the address of this location into the MAR and indicates that a write operation is involved by setting R/W line to 0.
 - Block transfer operation → It consists of consecutive memory location ~~across~~ ⁱⁿ Read / write operations.
 - Memory access time :- Time elapses between the initiation of an operation and the completion of that operation. Ex: time between Read and MFC signal.
 - Memory cycle time :- minimum time delay between the initiation of two successive memory operations.

Random access memory: - If any location can be accessed for a Read or write operation in some fixed amount of time.

Cache memory:-

- to reduce the memory access time
- It is a fast, small memory that is inserted between the larger, slower main memory and the processor.
- It holds the currently active segments of a program and their data.

Virtual memory:-

- is used to increase the apparent size of the physical memory.
- CPU generates virtual or logical address (processor)
- Memory Management Unit is a mapping function which maps logical address to physical address.

Internal organization of memory chips

16 x 8 organization.

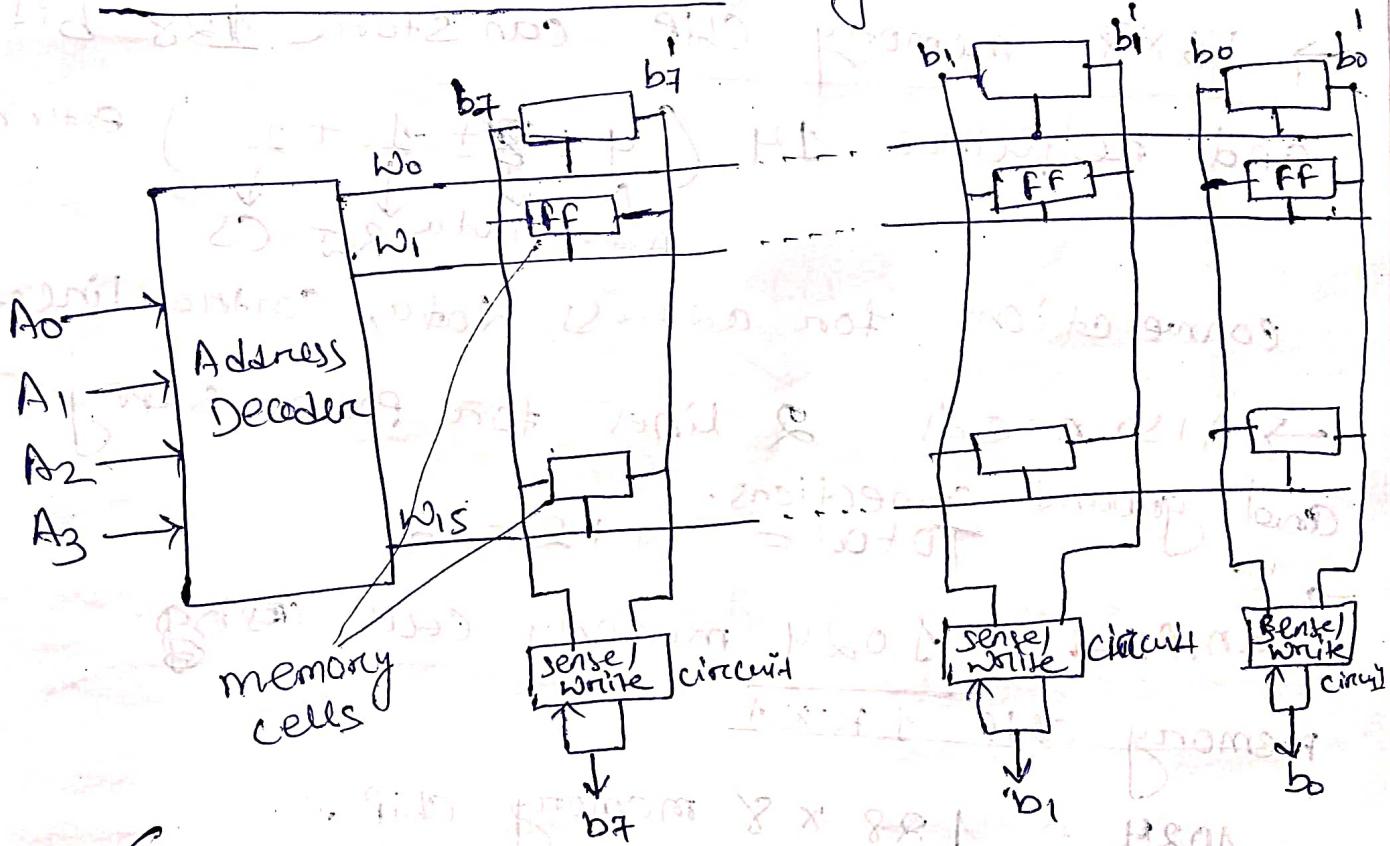
↓ ↗ each word 8 bits.

16 NO. of words

→ memory cells are usually organized in the form of an array, in which each cell is capable of storing one bit info.

→ Each row of cells constitutes a memory word and all cells of a row are connected to a common line referred to as the word line.

→ The cells in each column are connected to a sense/write circuit by two bit lines.



(Organization of 8 bit cells in a memory chip)

→ The sense/write circuits are connected to the data input/output lines of the chip.

→ During a read operation, these circuits sense or read the information stored in the cells selected by a word line and transmit this information to the output data lines.

→ During a write operation, the sense/write circuits receive input information and store it in the cells of the selected word.

Example:- 16×8 memory chip organization.

- Two control lines, R/W and CS are provided in addition to address and data.
- The R/W (Read/Write) input specifies the required operation, and the ~~data~~ output.
- CS (Chip select) input selects a given chip in a multichip memory system.
- 16×8 memory chip can store $128 \rightarrow 128$ bits and require 14 ($4 + 8 + \frac{1}{2} + 1$) external connections for address, data, control lines.
 $\downarrow \quad \downarrow \quad \downarrow \quad \downarrow$
Address Data R/W CS
- Also needs 2 lines for power supply and ground connections.
- Total = $14 + 2 = 16$

Example:- 1024 memory cells ~~containing~~
memory chip $1K \times 1$

$1024 = 128 \times 8$ memory chip.

External connection = $7 + 8 + 2 = 17$

$1024 = 1024 \times 1$

$1024 = 1K \times 1$ memory chip.

External connection = $10 + 1 + 4 = 15$

$1K \times 1, 128 \times 8, 64 \times 16$

$$1024 = 2^7 \times 2^3$$

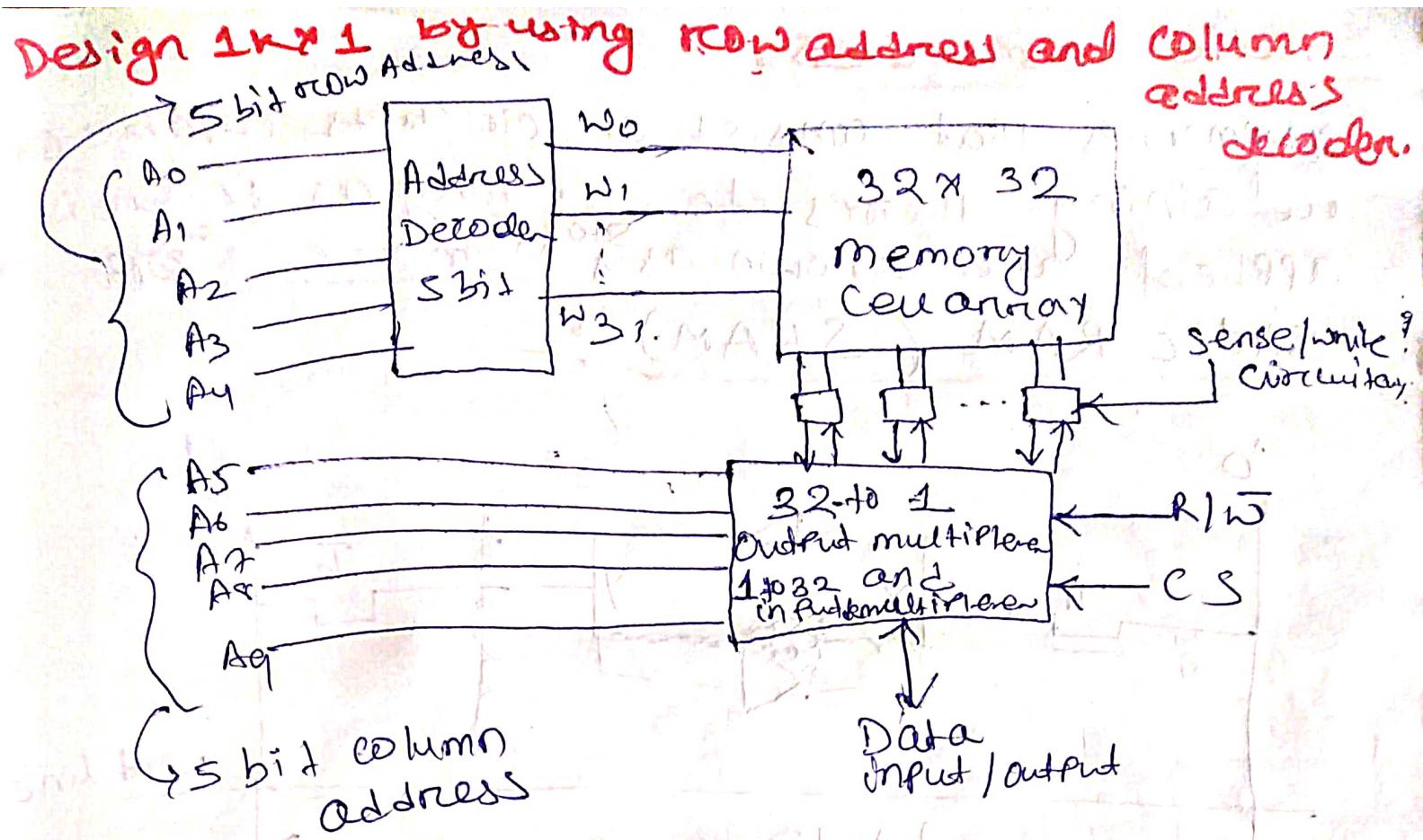
$$1024 = 2^6 \times 2^4$$

1024 ~~is~~ memory cell

$$1K \times 1 \Rightarrow 15$$

$$128 \times 8 \Rightarrow 19$$

$$64 \times 16 \Rightarrow 6 + 16 + 4 = 26$$



$$1024 = 2^{10}$$

10 bit
5 bit (row)
5 bit (column)

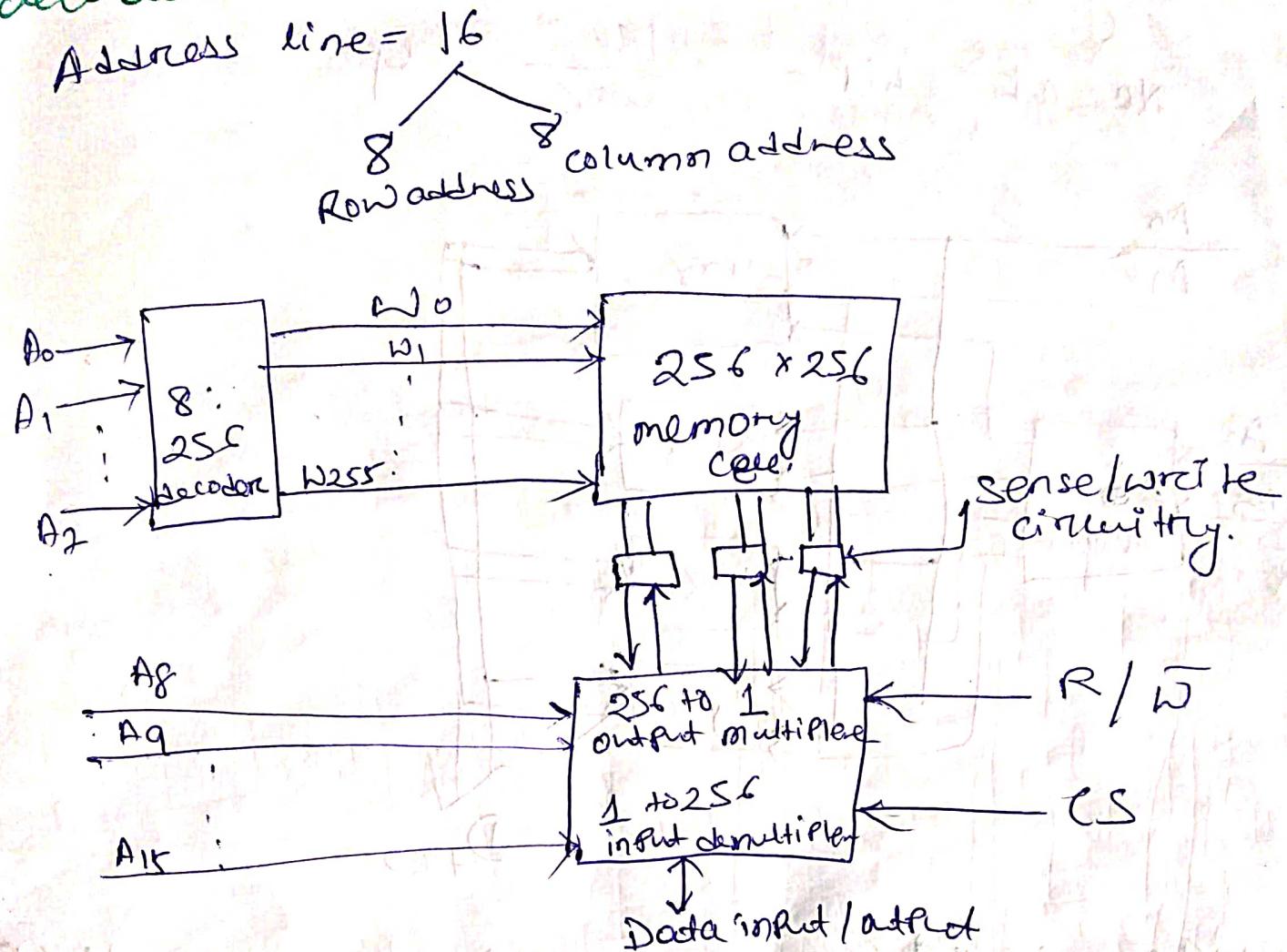
- A row address selects a row of 32 cells, all of which are accessed in parallel. However, according to the column address, only one of these cells is connected to the external data line by the output multiplexer and input demultiplexer.

Bn :- 4M bit chip

$$\begin{aligned} &= 4 \times 2^{20} = 512K \times 8 \\ &= 2^{22} = 2^{19} \times 2^3 = 2^{22} \end{aligned}$$

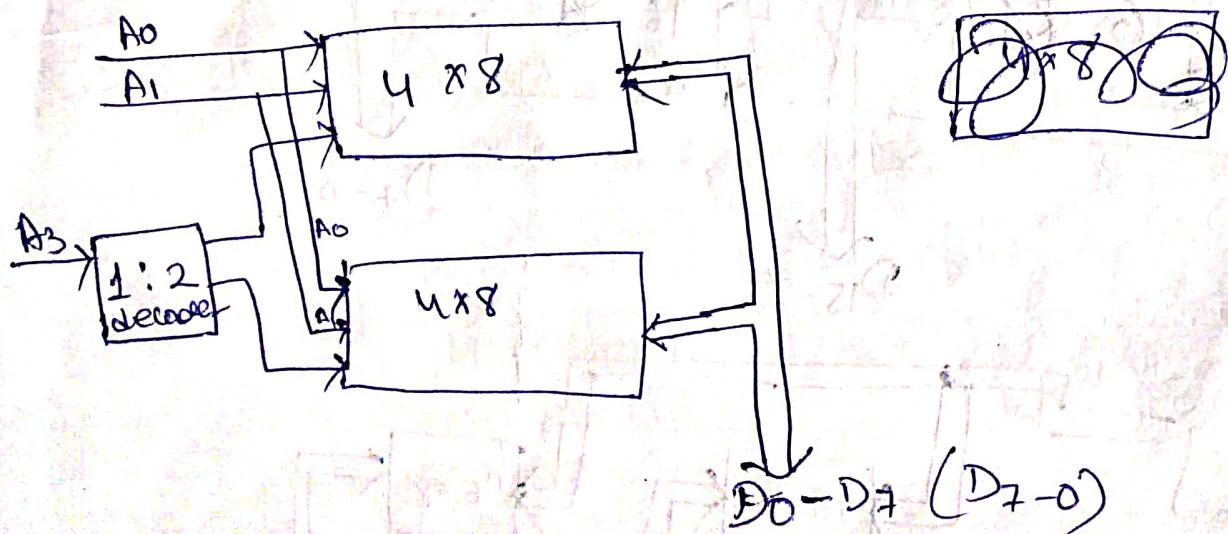
- 4M bit chip using $512K \times 8$ organization
- Total External connection = $19 + 8 + 4 = 31$
- Address bus = 19, Data bus = 8.

Design $64K \times 1$ using Row and column address decoder.



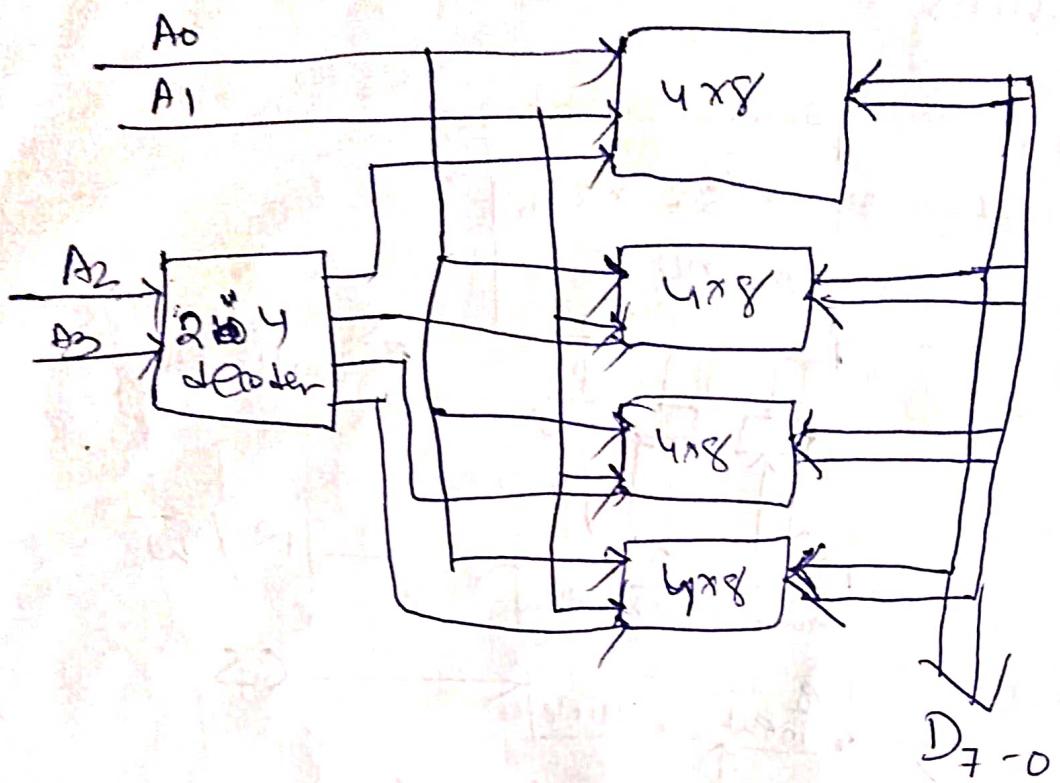
Design 8×8 memory by using 4×8 memory cell.

$$\text{No. of Chip} = \frac{8 \times 8}{4 \times 8} = 2$$



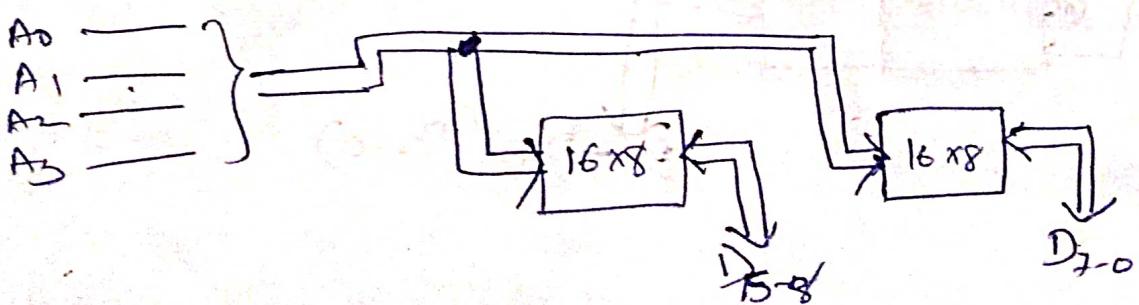
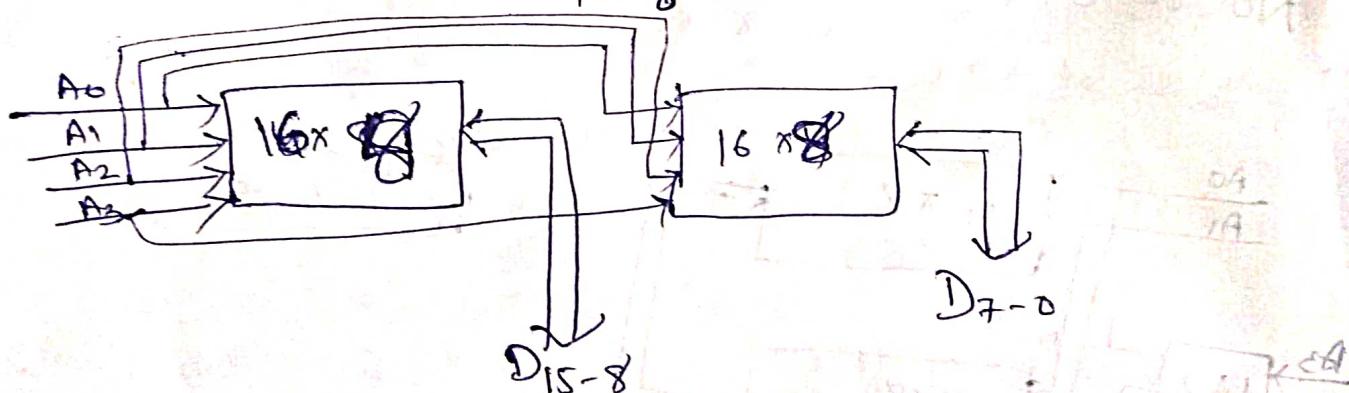
→ ~~16~~
 Design 16×8 memory ~~chip~~ chip by using 4×8
 memory chip

$$\text{No. of Chip} = \frac{16 \times 8}{4 \times 8} = 4$$



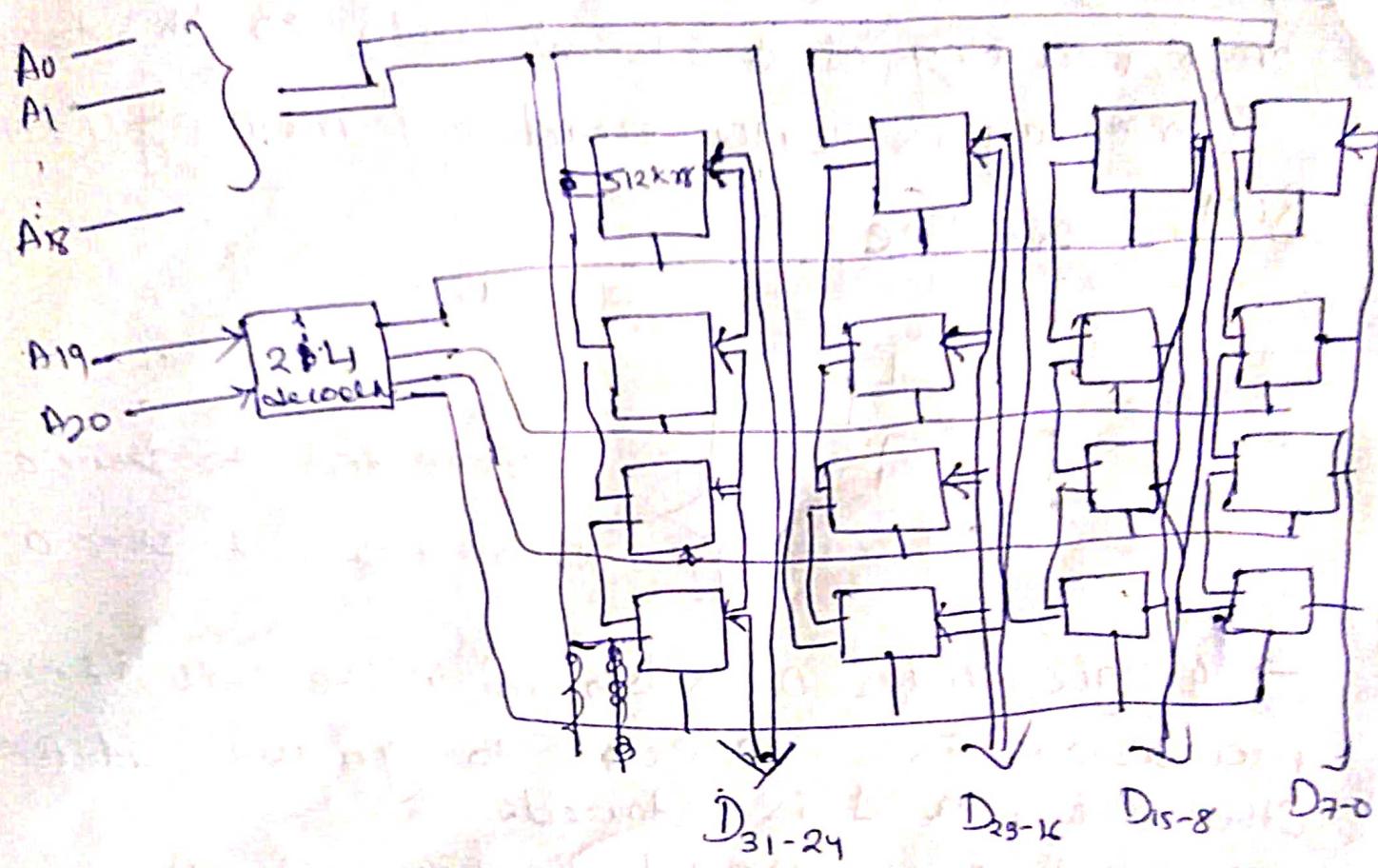
Design $16 \times 8 \times 6$ memory chip by using 16×8 memory chip.

$$\text{No. of Chip} = \frac{16 \times 8 \times 6}{16 \times 8} = 2$$

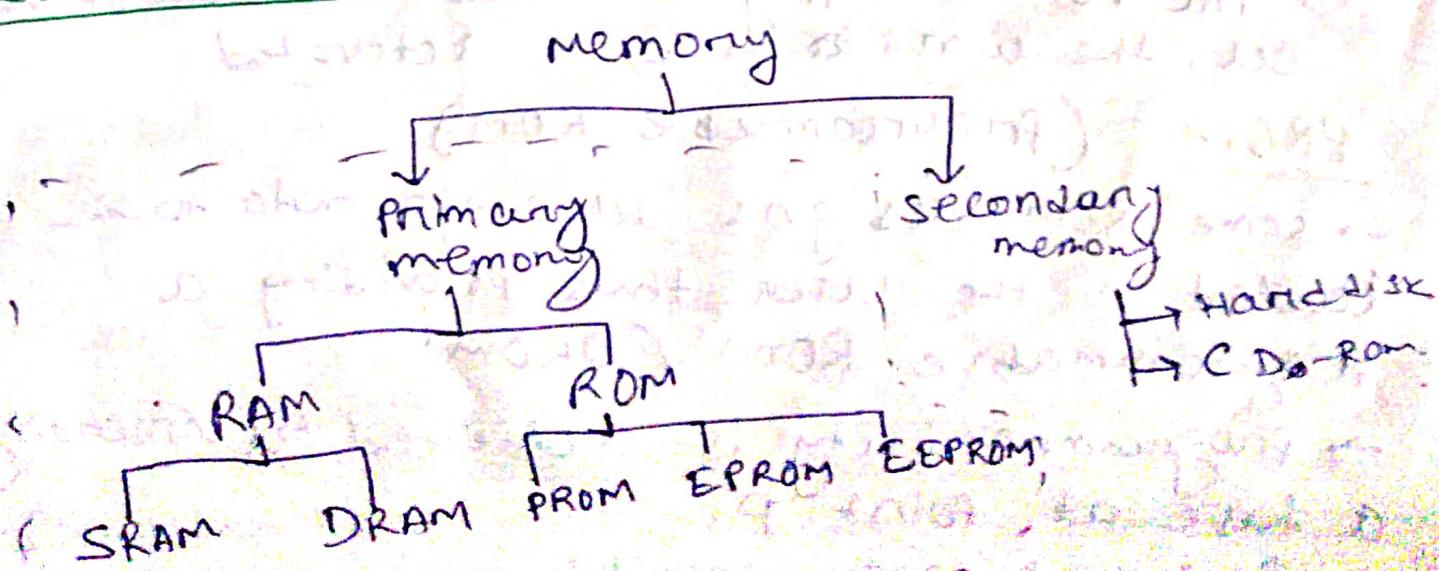


design a $2M \times 32$ memory chip, using $512K \times 8$ memory chip

$$\text{No. of chips} = \frac{2M \times 32}{512K \times 8} = \frac{2^{21} \times 2^5}{2^{19} \times 2^3}$$
$$= \frac{2^2 \times 2^2}{\text{in column}} = \frac{2^2}{\text{in a row}}$$

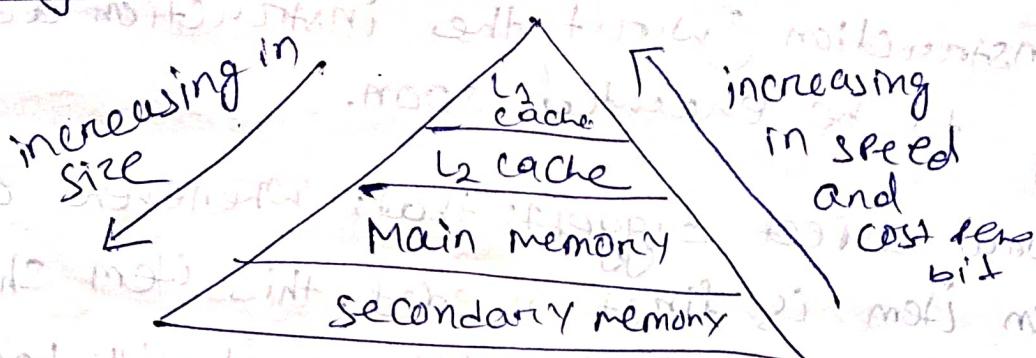


Classification of memory:



CACHE MEMORY:-

Memory Hierarchy :-



Why Cache ?

→ The processor spends more time to access data and instruction from main memory. A Cache memory is introduced in between CPU and main memory to improve average access time. It is a small and fast storage unit.

Locality of reference:-

- Analysis of program shows that most of their execution time is spent on routines in which instructions are executed repeatedly. (Simple loop, nested loop).
- Many instructions in localized areas of the program are executed repeatedly during some time period, and the remainder of the program is accessed relatively infrequently. This is referred to as locality of reference.

→ It manifests itself in two ways: temporal and spatial.

Temporal:

→ A recently executed instruction is likely to be executed again very soon.

Spatial:

→ Instructions in close proximity to a recently executed instruction (w.r.t the instruction address) are likely to be executed soon.

→ The temporal aspect suggests that whenever an information item is first needed, this item should be brought into the cache where it will hopefully remain until it is needed again.

→ The spatial aspect suggests that instead of fetching just one item from the main memory to cache memory, it is useful to fetch several items that reside at adjacent address as well.

→ Block refers to a set of contiguous address locations of some size. A cache block is also called cache line.

Read Request and Write Request

- When a read request is received from the processor, the contents of a block of memory words containing the location specified are transferred into the cache one word at a time.
- Subsequently, when the program references any of the locations in this block, the desired contents are read directly from the cache.
- The cache memory contains less no. of blocks compared to the total no. of blocks in main memory.
- A mapping function specifies the correspondence between the main memory blocks and those in cache memory blocks.
- When a cache is full and the requested memory word is not present in the cache, the cache control hardware must decide which block should be removed to create space for the new block. This decision is taken by replacement algorithm.
- The cache control circuit determines whether the requested word currently exists in the cache. If it does, the read or write operation is performed on the appropriate cache location. In this case, read or write hit is occurred.
- In read operation, the main memory is not involved.
- In write operation, the system can proceed in two ways
 - Write Through Protocol
 - Write Back Protocol

- In write-through protocol, the cache location and main memory location are updated simultaneously.
- In write-back protocol, the cache location is updated only and marked it as updated with an associated flag bit, called dirty or modified bit.
- The main memory location is updated later when that block is selected for replacement. This protocol is also called copy-back protocol.
- Write-through protocol is simple but it results in unnecessary write operations because a cache block is written back in the main memory when a given cache word is updated several times during its cache residency.
- Write-back protocol also results unnecessary write operation because when a cache block is written back to the memory all words of the block are written back, even if only a single word has been changed while the block is in cache.
- When the addressed word in a read operation is not in the cache, a read miss occurs. The block of words that contains the requested word is copied from the MM into the cache. After the entire block is loaded into the cache, the particular word requested is forwarded into the processor.
- Alternatively, this word may sent to the processor as soon as it is read from MM. This approach is called load-through or early read, reduces the processor's waiting period.
- During write operation, if the desired address is not present in the cache, a write miss occurs. In write-through protocol, directly written into the MM and in write-back, first bring into cache then update it.

- Hit rate and Miss Penalty: — ⑥
- A successful access to data in a cache is called a ^{Cache} hit.
 - An unsuccessful access to the data in a cache is called cache miss.
 - The no. of cache hit stated as a fraction of all attempted accesses is called the hit rate.
 - The no. of cache miss stated as a fraction of all attempted accesses is called the miss rate.
 - When the cache miss occurred, the information is not present in cache. So it take entire time to bring the data from main memory to cache memory is called miss penalty.

$$\text{Hit rate} = \frac{\text{No. of successful access}}{\text{Total no. of attempted access}}$$

$$\text{Miss rate} = \frac{\text{No. of unsuccessful access}}{\text{Total no. of attempted access}}$$

Average access time (Tavg)

$$T_{avg} = h * C + (1-h) M \quad \text{or } h * T_C + (1-h) * T_M$$

where h = hit rate

C = cache hit time

$$M = \frac{\text{miss penalty}}{\text{miss hit ratio}}$$

$$\text{Improve in performance} = \frac{\text{Time without cache}}{\text{Time with cache.}}$$

Questions

① What is the hit ratio if MM access time is 30 nsec and cache memory access time is 30 nsec and average access time is 42 nsec.

Ans:- $h=? \quad T_{avg} = h \times C + (1-h) \times M$

$$\Rightarrow 42 = h \times 30 + (1-h) \times 150$$

$$\Rightarrow 42 = 30h + 150 - 150h$$

$$\Rightarrow 120h = 150 - 42$$

$$\Rightarrow h = \frac{108}{120} = .9$$

$$= \frac{9}{10} = .9$$

② Hit rate of instruction is 95%, Hit rate of data is 90%. Assume that from 130 fetch, 100 are instruction fetch and 30 are data fetch find the improvement in performance if without cache is 10 cycle needed for read access and with cache

1 cycle for read access in case of hit and 17 cycles for read access in case of miss.

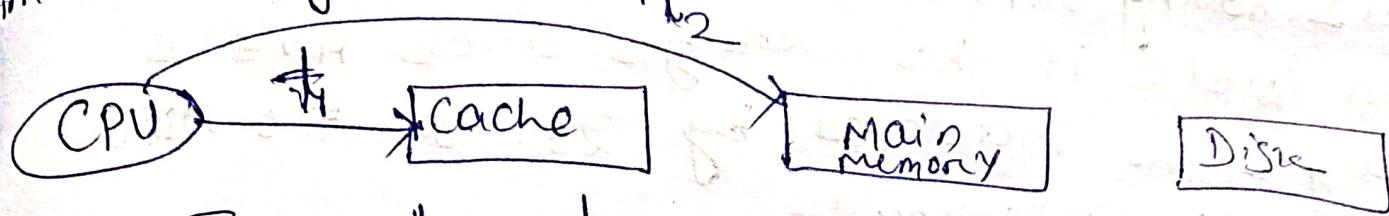
Ans

$$\text{Performance} = \frac{130 \times 10}{\text{Time with cache}} = 5.038$$

~~T_{avg}~~ $\rightarrow T_{avg}$ for data + T_{avg} for instn

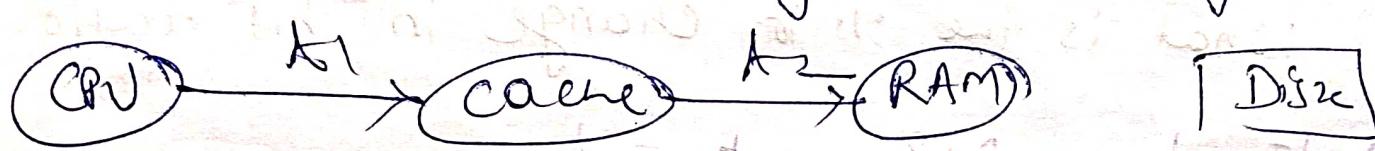
$$= \frac{1}{30} (9 \times 1 + 1 \times 17) + (0.95 \times 1 + 0.05 \times 17) \times 100$$

case-1:- Processor having access to both the levels simultaneously.



$$T_{avg} = h * t_1 + (1-h) * t_2$$

case-2:- Processor is having access to only 1 level.



$$T_{avg} = h * t_1 + (1-h) * (t_1 + t_2)$$

case-3:- process when cache miss, a block of data is transferred.

$$T_{avg} = h * t_1 + (1-h) * (t_1 + t_B)$$

t_B : time for block of data

multi level cache:-

$$T_{avg} = h_1 * t_1 + (1-h_1) h_2 * t_2 + (1-h_1)(1-h_2) * m$$

where h_1 : hit rate in L_1 cache

h_2 : " " " L_2 "

t_1 : access time to L_1

t_2 : access time of L_2

m : access time of main memory

$$Q = \frac{28.8E}{20.11E} = 1.4$$

Q. Consider a two-level memory system such that level 1 having hit ratio 0.8 and 5 times faster than level 2 if the average access time is changed or increased by 20% of 50 nsec. Compute the following :-

- What is the access time of level 1 memory (t_1)
- What is the new hit ratio (h_{new})
- What is the % change in hit ratio.

Ans :- $h_1 = 0.8$, $t_1 = \frac{t_2}{5}$

$T_{avg} = 50 \text{ nsec.}$

$T_{avg} = 50 + \frac{20 \times 50}{100} = 50 + 10 = 60 \text{ nsec.}$

$T_{avg} = h_1 \cdot t_1 + (1-h_1) \cdot t_2$

$50 = 0.8 \cdot t_1 + (0.2) \cdot 5t_1$

$\Rightarrow 50 = 0.8 \cdot t_1 + t_1$

$\Rightarrow t_1(1.8) = 50 \Rightarrow t_1 = \frac{50}{1.8} = 27.77 \text{ nsec}$

$T_{avg new} = h_{new} \cdot t_1 + (1-h_{new}) \cdot 5t_1$

$\Rightarrow 60 = h \cdot 27.77 + (1-h) \cdot 5 \cdot 27.77$

$\Rightarrow 60 = h \cdot 27.77 + 138.85 - 138.85h$

$\Rightarrow (138.85 - 27.77)h = 138.85 - 60$

$\Rightarrow h = \frac{111.08}{111.08} = 0.709$

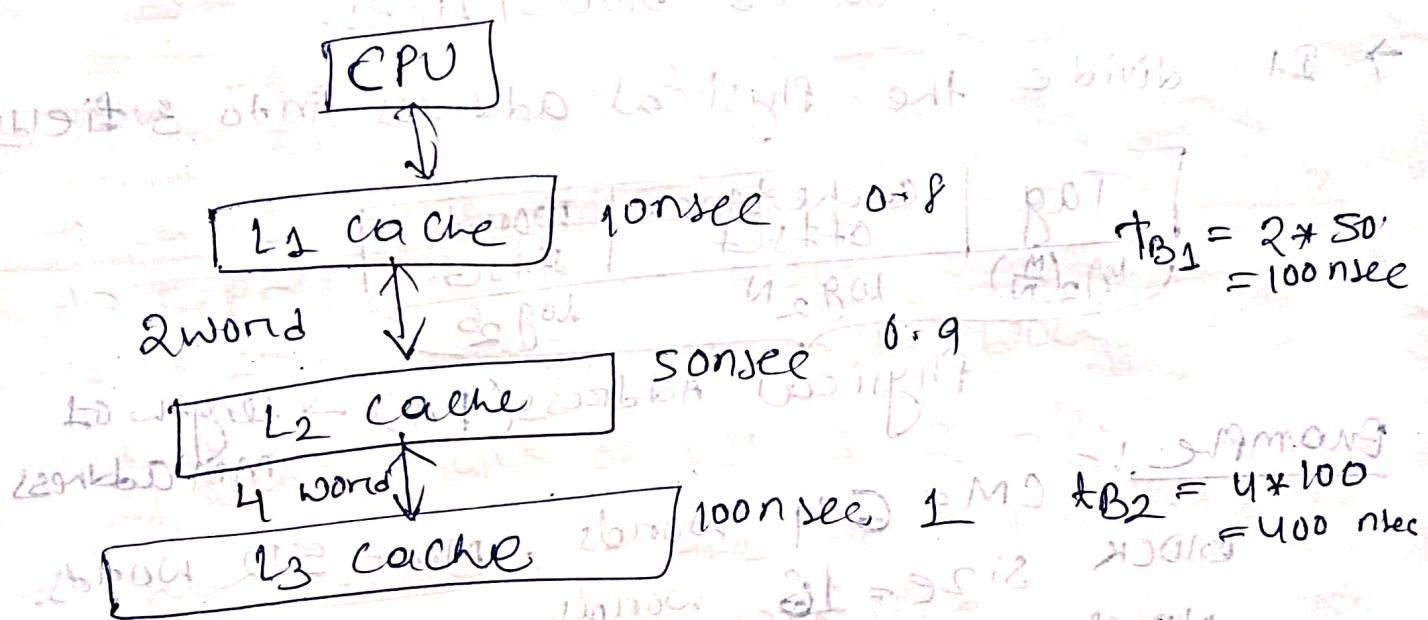
~~$h_{new} = 0.709$~~ $h_{new} = 0.709$ (1ij)

(ii) $h \rightarrow h_{new}$

$0.8 \rightarrow 0.709$
 0.71

$\frac{0.709 - 0.8}{0.8} \times 100 = \frac{-0.09}{0.8} \times 100 = -11.25\%$

Q. Consider a 3-level memory system with access time on word is 10 nsec, 50 nsec, 100 nsec. Hit ratios are 0.8, 0.9 and 1 if referenced word is not available in L-1. Get the two word block from L-2 to L-1 and supply the desired word to the processor. If it is not available in L-2 then get a 4 word block from L-3 to L-2 and transfer the associative block pointers to L-1 and handover the desired word to processor L-1. What's the average access time.



$$\begin{aligned}
 T_{avg} &= h_1 * t_1 + (1-h_1) * h_2 (t_1 + t_{B1}) + \\
 &\quad (1-h_1)(1-h_2) (t_1 + t_{B1} + t_{B2}) \\
 &= 0.8 * 10 + (0.2) * 0.9 * (10 + 100) + \\
 &\quad 0.2 * 0.1 * (10 + 100 + 400) \\
 &= 8 + 0.18 * 110 + 0.02 * 510 \\
 &= 8 + 19.8 + 10.2 = 38 \text{ nsec.}
 \end{aligned}$$

Mapping Function or Address Mapping

- It is a technique to place main memory block into cache memory.
- 3 type of mapping functions.
 - * Direct mapping
 - * Associative "
 - * Set-associative "

Direct Mapping:-

→ Selects $(K \bmod N)^{th}$ cache position for accommodating K^{th} main memory block and N is the no. of cache blocks/lines.

→ It divide the physical address into 3 fields

Tag	cache block offset	word offset
$\log_2(\frac{M}{N})$	$\log_2 N$	$\log_2 B$

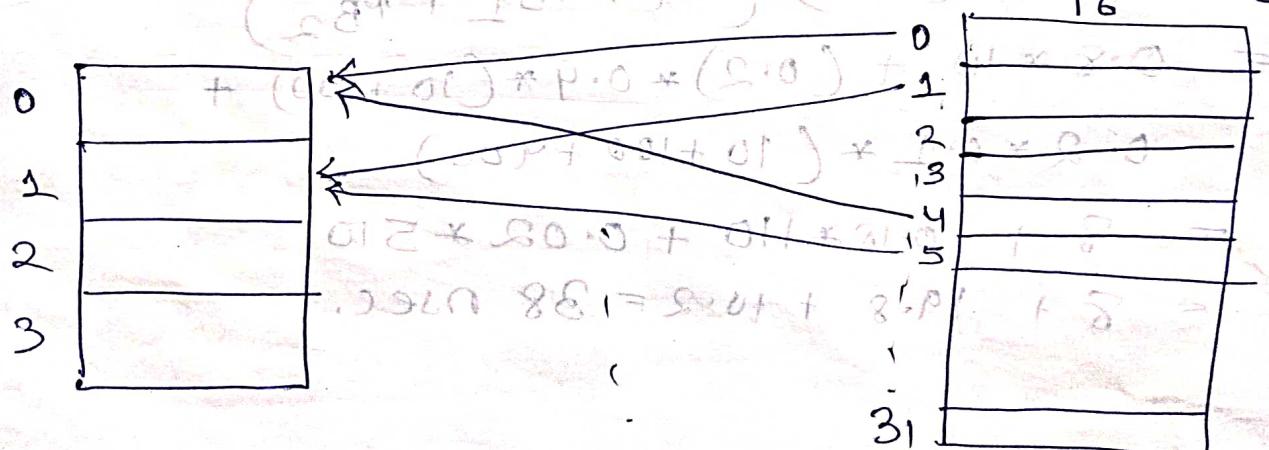
Physical Address (PA) → length of MM address

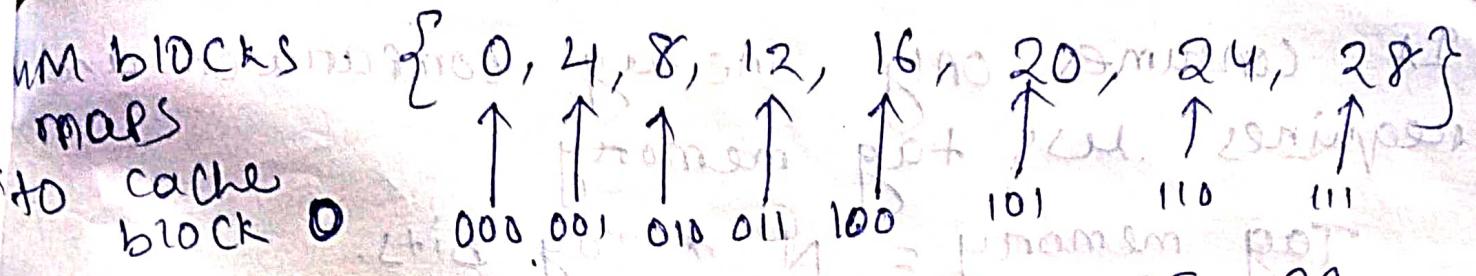
Example:-

CM = 64 words, MM = 512 words.
BLOCK size = 16 words.

NO. of cache blocks = $\frac{64}{16} = 4$ - ~~blocks~~ ~~16~~ ~~64~~

NO. of main memory blocks = $\frac{512}{16} = 32$





Total 8 blocks 1, 5, 9, 13, 12, 21, 25, 29

maps to one C.B. 2, 6, 10, 14, 18, 22, 26, 30

= 3 bit require 3, 7, 11, 15, 19, 23, 27, 31

~~Tag~~ → Tag specifies how many comparators for one cache block.

$$\frac{\text{No. of MM block}}{\text{No. of CM block}} = \frac{32}{4} = 8$$

⇒ 8 comparators means 3 bits requiring so Tag bits = 3.

→ cache block offset is the no. of bits required to represent cache block.

In this example, 14 cache blocks means cache block offset is 2 bit.

→ word offset is the no. of bits required to represent the block size.

In this example, block size = 16 words.

so word offset is 4 bit.

Position	001	100	101	110	111
Tag	3	2	1	0	1
C.B					
Offset					
Word					
Offset					

→ It consumes only one tag comparator and requires less tag memory.

$$\text{Tag memory} = N * \text{Tag bits.}$$

→ No. of cache blocks.

→ The conflict penalties makes more blocks movement and hence the direct mapping is the slowest mapping technique.

Q. Check whether the word 313 of main memory is present in cache or not. \rightarrow
(Same data of previous example)

Block size = 16 words.

M.M size = 512 words

C.M size = 64 words

313rd word belongs to 19th block no.

$$\frac{313}{16} = 19.56$$

① M.M block 19 maps to $19 \div 4 = 3$
cache
No. C.B blocks no.

② Tag bit will be 100 (see the previous page, 19 comes under 100)

OR

$313_{10} =$

100	11	1001
Tag	C.B offset	word offset

$$(313)_{10} = 100111001$$

Divide the bits into 3 fields (indicated in previous page)

If tag bits of cache block no. 3 is 100 then 313 word of M.M is present in Cache memory.

8. Check whether 170 word of MM is present in cache or not.

$$(170)_{10} = \boxed{010101010}$$

Tag of word of first

If Cache block No. 2 is 010 then 170 word of M.M is present in cache

Limitation of Direct mapping :-

- Due to $(K \bmod N)^{th}$ rule, the direct mapping may result conflict penalties.
- Consider a cache with 4 blocks, which is initially empty has to support the main memory block references.

MM : 0, 8, 0, 4, 8, 0, 8, 4, 12, 4, 12
ref. M

0	0	8
1		
2		
3		

0 4 8

Hit = 0.

Other blocks
are empty
still can't put block.

→ It is called conflict

penalty and the movement of block is increase

- To overcome this problem, a technique is suggested called Associative mapping.

Associative mapping:-

- avoids the conflict penalty by relaxing the mapping rule.
- Any block of main memory can be placed anywhere in the cache.
- with the previous reference

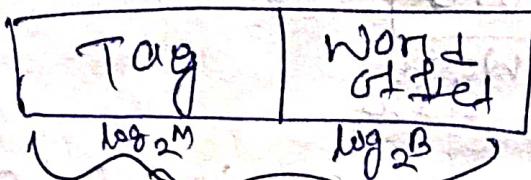
0
8
12

0, 8, 0, 4, 8, 0, 8, 4, 12, 4, 12

M M H M H H H H M H H

No. of miss = 4, No. of hit = 7

→ Associative mapping divide the physical address into 2 portion.



physical address.

Tag bits = $\log_2 M$ M : No. of MM blocks

word offset = $\log_2 B$ B : Block size.

Example:- C.M size = 64 words,

M.M Size = 512 words,

Block size = 16 words

Determine tag bits and word offset.

word offset = $\log_2 64 = 6$ M = 512
~~Tag bits = $\log_2 512 = 9$ N = 16
= 32~~

~~Tag = $\log_2 32 = 5$~~ word offset = $\log_2 16 = 4$

→ since the no. of tag bits are more, the associate memory required more tag memory.

→ since all the cache blocks are fetched

for hit or miss, therefore ~~at~~ tag

No. of tag comparitor are required.

→ No. of cache blocks.

Set-Associative mapping: ~~got to do with~~

- contains the advantages of Direct mapping and associative mapping.
- cache is divided into sets, and the no. of blocks per set depends on the association used.

$S = \text{no. of set in a cache block}$

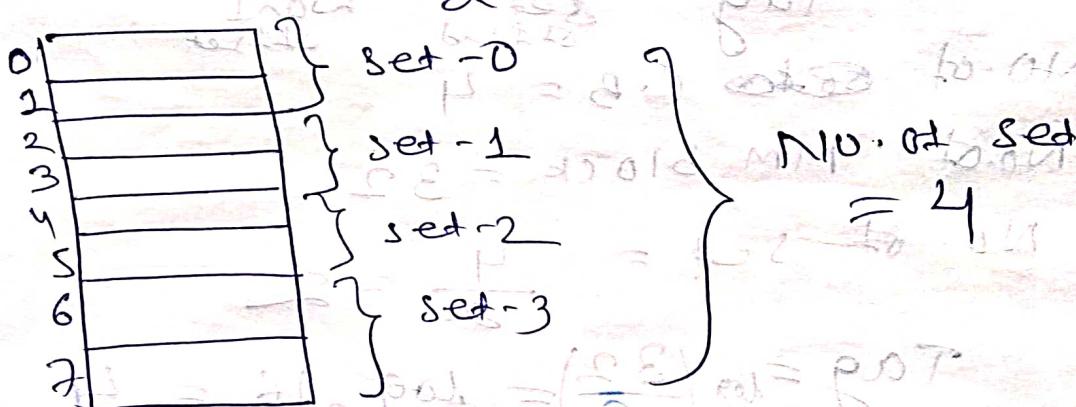
$$S = \frac{N}{P} \quad \text{where } N: \text{No. of cache block}$$

p-way set association

Ex: If no. of cache block,

and 2-way set association.

$$\text{No. of set} = \frac{8}{2} = 4$$



→ k^{th} -block of MM is mapped to $(k \bmod S)$ set in cache.

Within the set, it can be placed anywhere.

→ The physical address is divided into 3 partitions.

Tag	Set offset	word offset
-----	------------	-------------

$$\log_2(\frac{M}{S}) \quad \log_2 S \quad \log_2 B$$

Physical address.

→ The no. of tag comparators in set-associative mapping is the no. of sets.

If 2-way then 2 no. of tag comparators

If 4-way then 4 no. of tag comparators

Example -

MM = 512 words

CMF = 64 words

DS = 16 words

2-way set associative,
calculate tag, set offset and word offset

4	1	4
Tag	Set offset	Word offset

No. of C.B. = 4

No. of MM block = 32

$$\text{No. of sets} = \frac{4}{2} = 2$$

$$\text{Tag} = \log_2 \left(\frac{32}{2} \right) = \log_2 16 = 4$$

Questions On mapping functions.

(9)

Q) A cache consists of 128 blocks of 16 words each, for a total of 2048 (2K) words.

Main memory is addressable by a 16-bit address. The main memory has 64K words.

Obtain Tag bits in different mapping technique.

Ans 1 - M.M Size = 64K words, Block size = 16 words

$$\text{No. of MM blocks} = \frac{64K}{16} = \frac{2^{16}}{2^4} = \frac{2^{12}}{2^4} = 4K$$

C.M Size = 2048 words = 2K

$$\text{No. of Cache blocks} = \frac{2048}{16} = \frac{2^{11}}{2^4} = \frac{2^7}{2^4} = 128$$

Direct mapping

5	7	4
---	---	---

$$4K = \frac{2^{12}}{128} = \frac{2^{12}}{2^7} = 2^5$$

Associative mapping

12	4
----	---

$$4K = 2^{12}$$

Set-associative mapping: 2-way set associative

$$\text{No. of set} = \frac{128}{2} = 64$$

$$\frac{4K}{64} = \frac{2^{12}}{2^6} = 2^6$$

6	6	4
---	---	---

Tag Set offset Word offset



A cache consist of 64 blocks and MM consist of 4K blocks, each block consists of 128 words.

- (a) find size of MM.
- (b) find size of CM
- (c) find length of main memory address
- (d) how many bits are there in each of the tag in case of direct mapping.
- (e) " in case of associative mapping
4-way set associativity
- (f) " "

(3) consider 2GB main memory and 1MB cache partitioned into 256 word blocks. Each word consumes 32 bits.

- (i) how many bits are needed to physical address.
- (ii) how many blocks are present in MM and CM.
- (iii) Identify no. of tag bit, tag comparison and tag memory for the following
 - (i) Direct mapping
 - (ii) Associative mapping
 - (iii) 4-way set associative mapping?

SOL - Hit (memory ~~must~~ be convert into words).

Performance Consideration

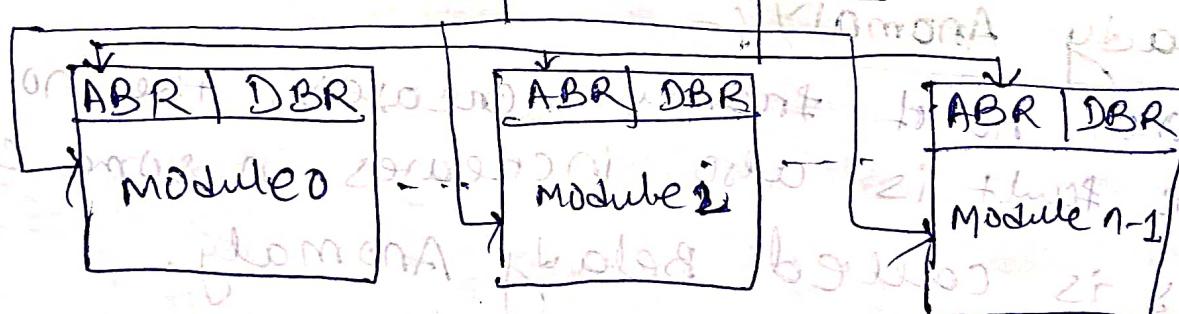
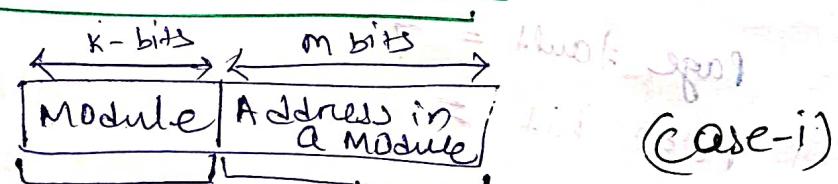
Objective: Best possible performance at the lowest cost.

- To improve performance without increasing the cost.
- Measure is the price or performance ratio.
- To achieve high performance, parallelism is used in slower unit (memory). This parallelism can be achieved organized by memory interleaving.

Memory interleaving:-

- There are two methods of address layout.
 - consecutive words in a module
 - consecutive words in consecutive modules.

Consecutive words in a module:-



What is ABR and DBR :-

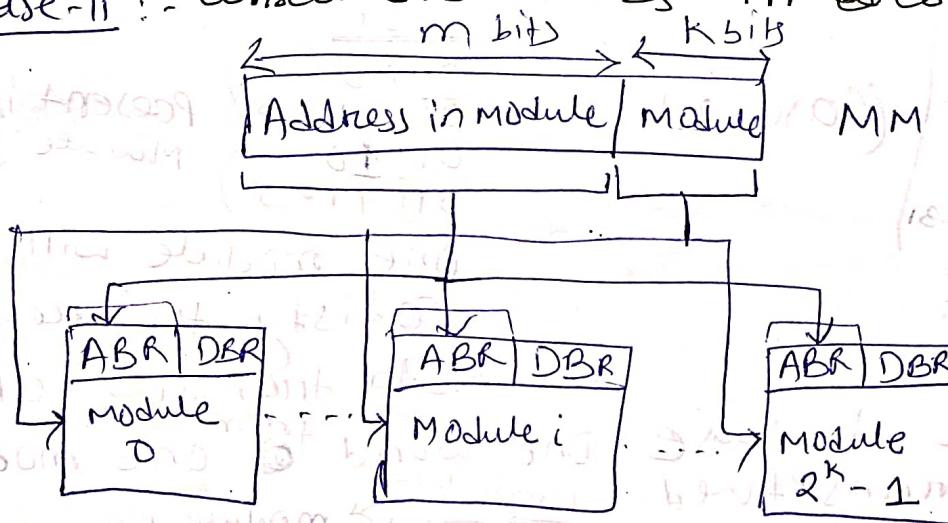
If the main memory of a computer is structured as a collection of physically separate modules, each with its own address buffer register (ABR) and data buffer register (DBR), memory access operation may proceed in more than one module at the same time.

(12) MS, the aggregate rate of transmission of words to and from the main memory can be increased.

Case-i: consecutive words in a module.

- The high order K bits points to the n modules.
- The low order m bits points to a particular word in that module.
- When consecutive locations are accessed, ~~one~~ a block of data is transferred to a cache, only one module is involved.
- At the same time, devices with direct memory access (DMA) ability may be accessing information in other memory address.

Case-ii :- consecutive words in ~~one~~ consecutive module



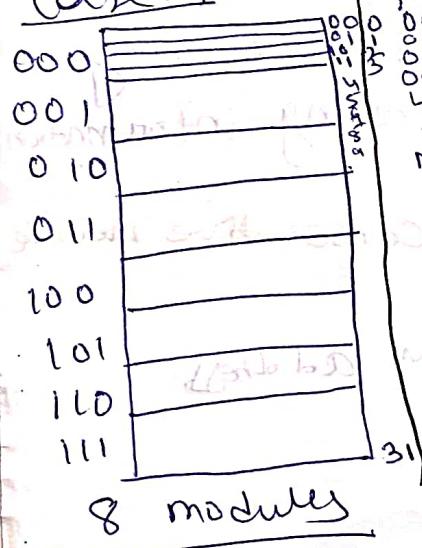
- This way of memory addressing is called memory interleaving.
- The low order ~~addres~~ K bits of the memory address select a module.
- The high order ~~one~~ m bits points to a location within that module.
- In this way, the consecutive addresses are located in successive modules.
- ~~One~~ Any components of the system that generate requests for access to consecutive memory locations can keep several modules busy at any one time.

→ This results in both faster access to a block of data and higher average utilization of the memory system as a whole.

→ To implement the interleaved structure, there must be 2^k modules, otherwise there will be gaps of nonexistent location in the memory address space.

Examples of case-i and case-ii is as follows:

Case-i



Block size = 4 words
Mem size = 32 words
No. of module = $32/4 = 8$
particular word in memory module

(Case-ii)

Suppose Request

$$\begin{array}{l} 01100 = 12 \\ 01101 = 13 \\ 01110 = 14 \\ 01111 = 15 \end{array} \left. \begin{array}{l} \text{present in} \\ \text{module 3 (01)} \end{array} \right\}$$

One module will satisfy the request

But data transfer takes one word from one module

particular word → module No.

00000 → to word 0

00001 → to word 8

00000 → to word 16

00000 → to word 24

in MODULE 0

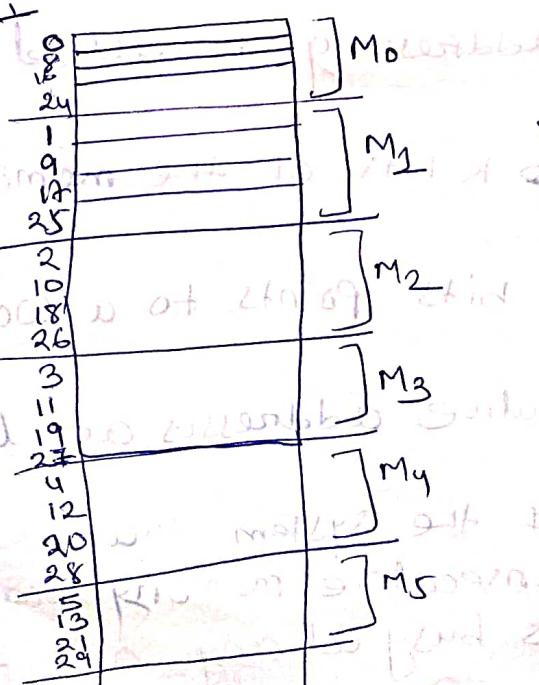
00 001 = 1

01 001 = 9

10 001 = 17

11 001 = 25

Case-ii



in MODULE 1

00 001 = 1

01 001 = 9

10 001 = 17

11 001 = 25

Suppose Request

12, 13, 14, 15

01100 01101 01110 01111
Mu M₅ M₆ M₇

Present in Mu, M₅, M₆, M₇. memory module

At same time these data ~~can~~ can be transferred because present in different module.

Answers

problem-1 :-

Let a cache with 8-word block. A read miss occurs so that entire block is to be loaded. Find the time required to access the block from memory without Interleaved and with Interleaved (4 module).

If the hardware properties are:-

- it takes one cycle to send address to MM.
- first word of the block will be transferred by 8 clock cycles. and other words will take 4 clock cycle using ~~fast~~ page mode.
- 1 clock cycle is needed to send one word to the cache memory.

Ans :-

Without Interleaved :-

Send address to MM = 1 CLK cycle.

1st word access to DBR = 8 CLK.

Rest 7 word access to DBR and simultaneously 1st 7 words from DBR to cache = 7*4 = 28 CLK cycle

Load word after access transfer from DBR to cache = 1 CLK

38 CLK cycles

Interleaved :-

1st word address send to 4 module of MM = 1 CLK

1st word access simultaneously from 4-module to corresponding DBR = 8 CLK cycles.

2nd word access simultaneously from 4 module to corresponding DBR = 4 CLK cycles.

from 4 DBR one by one to cache memory

Transfer = 4 CLK

(Reduces the block transfer time by more than a factor of 2)

Total = 17 CLOCK cycles

The time required to load the desired block into cache is

Problem

① with interleaving:

1 clk cycle \rightarrow to send address to all modules

in M.M. 1st word all 4
8 clk cycle \rightarrow to send data from modules to DBR.

\rightarrow These ~~block~~ words are transferred to the cache one word at a time, during the next 4 clock cycles.

\rightarrow During this time, the next word in each module is accessed and stored into DBR.

\rightarrow Then it takes another 4 cycles to transfer these words to the cache.

\rightarrow Therefore the total time required

$$1 + 8 + 4 + 4 = 17 \text{ clock cycles}$$

\rightarrow Thus interleaving reduces the block transfer time by more than a factor of 2.

Without interleaving

1st CLK: Send address to MM.

9th CLK: 1st word data present in the DBR

After 10th CLK: Data placed into the cache

After 13th CLK: 2nd word data present in the DBR

After 14th CLK: Data placed into the cache

After 18 \rightarrow 3rd complete; place in cache

22 \rightarrow 4th

26 \rightarrow 5th

30 \rightarrow 6th

34 \rightarrow 7th

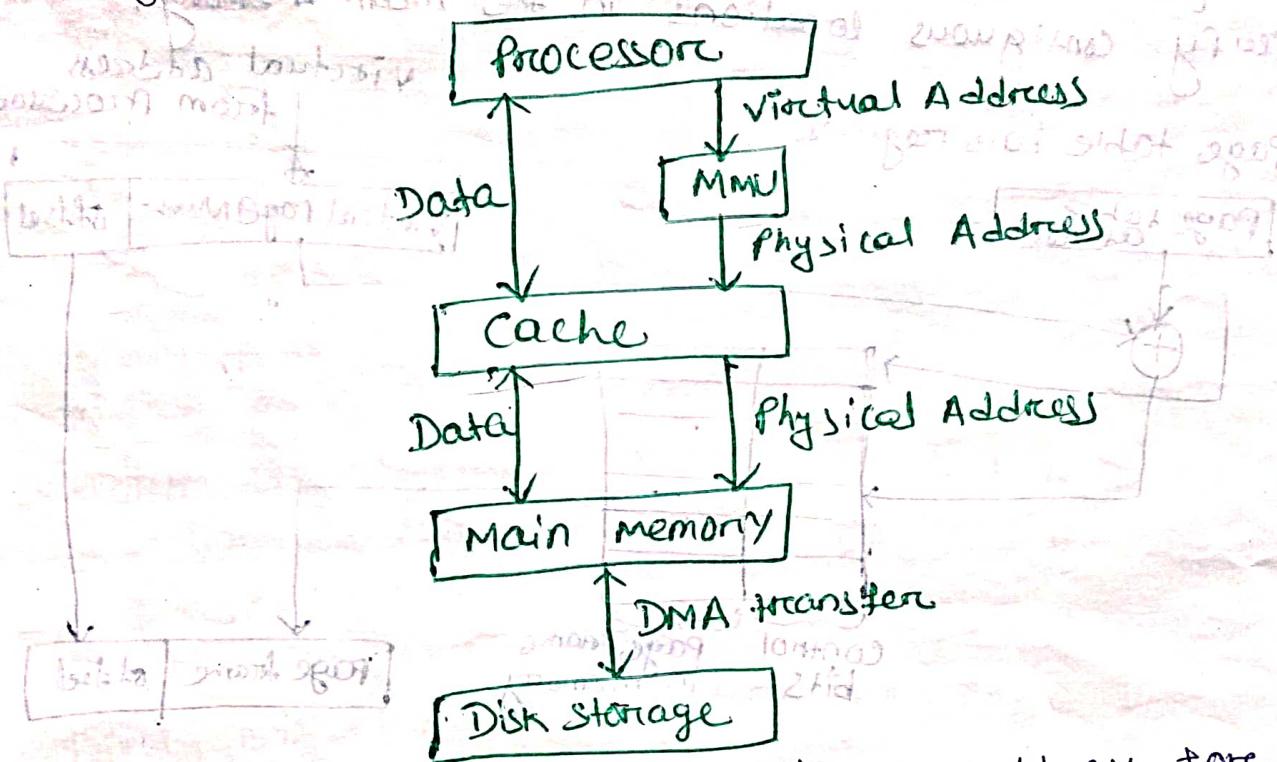
38 \rightarrow 8th

Virtual Memories :-

- It is a technique that allows the execution of program that may not be completely in memory i.e. program size can be greater than the physical memory size.

Advantage:

- Program size can be ~~more~~ larger than the physical memory.



- The processor issues the binary address for the instruction and data is called virtual address. This address is independent to the available physical main memory space.
- Virtual address is an imaginary address generated by the processor.
- A memory management unit (MMU) is a special hardware unit, translates virtual address into physical address.
- The desired data are in the main memory, these data are fetched.
- If the data is not present in the main memory, the MMU causes the operating system to bring the data into the memory from disk.

Transfer of data between the disk and the main memory is performed using the DMA scheme.

Address translation :- (Paging)

→ translation of virtual addresses into physical address.

→ page: All program and data are composed of fixed length units is called pages.

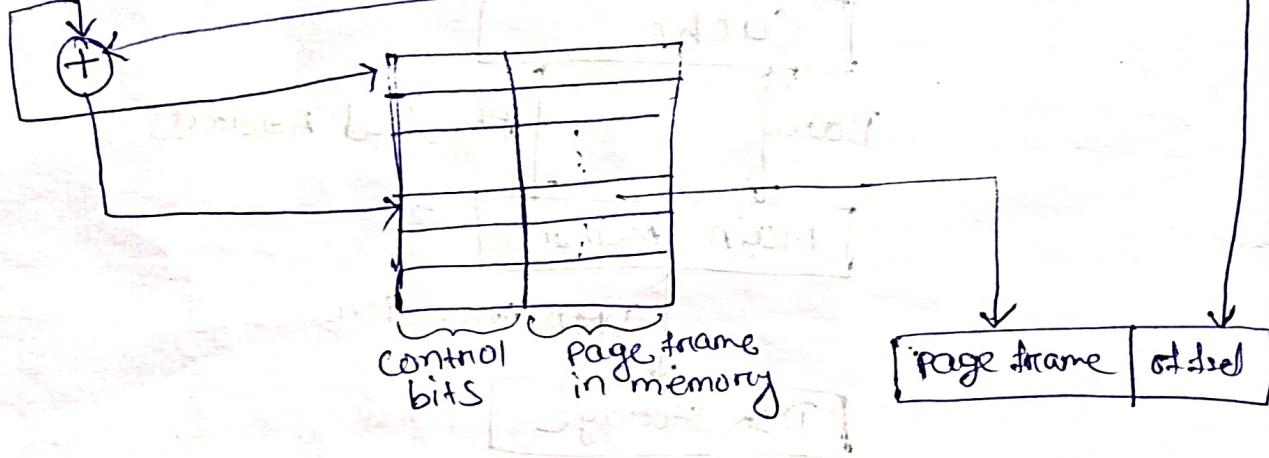
→ Page is consists of a block of words that occupy contiguous locations in the main memory.

Page table base register

Page table address

virtual address
from processor

Virtual page Number offset



→ The cache bridges the speed gap between the processor and the main memory, and is implemented in hardware.

→ The virtual memory mechanism bridges the size and speed gaps between the main memory and secondary storage and is implemented in part by software techniques.

→ Processor generates virtual address, a virtual page number (high-order bits) followed by an offset (low-order bits) that specifies the location of a particular byte (or word) within a page.

→ Page table: Contain the information about the main memory locations of each page. This information includes the main memory address where the page is stored and the current status of the page.

→ Page frame: is a main memory location that can hold one page.

→ Page table base register: The starting address of the page table is kept in a page table base register (PTBR).

→ The corresponding entry in the page table is obtained by adding the virtual page number to the contents of PTBR.

→ The contents of this location give the starting address of the page if that page currently resides in the main memory.

→ Valid bit indicates the validity of the page that is whether the page is actually loaded in the main memory.

→ This bit allows the OS to invalidate the page without actually removing it.

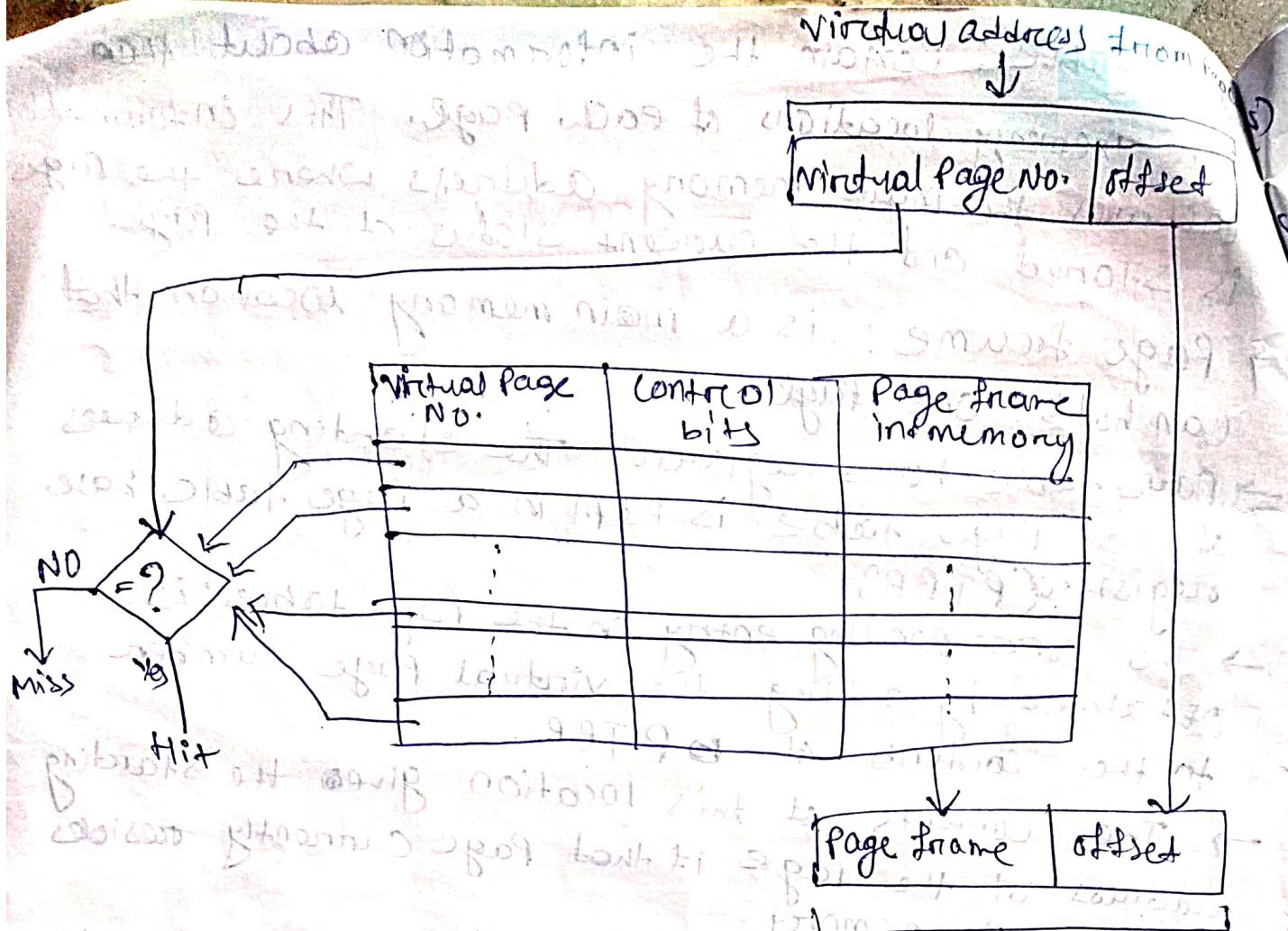
→ Dirty bit indicates whether the page has been modified during its residency in the memory. This information is needed to determine whether the page should be written back to the disk before it is removed from the main memory to make room for another page.

→ Page table is kept in the main memory.

→ A small cache, called the Translation look aside buffer (TLB) is incorporated into the MMU.

It contains the page table entries that correspond to the most recently accessed page.

→ Associative mapping technique is used in TLB.



Address translation procedure:

- Given a virtual address, the MMU looks in the TLB for the referenced page.
- If the page table entry for this page found in the TLB, the physical address is obtained immediately. (TLB hit).
- If there is Miss in TLB, then the required entry is obtained from the page table in the main memory and TLB is updated. (TLB miss then searching in page table)
- If the requested page is not in the main memory, a page fault is ~~said to have~~ occurred.
- The whole page must be brought from the disk into the memory before access can proceed.

- (ii) When a page fault occurs, the MMU asks the OS to intervene by raising an exception (interrupt).
- (iii) The processing of active task is interrupted and control is transferred to the OS. The OS copies the requested page from the disk into the main memory and returns control to the interrupted task.
- (iv) The execution of interrupted instruction must continue from the point of instruction or the instruction must be restarted.

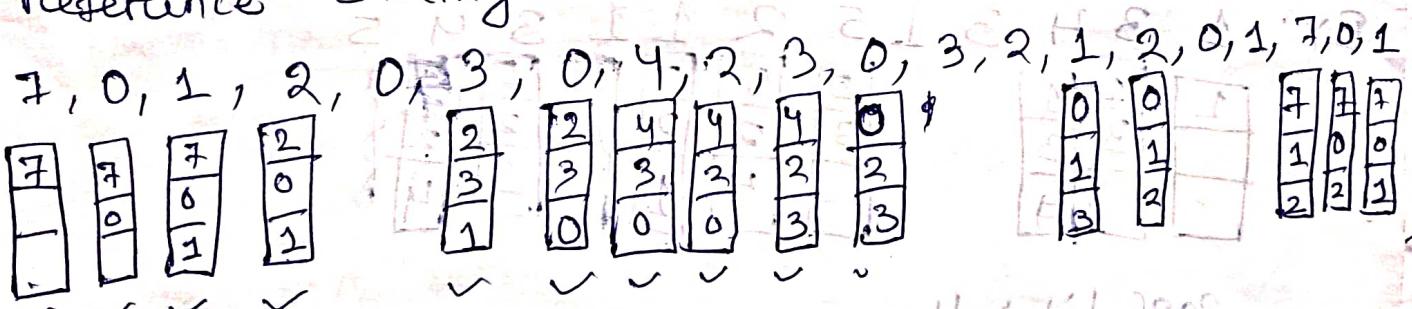
Replacement Algorithm:-

→ If a new page is brought from the disk and the main memory is full, the replacement algorithm selects a victim page in the main memory and replace the victim page with the new page from the disk.

i) FIFO page replacement algorithm:-

→ The page that will come first will be replaced first.

Ex: Find the no. of page fault for the following reference string with frame size = 3

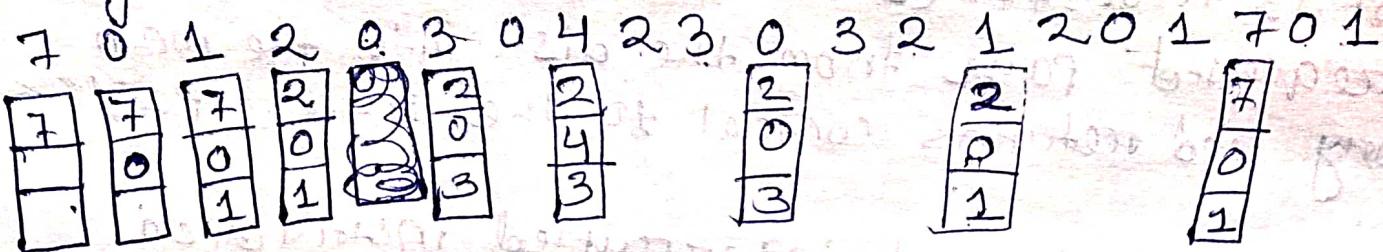


Page faults = 15.

$$\text{page hit} = 5 \quad \text{hit rate} = \frac{5}{20}$$

(i) Optimal Page replacement algorithm :-

→ Replace the page that will not be used for the longest period of time. (future will check).

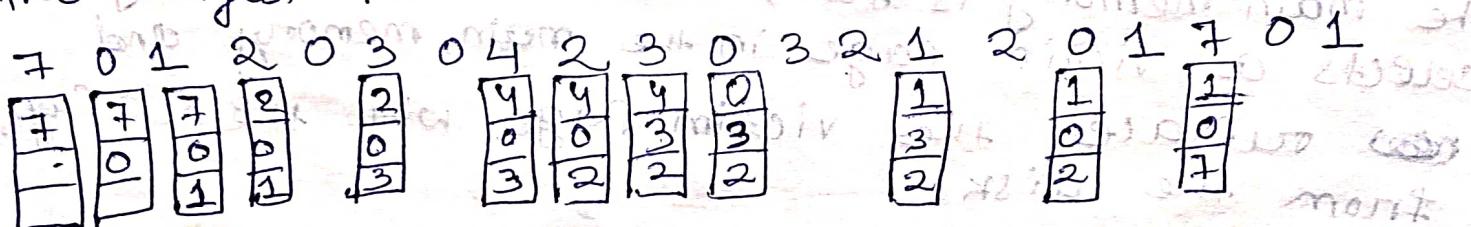


No. of Page fault = 9

Page hit = 11 hit rate = 11/20

(ii) Least Recently Used (LRU) Algorithm :-

→ Replace the page that has not been used for the longest Period of time.

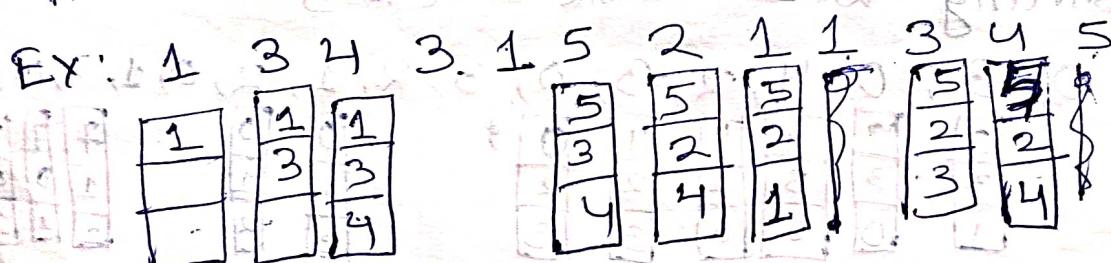


No. of Page fault = 12 hit rate = 8/20.

Page hit = 8

(iv) Most frequently used (MFU) Algorithm :-

→ The block which was used maximum no. of times that block will replace first.



Page hit = 4

Page fault = 8

hit rate = $\frac{4}{12}$

Least Frequently Used (LFU) :-

The block which was used less no. of times will be replaced first.

Ex:	1	3	4	3	1	5	2	1	1	3	4	5
	1	1	1		1	1	1	1	1	1	1	1
	3	3	3		3	3	2		3	3	4	5
	4				5	2						

$$\text{Page hit} = 5$$

$$\text{Page fault} = 7$$

(vi) Most Recently Used (MRU) :-

The block which was used in near past must be replaced first.

1	3	4	3	1	5	2	1	1	3	4	5
1	1	1		1	1	1	1	1	1	1	1
3	3	3		3	3	2		3	3	4	5
4				4	3	4		4			

$$\text{Page fault} = 7$$

$$\text{Page hit} = 5$$

Belady Anomaly :-

When no. of frames increases, the no. of page fault is also increases in some case.

This is called Belady Anomaly.

Example: In FIFO it occurs.

1, 2, 3, 4, 1, 2, 5, 1, 2, 3, 4, 5

1	2	1	4	4	4	5	5	5	3	3	2	9
2	2	2	2	1	1	1	1	1	2	2	2	
3	3	3	3	3	2	2	2	2	3	3	3	

1	2	1	2	1	2	3	4	5	4	5	4	10
2	2	2	2	1	1	1	1	1	2	2	2	
3	3	3	3	3	2	2	2	2	3	3	3	