

Unit 4: Memory

Memory: Basic concept and hierarchy, semiconductor RAM memories, 2D & 2 1/2 D memory organization, ROM memories, Cache memories: concept and design issues and performance, address.

Mapping and replacement Auxiliary memories: magnetic disk, magnetic tape and optical disks.

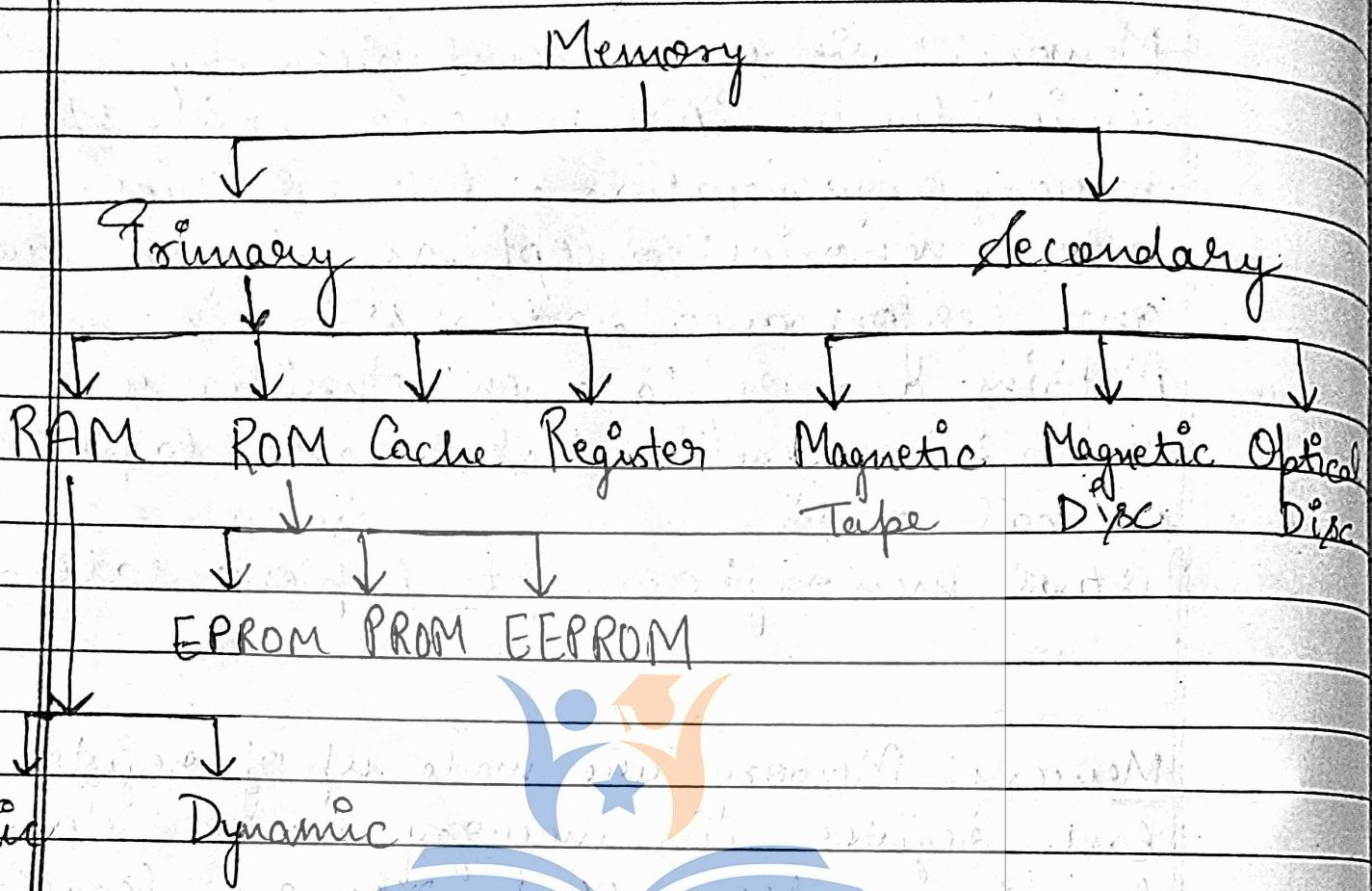
Virtual memory: concept implementation.

Memory: Memory are made up of register. Each register in memory is one storage location also called "Memory location".

Each memory location is identified by an address the no. of storage locations can vary from a few in some memories to hundred of thousand in other.

The total no. of bits that a memory can store is its capacity.

Types of Memory



Primary: It includes ROM and RAM and it is located close to CPU on the computer motherboard enabling the CPU to read data from Primary memory very quickly indeed.

Secondary: It includes physically located with a separate storage device such as a hard disk drive or solid state drive (SSD) which is connected to the computer system either directly or over the network hardware.

RAM or Random Access Memory

Volatile memory

⇒ Memory in which any location can be reached in a short fixed amount of time after specifying its address.

Types

① SRAM

- * It is less memory cell per unit area.
- * Access time is less so faster.
- * Consists of flip flop register.

↳ Cost is more

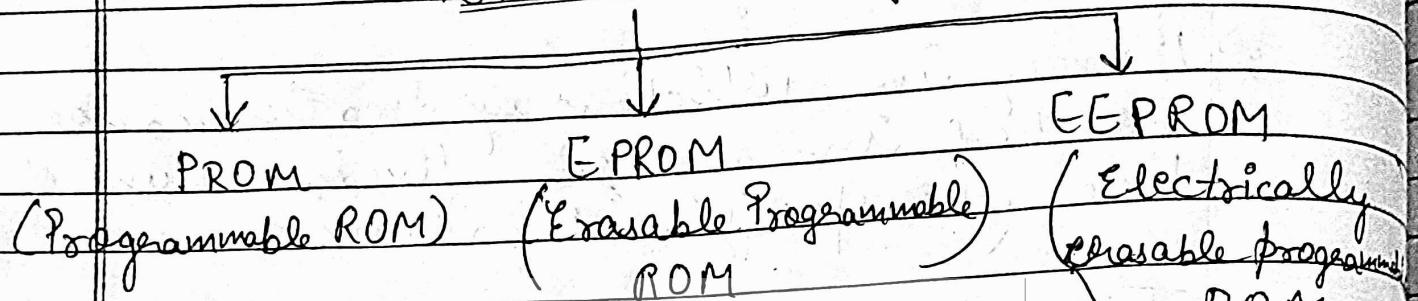
↳ Need less power

DRAM

- * It is more memory cell per unit area.
- * Access time is more than SRAM, so slower.
- * Consists of pair of capacitor.
- * Cost is less.
- * More power is required.

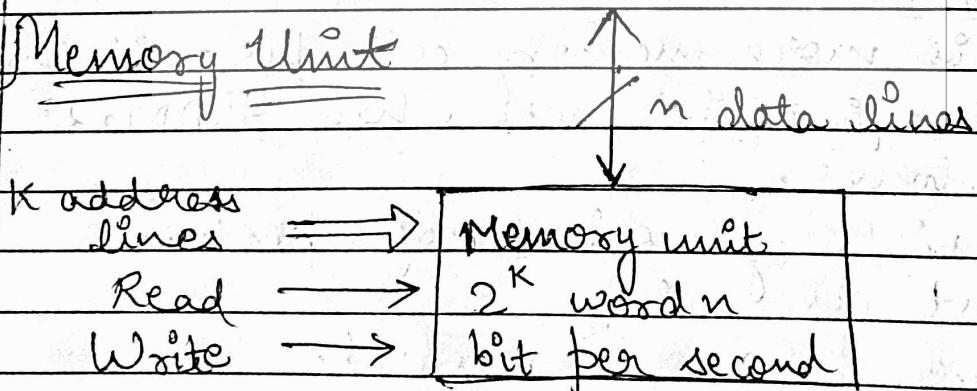
ROM (non-volatile) or (Read Only Memory)

Permanent Memory



- * Instruction can be stored only once
- * It can be stored multiple times
- * Errors in written instruction cannot be changed or erased.
- * Written instruction can be changed or erased with help of UV Ray.
- * Once written, it is not possible to rewrite the instruction.
- * EEPROM ROM has instructions already written in EEPROM.

Memory Unit



The Block diagram of memory unit, then the lines provide the information to be stored in memory and the K address lines specify the particular word.

When ~~it has~~ there are k address lines we can access 2^k memory words.

$$\text{If } k = 10$$

$$2^k = 2^{10} = 1024$$

$$\text{Total memory} = 1024$$

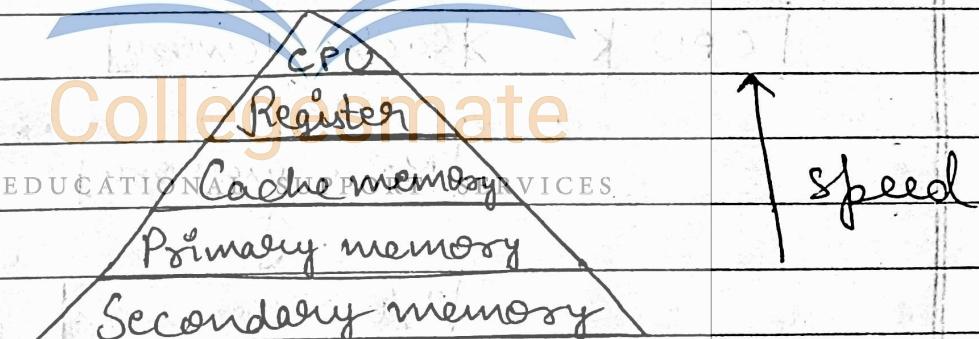
Ques:

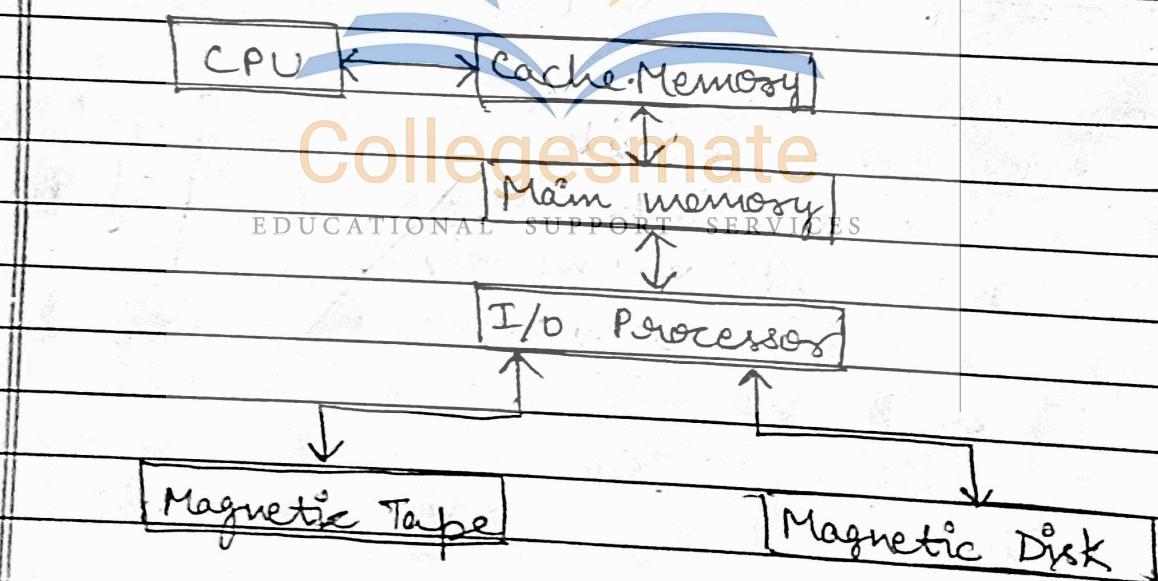
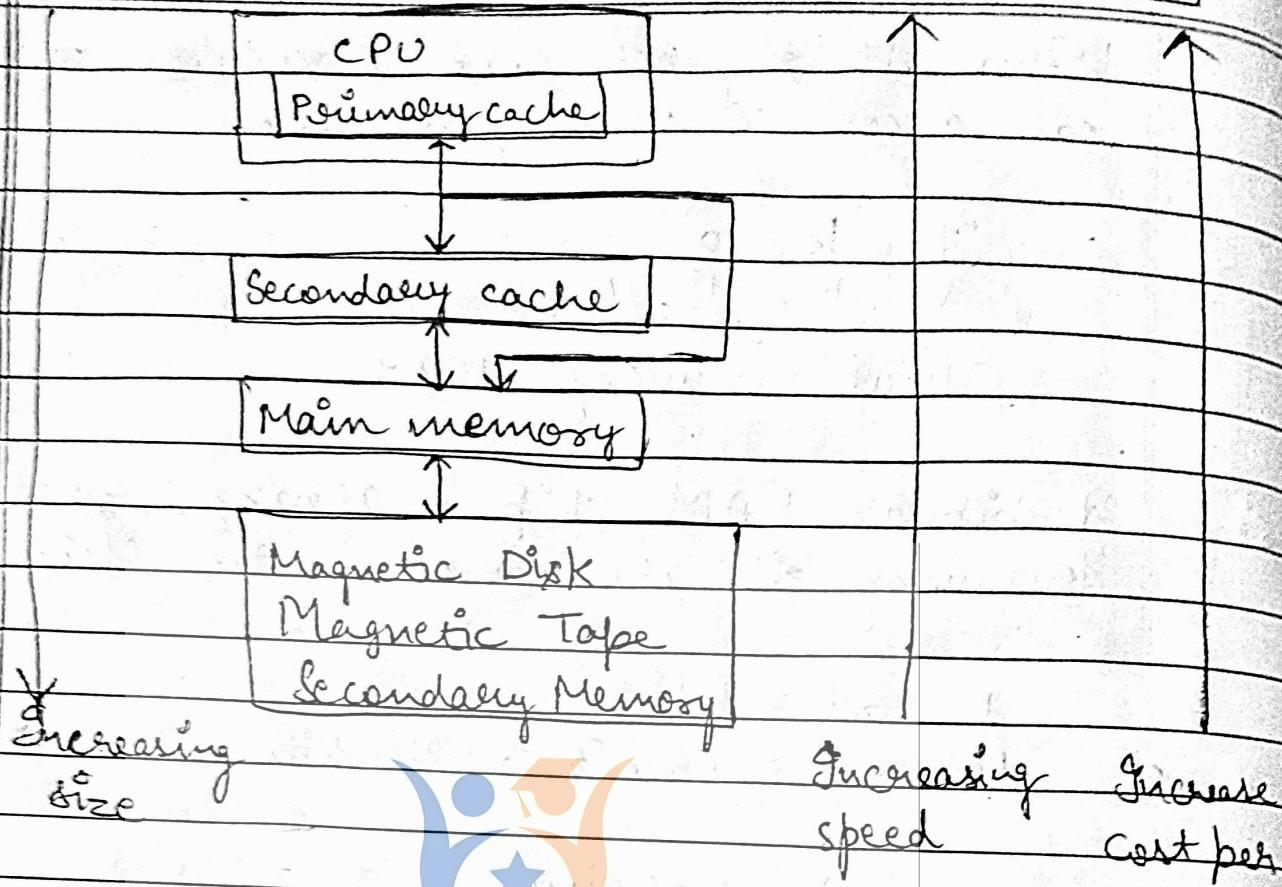
A bipolar RAM chip is arranged as 16 words. How many bits are stored in the chip.

$$\Rightarrow 1 \text{ word} = 8 \text{ bits}$$

$$16 \text{ word} = 16 \times 8 = 128 \text{ bits}$$

Memory hierarchy





Computer memory should be fast, large and inexpensive.

Increased size, speed are achieved at increase cost.

Very fast memory system can be achieved if SRAM (Static RAM chips are used).

Processors fetch decode and data from the main memory to execute the program.

The DRAM which form the main memory are slower device.

The program which is to be executed is loaded in the main memory.

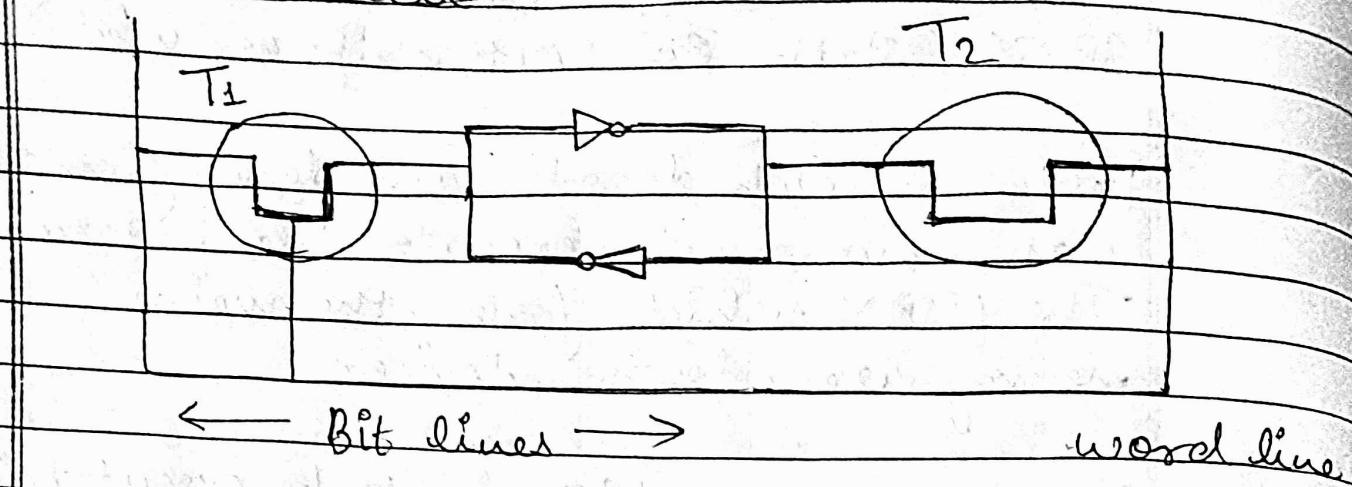
The cache memory just discussed is called secondary cache, most processor have the built in cache memory called primary cache.

Semiconductor RAM

EDUCATIONAL SUPPORT SERVICES

Memory that consist of circuit capable retaining these states as long as power is applied are known as static memory. These are Random Access Memory so called static RAM.

Static RAM cell



This figure shows the implementation of SRAM cell. It consists of 2 cross clips and 2 transistors T_1 and T_2 which act as switch. The latch is connected to two bit lines transistor T_1 and T_2 .

When the word line is zero level, the transistors are off and the latch is stable state in Read operation.

Read op: If $b=1, b'=0$ then T_1 transistors is on and T_2 is also on and word line off.

Write op: Only word line is activated and transistor T_2 is activated and $b=0, b'=1$ and T_1 is off.

Page No.	
Date:	

Dynamic RAM

Dynamic RAM store the data as a charge on the capacitor. It contains 1000 of such memory cells. When column and row line is high the mosfet closes and charges the capacitor.

When the column and row lines go low the mosfet opens and the capacitor maintains its charge in this way it stores one bit.

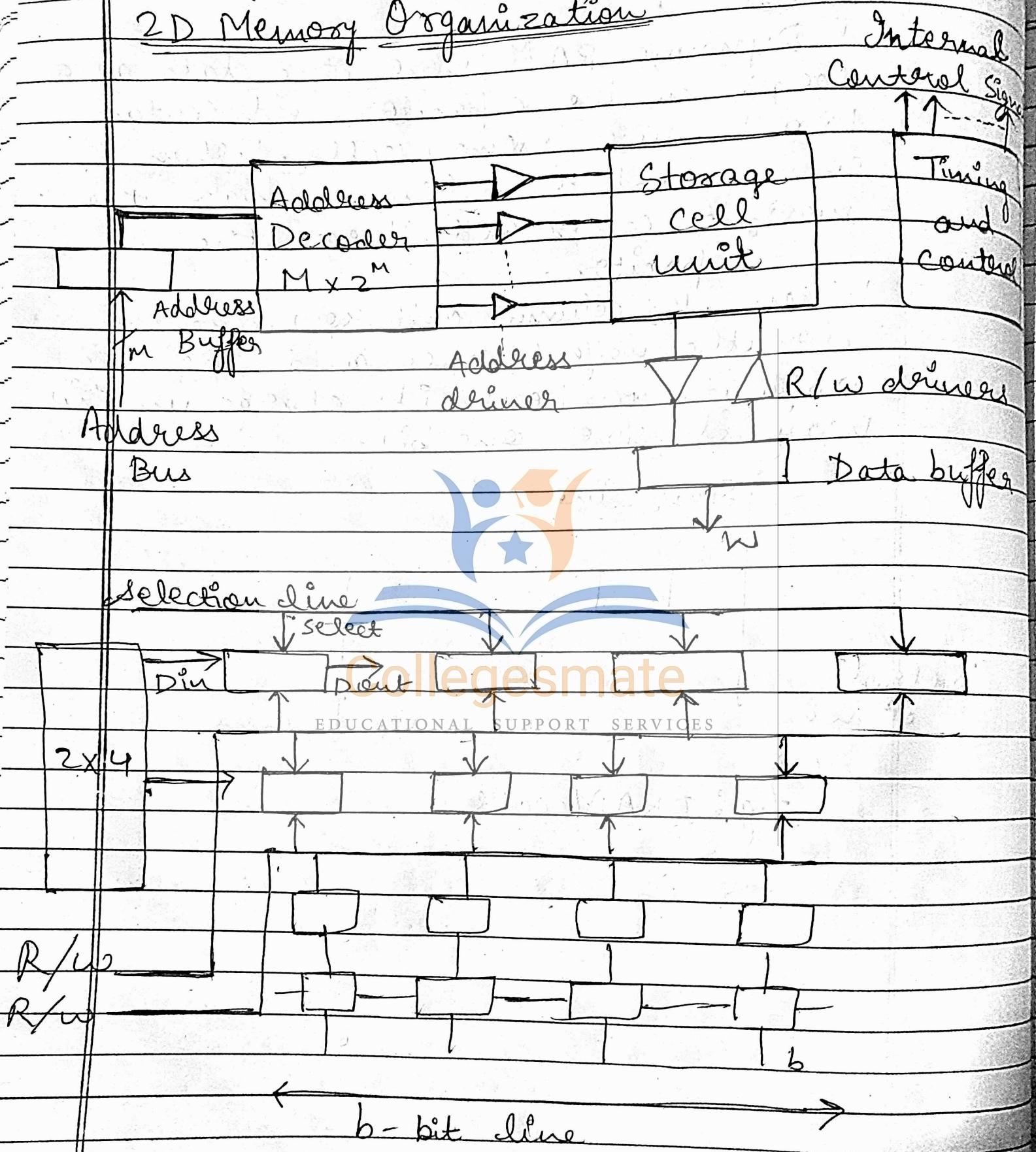
Sense line

EDUCATIONAL SUPPORT SERVICES

Control
line

fig: DRAM cell

2D Memory Organization



This is a simplest organization as shown in fig the cells are organized in the form of a 2D array with row and column.

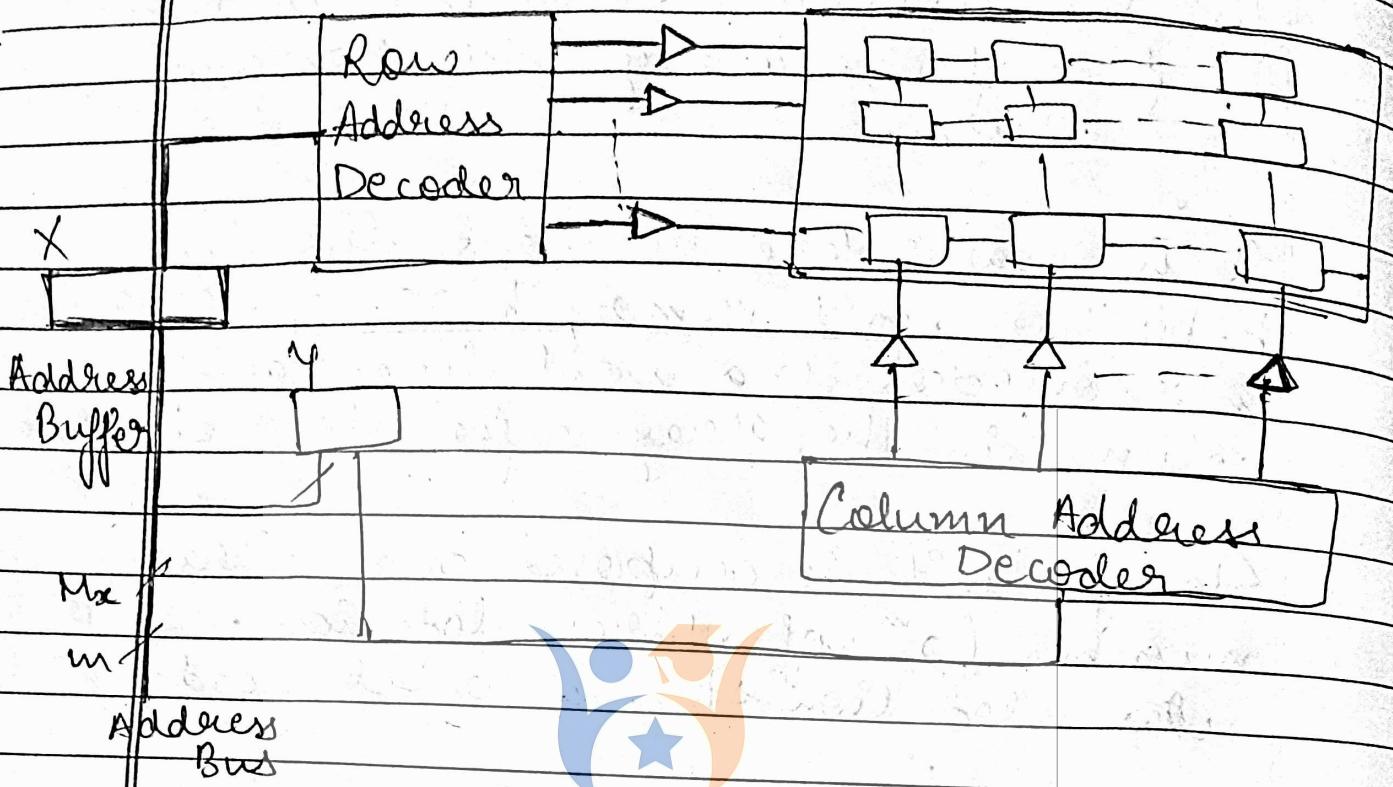
Each row refers to word line for b-bit per word memory b no. of cells are interconnected to a word line. Each column in the array refers to a bit line.

Storage unit is composed of a large number (2^m) of address location. Each location store a w bit word.

An m bit address is generated by the processor. It is stored into an address buffer. Output of address buffer is fed to the address decoder.

The address decoder select one of the address decoder in the memory the output of address decoder passes through byte state buffer.

2 1/2 D Memory

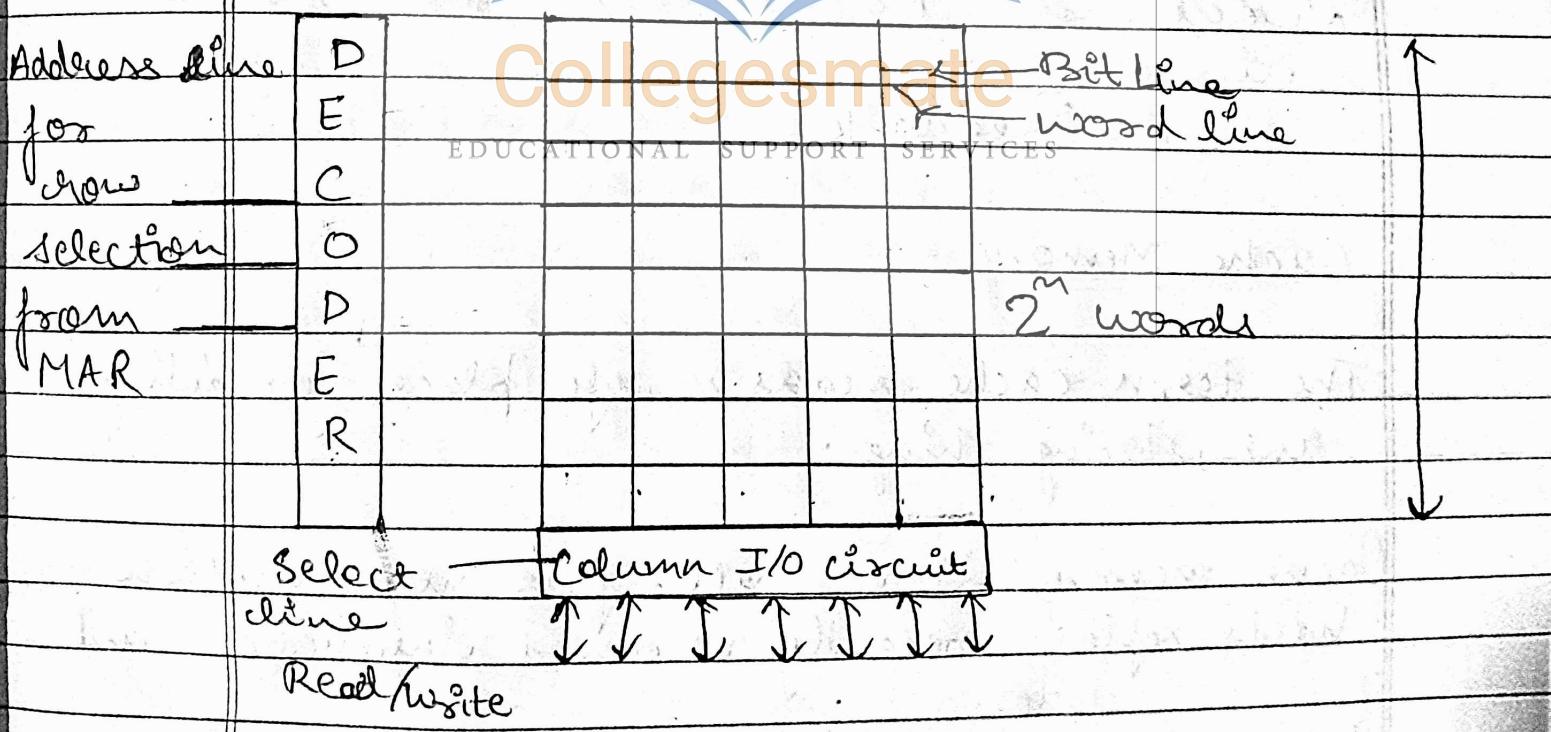


In 2.5D organization the scenario is the same but we have 2 different decoders one is column decoder and another is row decoder. Column decoder is used to select the column and row decoder is used to select the row. The address from the MAR will go in decoder's input.

Decoder will select the respective cell through the bit outline, the data from the location will be written at that memory location.

Read and Write of

1. If the select line is in Read mode then the word/bit which is represented by the MAR that will be coming out to the Data lines and get read.
 2. If the select line is in write Mode then the data from memory data Reg will go to the respective cell which is address by the MAR.
 3. With the help of select line the data will get selected where the read and write of " will take place.



2D Memory Organization

Decoder
Selection
For
Row

Address line
for column
selection

Select
line

Decoder of column
selection

Bit in Bit out

Read/write

EDUCATIONAL SUPPORT SERVICES

Cache Memory

The term cache means a safe place for hiding and storing things.

Cache memory is a small, fast memory which holds copies of recently accessed instruction and data.

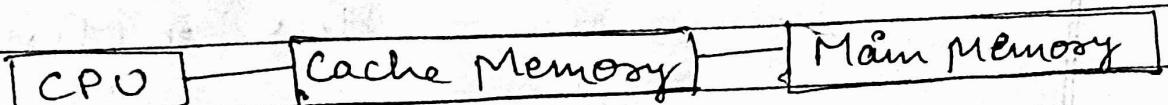
When the processor makes a request for memory reference the request is first sought in the cache.

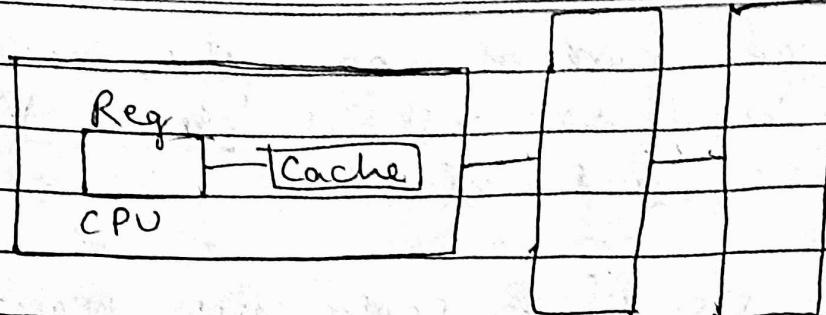
If we get that memory reference which is requested we call it 'Cache Hit' otherwise 'Cache Miss'.

In the case of a cache miss requested element is brought from a subsequent memory level from memory hierarchy and placed in cache.

A block of elements are transferred from main memory to cache memory by expecting that the next requested element will be residing in the neighboring locality of current requested element (spatial locality) and this has to happen under one main memory Access time.

- * Cache memory is organized not in bytes, but as blocks of cache lines with each lines containing some no. of bytes (16-64).
- * Cache lines do not have fixed addresses which enables the cache system to populate each cache line with a unique (non-contiguous) address.





RAM SM (Secondary Memory)
(HD)

Cache Memory 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.

Collegesmate

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

(H)

DR

= no. of hits

Total no. of CPU references

Miss Ratio = $\frac{\text{No. of Miss}}{\text{Hits} + \text{Miss}}$

or

= no. of Miss

Total no. of CPU references

Cache Access time

If cache hit when access the cache that time is called cache access time it is also called (cache hit time).

Time required to access word from the cache.

Miss Penalty

Cache Miss Time Penalty: It is the time required to fetch the required block from main memory.

$$(\text{Cache access} + \text{M M Access time})$$

Average Access time of CPU

$$= \text{hit ratio} \times \text{Cache Access time} + \\ (1 - \text{hit ratio}) \times \text{Miss Penalty}$$

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

Ques: The access time Cache Memory is 100 ns and that of main memory 1000 ns. It is estimated that 80% of memory request are for read and the remaining 20% for write. The hit ratio for read access only is 0.9. A write through procedure is used.

- (a) What is the avg access time of the system considering only memory read cycle?
 - (b) What is the avg access time of the system for both read and write cycle?
 - (c) What is the hit ratio taking into consideration time the write cycles.
- Total avg access time = $0.9 \times 100 + (1 - 0.9) \times (100 + 1000)$

$$\begin{aligned} &= 0.9 \times 100 + (1 - 0.9) \times 1100 \\ &= 200 \text{ ns} \end{aligned}$$

$$\begin{aligned} \text{Total write} &= 1 \times \text{Max}(100, 1000) \\ &= 1000 \text{ ns} \end{aligned}$$

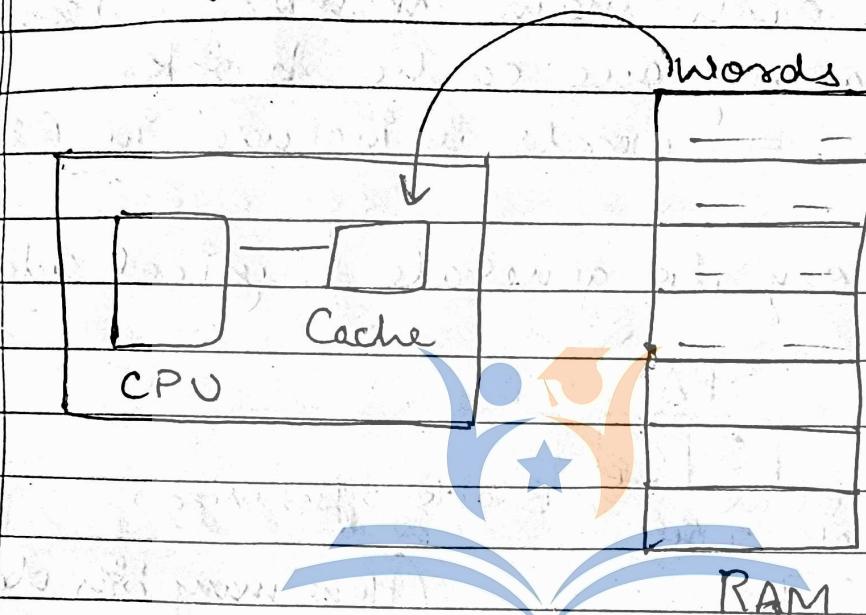
$$\begin{aligned} (d) \quad \text{Avg Access time for Booth read and write} \\ &= 80\% \text{ are read request} + 20\% \\ &\quad \cancel{\text{write request}} \\ &= 0.8 \times 200 + 0.2 \times 1000 \\ &= 160 + 200 \\ &= 360 \text{ ns} \end{aligned}$$

Cache Mapping and its types

Direct
Mapping

Fully
Associative

Set Association



which line contains in the cache memory.
So this is called a cache mapping and
How to fetch the Block or reference of cache
data by CPU.

Direct Mapping

Cache		
L0	B0	B4
L1	B1	B5
L2	B2	B6
L3	B3	B7

W0	W1	W2	W3	B0
W4	W5	W6	W7	B1
W8	W9	W10	W11	B2
W12	W13	W14	W15	B3
W16	W17	W18	W19	B4
:	:	:	:	:
W124	W125	W126	W127	B31

Line size = Block size
 $\frac{1024}{4} = 4$

$\frac{32}{4} = 32$ // 128 Words

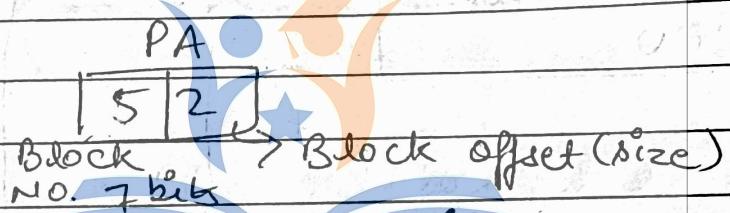
($k \bmod n$ formula used)
 \downarrow
 $0 \bmod 4 = 0$

where k is block number
 n is no. of lines

A direct-mapped cache is the simplest approach each main memory address maps to exactly one cache block.

For example: words is include in Block.

Main memory to generate physical address.



$$0 - 127 = 7$$

$$0 - 3 = 2 \text{ bits}$$

(How many bits to be used to represent a word?)

Address $0001010 = 10$ Word

5 bit (Block no.) $\rightarrow 2$

10 (Block offset) $\Rightarrow 2$ (word)

~~P(A) = 7~~

3	2	2	
Tag	Line no.	(Block offset)	

No. of lines = 4 = ②

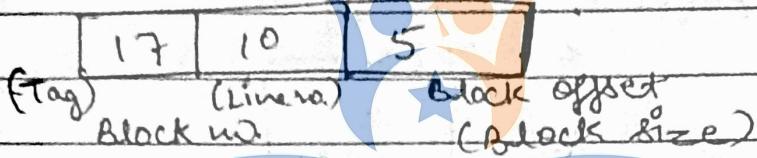
Ex: Cache memory available to represent a

Tag	010	00	00 = 32
			01 = 33
			10 = 34
			11 = 35

$\text{Tag} = 0010$	$\text{Tag} = 001\ 01$	00 20
	B8	01 21
010	L0	10 22
001	L1 + B5	11 23

Ques: Consider a direct mapped Cache of size 32 kB with Block size 32 bytes. The CPU generates 32 bit address. The no. of bits required for Cache Index and tag bits respectively.

$$PA = 32 \text{ bit}$$



Line size = Block size

EDUCATIONAL SUPPORT SERVICES

1023

32 KB

3x8

Fully associative

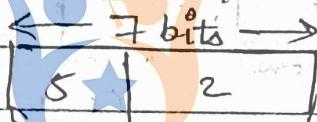
L ₀	B ₁
L ₁	B ₀
L ₂	B ₂
L ₃	B ₃

16 words

W ₀ , W ₁ , W ₂ , W ₃	B ₀
W ₄ , W ₅ , W ₆ , W ₇	B ₁
W ₈ , W ₉ , W ₁₀ , W ₁₁	B ₂
W ₁₂ , W ₁₃ , W ₁₄ , W ₁₅	B ₃
W ₁₂₄ , W ₁₂₅ , W ₁₂₆ , W ₁₂₇	B ₃₁

128 words

Line size = Block size



Block no. (Block offset)

Cache

Collegesmate

EDUCATIONAL SUPPORT SERVICES



Set Associative: It is combination of Direct and fully Association.

L ₀	B ₀		S ₀
L ₁	B ₀		
L ₂	B ₁		S ₁
L ₃	B ₁		

16 words

Cache

W ₀	W ₁	W ₂	W ₃	B ₀
W ₄	W ₅	W ₆	W ₇	B ₁
W ₈	W ₉	W ₁₀	W ₁₁	B ₂
W ₁₂	W ₁₃	W ₁₄	W ₁₅	B ₃
W ₁₂₄	W ₁₂₅	W ₁₂₆	W ₁₂₇	B ₃₁

128 words

$$\frac{16^4}{16} = 1$$

k way set associative are given means no. of set to calculate.

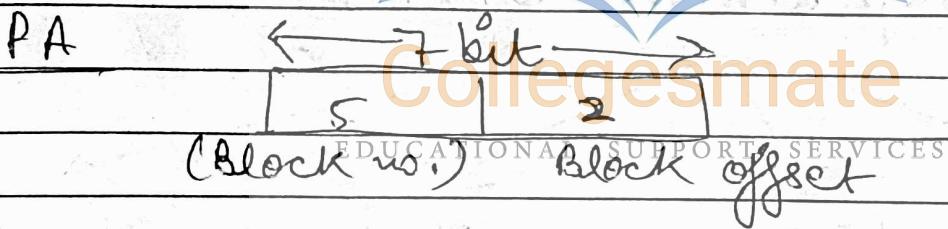
$$\begin{array}{l} \text{2 way} \\ \text{4 way} \\ \text{8 way} \end{array} = \frac{\text{no. of lines}}{k} = \frac{4}{2} = 2$$

Direct Mapping: $k \mod n$ where n
 $\underline{-} 0 \mod 2$ number of set
 $= (\text{Association})$

$$B_1 = 1 \mod 3$$

= ①

set 1

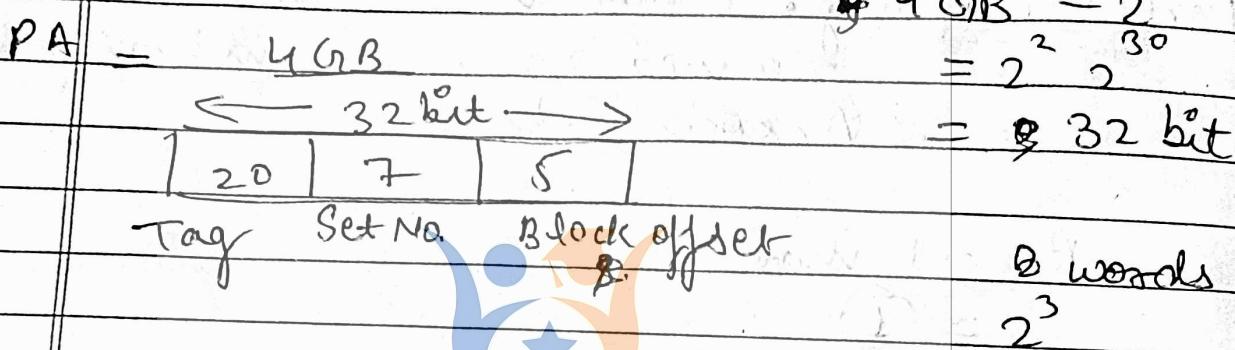


Cache memory

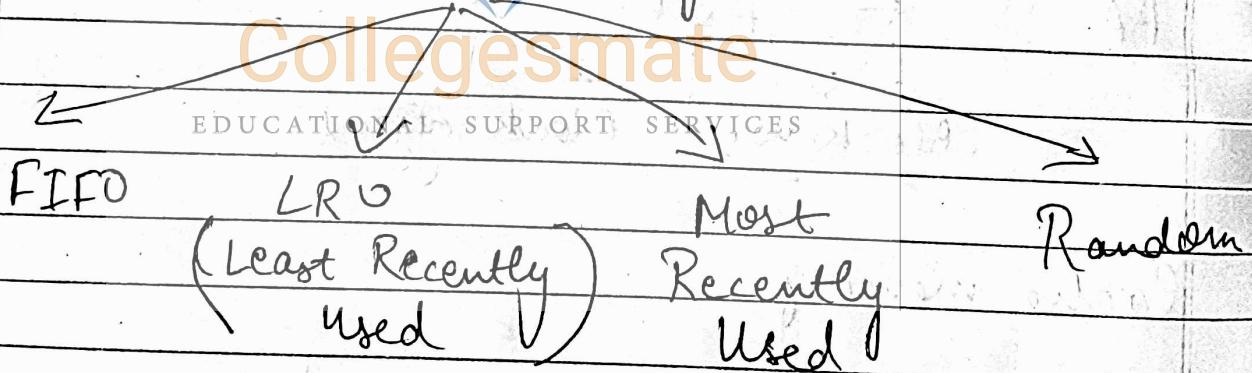
Tag	Set No.	Block offset
4	1	2

Eg:

If a 4 way set associative cache memory with a capacity of 16 kB is built with a capacity of 16 words using a Block size of 8 words. The word length is 32 bit. The size of PA space is 4 GB. The no. of bits for the tag field is



Cache Replacement Algorithm



Ques:

Consider a 4 way set associative cache with total 16 cache blocks main memory block. Request are

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

S_0	0	Set 1 16 lines	
S_1		Set 2 4 set	16/4
S_2		Set 3 no. of set	
S_3		Set 4	= 4
	1255	K mod m	
		0 mod 4	- 0

Ques: Consider a fully associative cache with 8 cache blocks (0-7) and the following sequence of memory block requests.

4, 3, 25, 0, 19, 6, 25, 8, 16, 35, 45, 22, 8, 3, 16, 25, 7

If LRU replacement Policy is used, which Cache Block will have memory block 7?



Collegesmate
EDUCATIONAL SUPPORT SERVICES

0	4	45
1	3	22
2	25	
3	8	
4	19	
5	6	
6	16	
7	35	

Cache.

Virtual Memory: A computer can address more memory than the amount physically installed on the system. This extra memory is actually called virtual memory and it is a section of hard disk that's set up to emulate the computer RAM.

Virtual Address

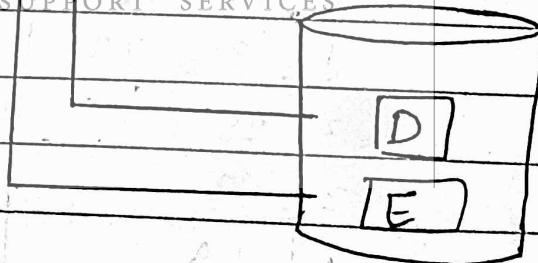
0	A
4K	B
8K	C
12K	D
16K	E

Physical Address

0	
4K	C
8K	
12K	
16K	B
20K	
24K	A

Collegesmate

EDUCATIONAL SUPPORT SERVICES



Secondary Memory