

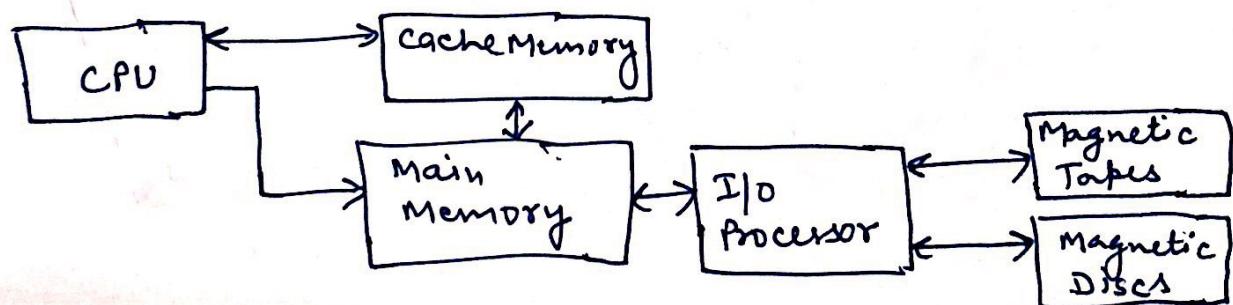
LECTURE

Content: Basic concept and hierarchy of memory

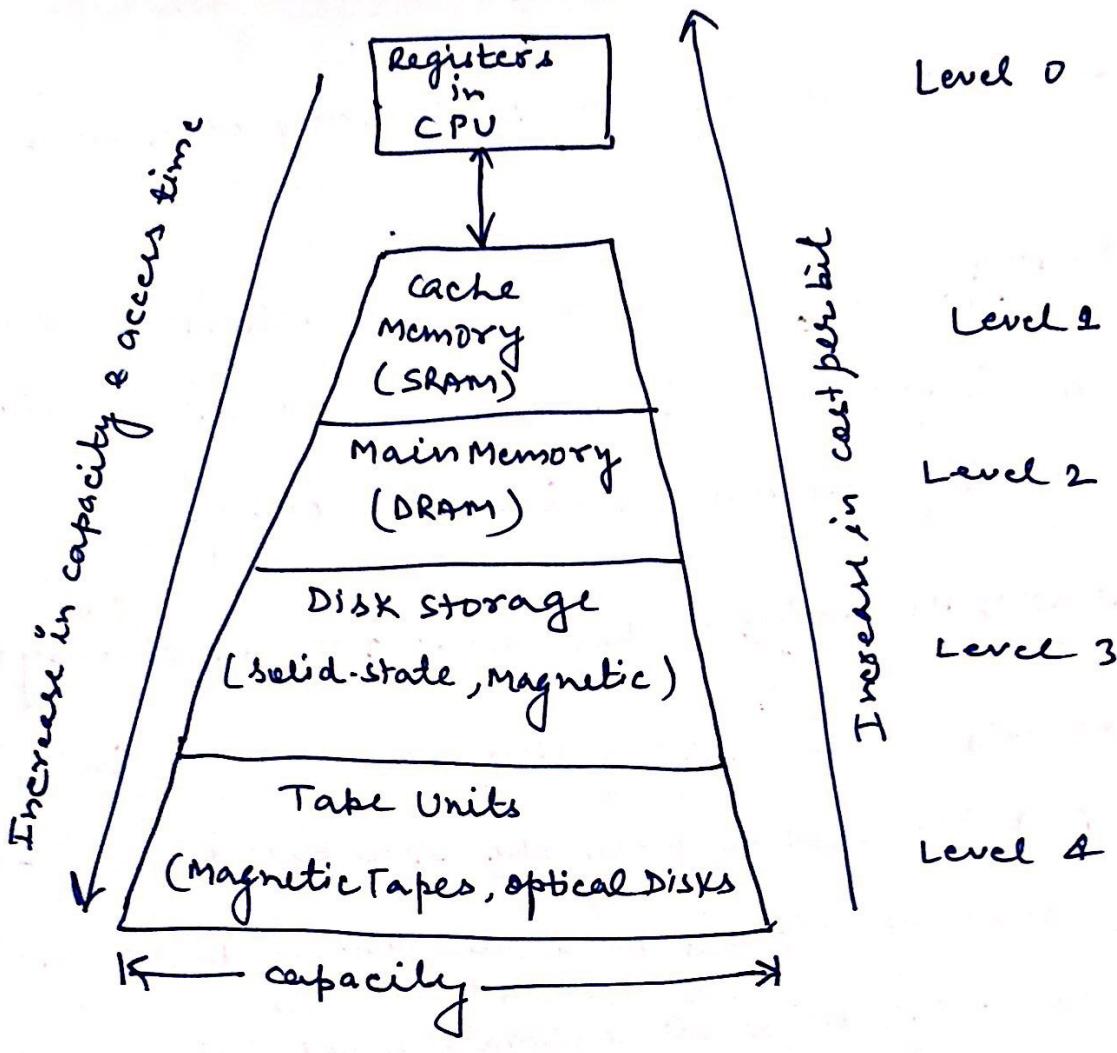
Memory: It is a storage unit which stores the data and information in a computer. It stores data and information in binary (0 or 1).

Memory Hierarchy: The memory hierarchy system consists of all storage devices employed in a computer system from the slow but high capacity auxiliary memory to a relatively faster main memory, to as even smaller and faster cache memory.

The advantage of memory hierarchy system is to increase the storage capacity, increase the processing speed and reduce the storage cost of the computer system.



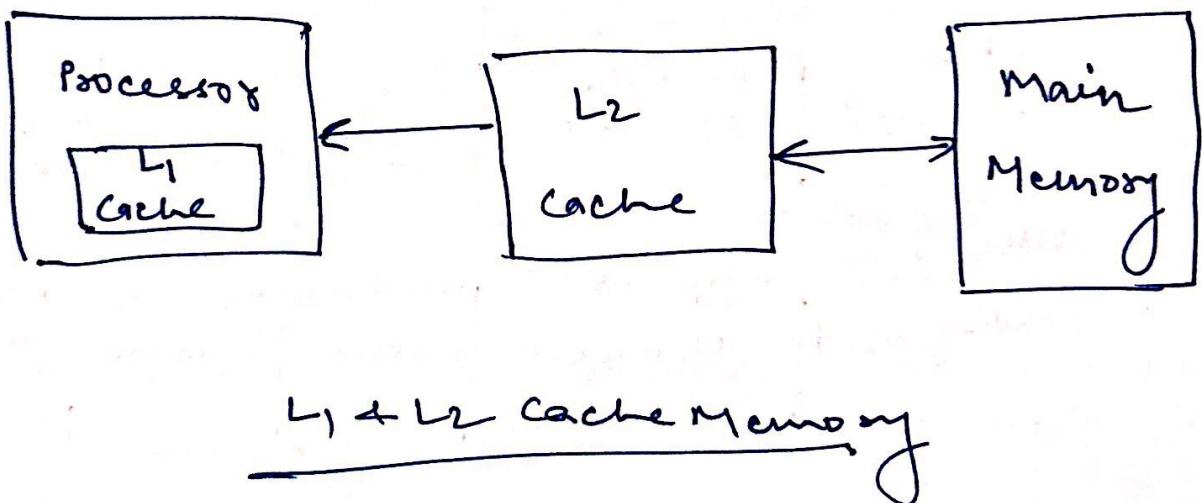
Memory hierarchy in a computer system.



A four level Memory hierarchy system

Cache Memory: Cache memory is defined as a high speed memory that is placed between the CPU and main memory to compensate the speed difference between the two cache is made up of SRAM.

It is also possible to place an even smaller cache between the cache and the processor thereby creating two level of cache shown in figure below.

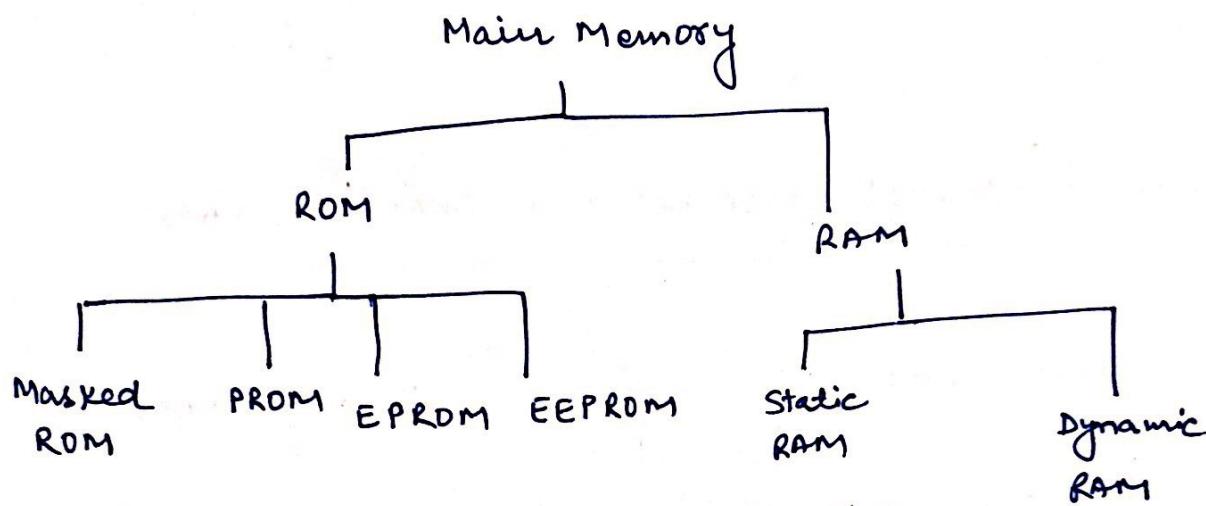


Memory level characteristics	Level 0 CPU register	Level 1 cache	Level 2 main memory	Level 3 disk storage	Level 4 Tape storage
Device Technology	ECL	256 K-bit SRAM	4 M-bit DRAM	1-6 bytes Magnetic disk unit	5-6 bytes magnetic tape unit
Access time, t_A	10 ns	25-40 ns	60-100 ns	12-20 ms	2-20 min
Capacity, s_i in bytes	512 bytes	128 Kbytes	512 Mbytes	60-220 Gbytes	512 Gbytes - 2 Tbytes
Cost, c_i (in cents/KB)	18000	72	5-6	0.23	0.01
Bandwidth, b_i (in MB/s)	400 - 800	250 - 400	80 - 133	3-5	0.18 - 0.23
Unit of transfer	4-8 bytes per word	32 bytes per block	0.5-1 Kbytes per page	5-512 Kbytes per file	Backup storage
Allocation Management	Compiler assignment	Hardware control	Operating system	Operating system/ user	Operating system/ user

LECTURE

Content: Semi-conductor RAM Memories

Main Memory Classification: The main memory is classified as:



Random Access Memory (RAM):

It is the main memory of a computer. It is used to store data and application that are currently in use. Both read and write operations are performed in RAM. In RAM data is accessed randomly (some time to access any word in the RAM at any location). RAM is volatile (retains data & information written in it till power supply is provided).

The two main classification of RAM are:

SRAM (Static RAM): It is made up of flip flops, so it holds data without external refresh as long as power is supplied to the circuit.

DRAM (Dynamic RAM): It is made up of capacitors, so external refreshing is done, many hundreds of times per second.

Comparison/Difference between SRAM & DRAM

Static RAM	Dynamic RAM
1. It contains less memory cells per unit area	1. It contains more memory cell as compared to SRAM per unit area.
2. It requires less access time hence faster memory	2. Its access time is greater than SRAM.
3. It is made up of flip flops.	3. It is made up of capacitor + MOSFET.
4. Refreshing circuit is not required	4. Refreshing circuitry is required to maintain the charge on the capacitors after every few milliseconds.
5. Cost is more	5. Cost is less.

Note: DRAM is the main component of the main memory because it has more memory cells per unit area as compared to SRAM.

Read Only Memory (ROM) :

ROM is a non-volatile memory which retains the data even when power is removed from this memory.

It is also random access. It is used for storing programs that are permanently resident in the computer.

ROM is used to store an initial program called a bootstrap loader. The function of this program is to start the computer operating system, when the power is turned on.

When power is turned on, the hardware of the computer set the program counter to the first address of the bootstrap loader. The bootstrap program loads a portion of the operating system from disk to main memory and control it then transferred to the operating system, which prepares the computer for general use -

The various kinds of ROM are -

Masked ROM:

- 1- Mask Programming : It is done by the company during the fabrication process of the unit. The procedure for fabricating a ROM requires that the customer fills out the truth table he wishes the ~~some~~ ROM to satisfy.

~~2 Proofs:~~ ROM contains all

- ① Masked ROM: The very first ROM's were hard-wired devices that contained a pre-programmed set of data or instructions. These kind of ROMs are known as masked ROM, which are inexpensive.
- ② PROM (Programmable ROM): PROM is read-only memory that can be modified only once by a user. The user buys a blank ROM and enters the desired contents using a PROM program.
- ③ EPROM (Erasable and Programmable ROM):
EPROM can be erased by exposing it to ultra-violet light for a duration of up to 40 minutes.
- ④ EEPROM (Electrically Erasable & Programmable ROM):
EEPROM is programmed and erased electrically. It can be erased and programmed about ten thousand times. Both erasing and programming take about 40 to 10 ms.

LECTURE

content: 2D & $2\frac{1}{2}$ D Memory organization.

RAM organizations:

To perform read and write operation in RAM memory can be organized by the following three methods:

1. One dimensional RAM (1-D RAM)
2. Two dimensional RAM (2-D RAM)
3. $2\frac{1}{2}$ dimensional RAM ($2\frac{1}{2}$ D RAM)

One dimensional (1-D) RAM :

This cell structure is of a one dimensional length. It has large number of addressable storage locations (2^m). each address of which stores n bit word. therefore the capacity of RAM is specified by -

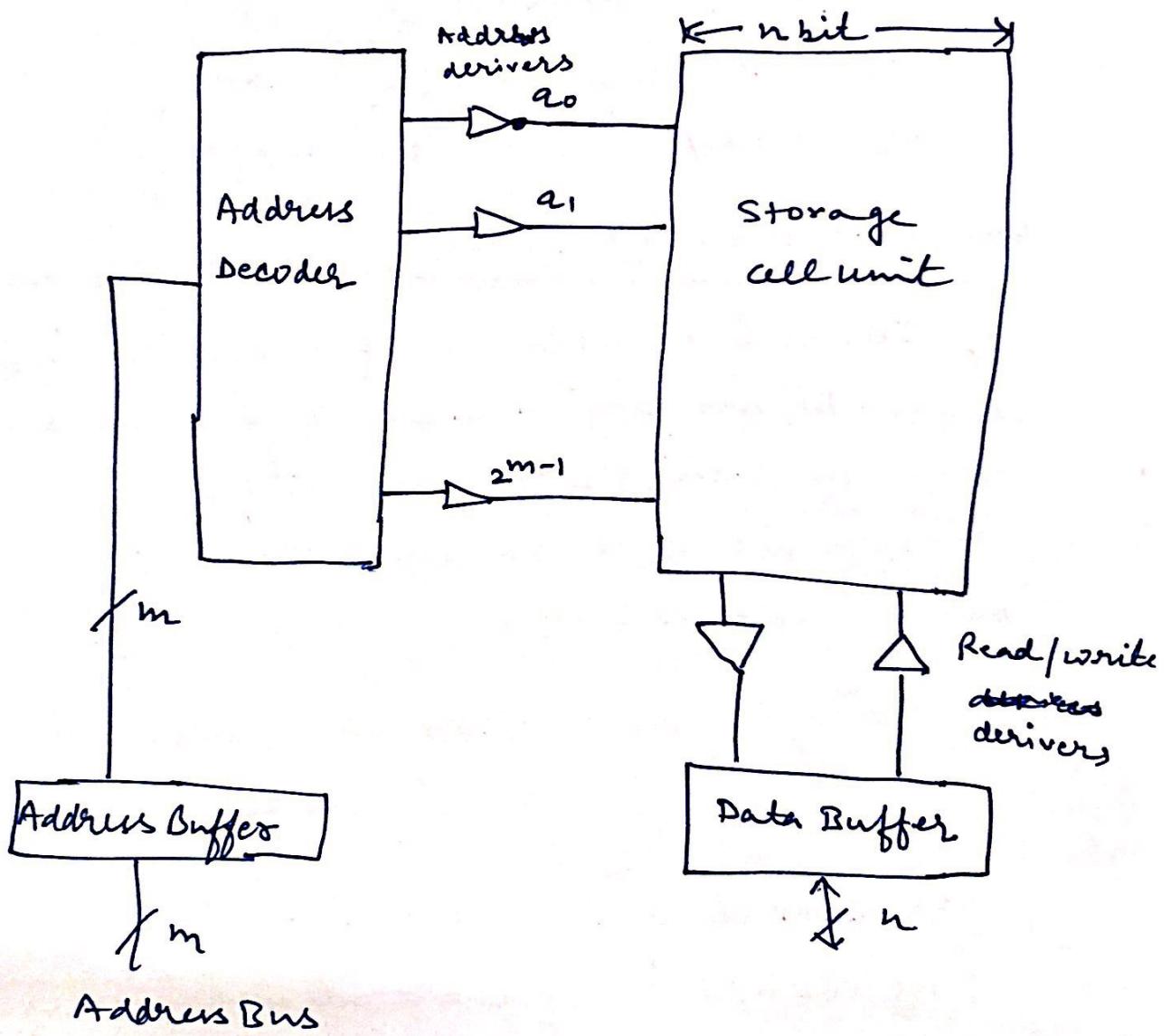
$$2^m \times n \text{ where } (2^m \rightarrow \text{no of address lines, } n \rightarrow \text{no of bits per address})$$

The RAM operates as follows -

First the address of the word to be accessed in RAM is transferred to address buffer of the RAM. The address is

the contents of the location L of the storage cell, RAM) is transferred to data buffer and from there to data bus. If the operation is write, the data from data bus is transferred to location specified in the storage cell, RAM.

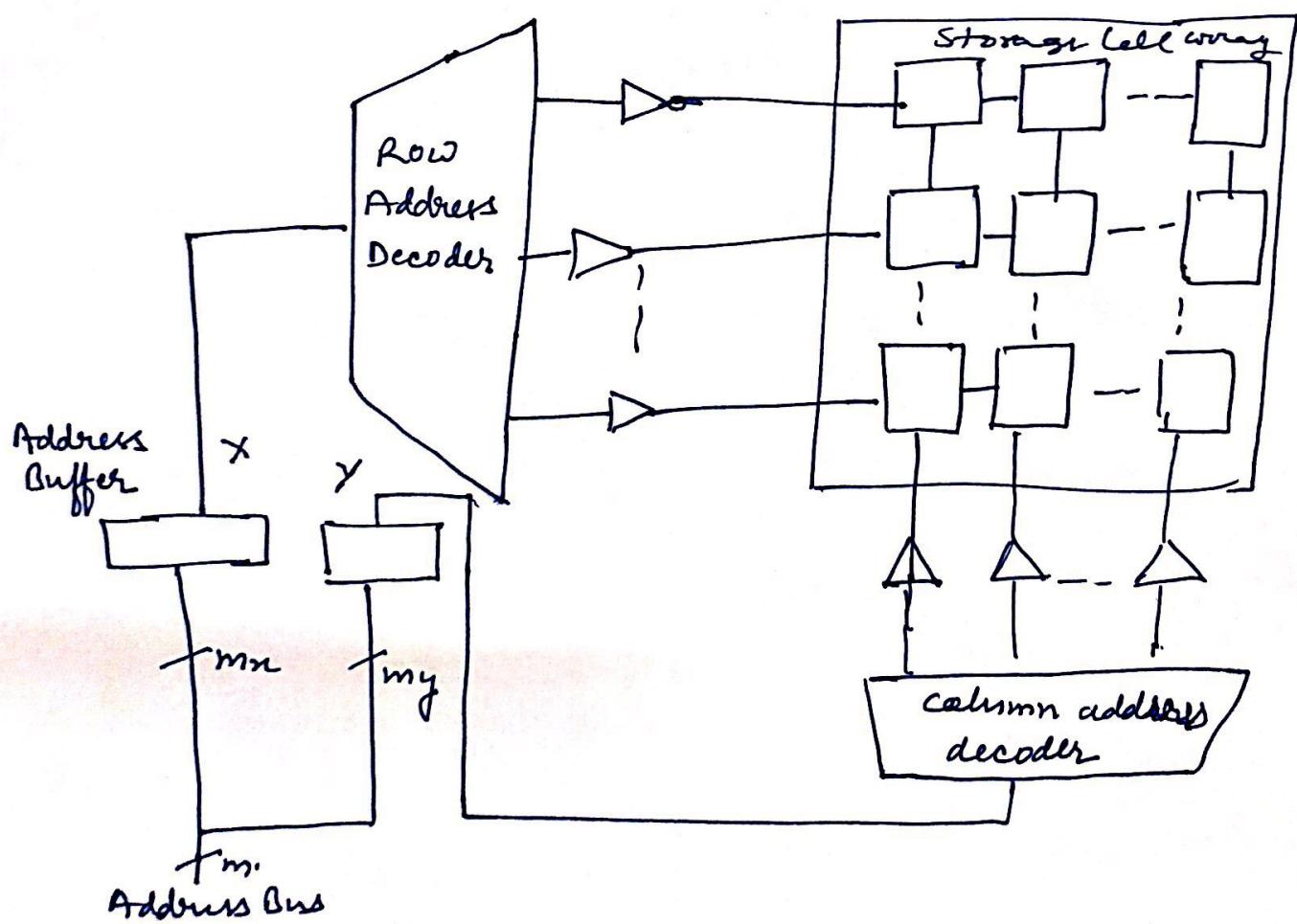
It is used for low memory capacity organization.



One Dimensional (1-D) RAM

Two Dimensional (2-D) RAM: In large memory capacity the chips are often arranged as row multiplexed column. Here the m -bit address word is divided into two parts x and y for address selection and column selection respectively. The cells are arranged in a rectangular array of N_x rows and N_y columns. (Where m_x, m_y are the number of bits for address and column line respectively as $N_x \leq 2^{m_x}$, $N_y \leq 2^{m_y}$). The total number of cell in 2D memory is $N_x \cdot N_y$.

- 2D is used for large memory organization.



Two dimension (2-D) RAM

$2\frac{1}{2}$ -Dimension ($2\frac{1}{2}$ -D) RAM :

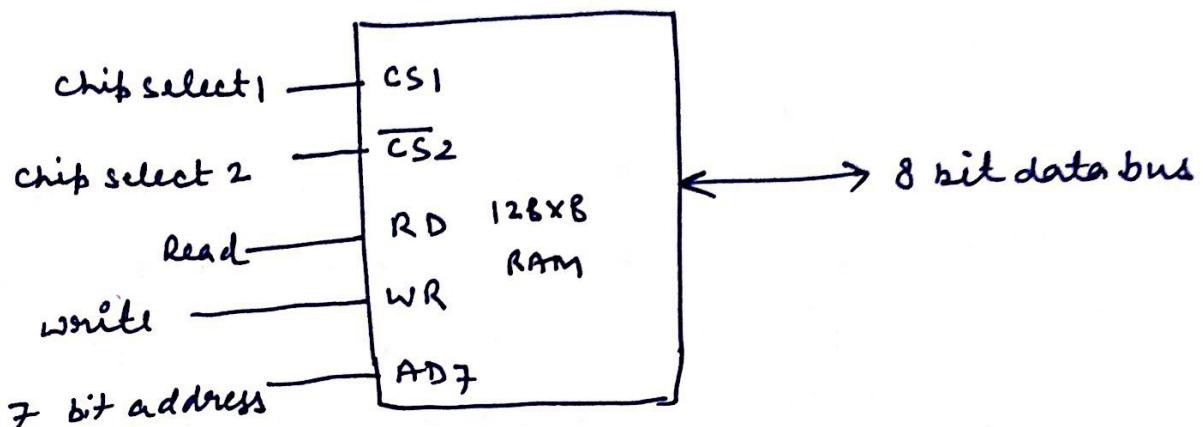
$2\frac{1}{2}$ -D is similar to 2-D RAM but in this the column lines and row lines are split into unequal size. (In 2-D the column lines and row lines are equal)

RAM and ROM chips:

RAM and ROM chips are available in a variety of size. If the memory needed for the computer is larger than the capacity of one chip. It is necessary to combine a number of chips to form the required memory size.

RAM Chip: The figure shows a RAM chip. It has one or more control inputs that select the chip only when needed.

- It operates only when $CS_1 = 1$ and $\overline{CS}_2 = 0$.
- If the chip select inputs are not enabled, or if they are enabled but the read or write inputs are not enabled, the memory is inhibited and its data bus is in a high impedance state.
- RAM chip has bidirectional data bus that allows the transfer of data either from memory to CPU (Read operation), or from CPU to memory (write operation).
- RAM chip has 7 No. of address lines which are used to select a specific word from the RAM chip of 128 bytes ($128 = 2^7$)

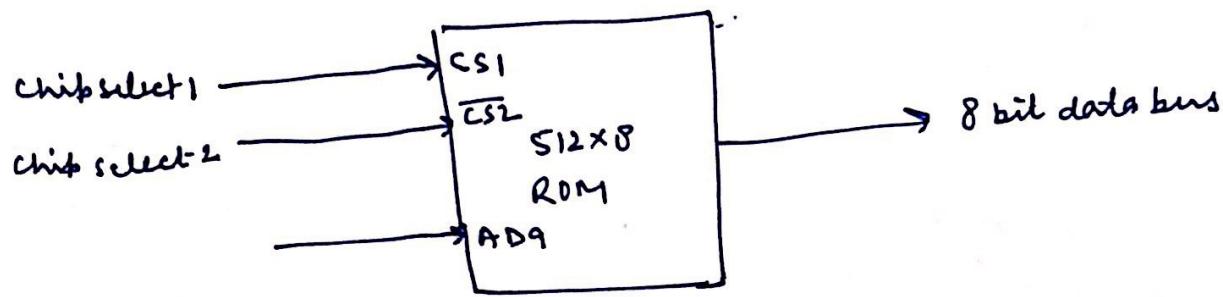


Typical RAM chip

ROM chip :

The figure shows a ROM chip. It has one or more control inputs that select the chip only when needed.

- It operates only when $CS_1 = 1$ and $\overline{CS}_2 = 0$
- It does not require read or write control because the unit can only read. Thus when the chip is enabled by the two select inputs, the data of the location selected by the address lines appears on the data bus.
- For the same size of RAM or ROM chip, it is possible to have more word in ROM, because the internal binary cells in ROM occupied less space than in RAM. For this reason, the diagram

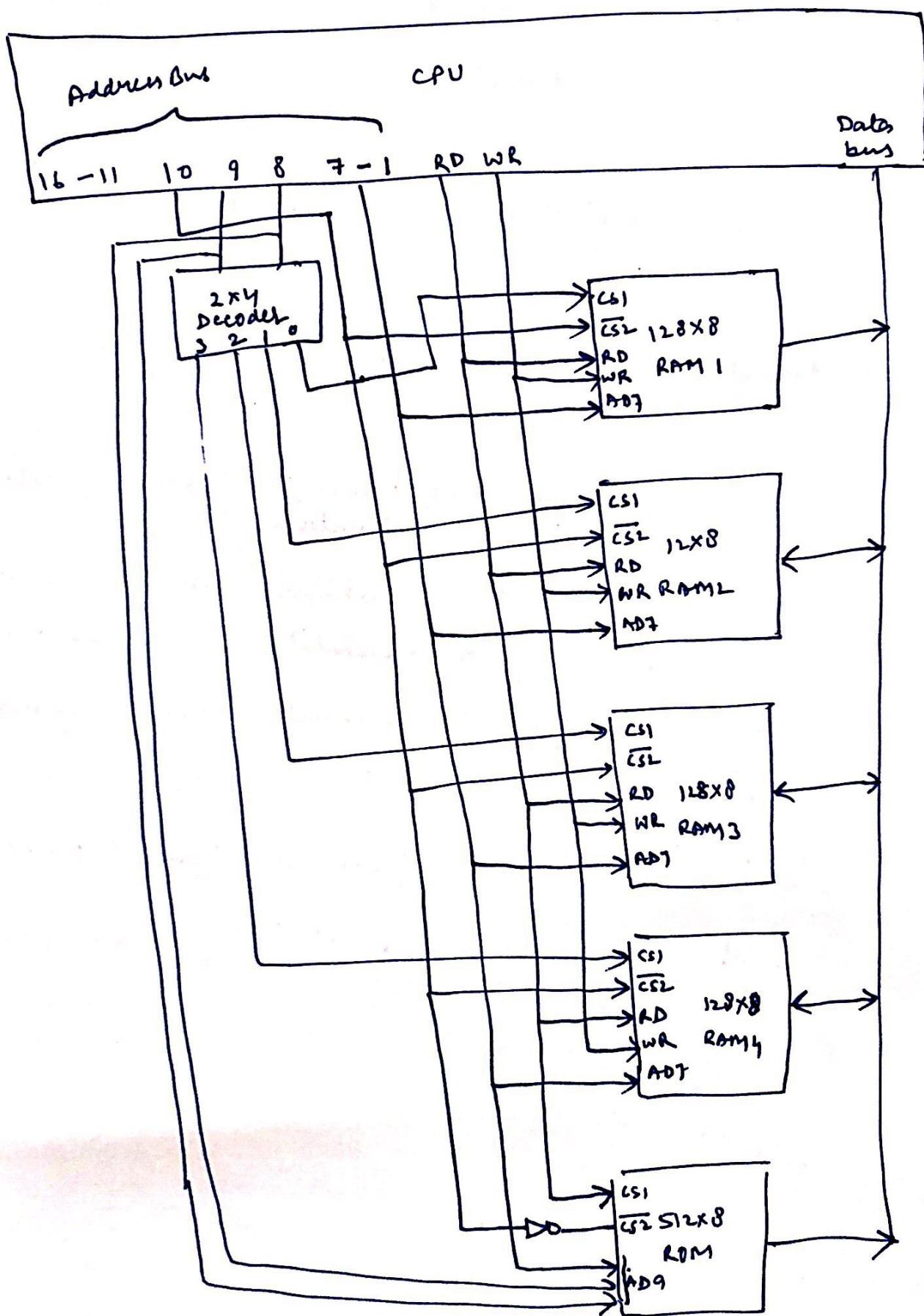


Typical ROM-chip - Block diagrams

Functional Table For RAM chip

CS1	$\overline{CS2}$	RD	WR	Memory Function	state of data bus
0	0	x	x	Inhibit	High Impedance
0	1	x	x	Inhibit	High Impedance
1	0	0	0	Inhibit	High Impedance
1	0	0	1	write	Input data to RAM
1	0	1	0	read	Output data from RAM
1	1	x	x	Inhibit	High Impedance

Memory Connection to the CPU :



Memory Address Map

Component	Hexadecimal address	Address bus								
		10	9	8	7	6	5	4	3	2
RAM1	0000 - 007F	0	0	0	x	x	x	x	x	x
RAM2	0080 - 00FF	0	0	1	x	x	x	x	x	x
RAM3	0100 - 017F	0	1	0	x	^	x	x	x	x
RAM4	0180 - 01FF	0	1	1	x	x	x	x	x	x
ROM	0200 - 03FF	1	x	x	x	x	x	x	x	x

- we need to provide a
- memory capacity of 2048 bytes?
- (b). How many lines of the address bus must be used to access 2048 bytes of memory? How many of these lines will be common to all chip.
- (c). How many lines must be decoded for chip select? specify the size of the decoders.

Ans: @ Number of RAM chips needed = $\frac{2048 \times 8}{128 \times 8}$

= 16 chips of 128×8 are needed.

① No. of address lines required (N)

$$= 2048 = 2^{11}$$

$$\boxed{N = 11}$$

No of lines common to all chips (m) = $128 = 2^7$

$$\boxed{M = 7}$$

② No. of lines for chip selection = 4 as there are 16 chips size of decoder 4×16 .

- Q. A computer uses RAM chips of 1024×1 capacity.
- How many chips are needed and how should their address lines be connected to provide a memory capacity of 1024 bytes.
 - How many chips are needed to provide a memory capacity of 16 K bytes? Explain in words how the chips are to be connected to the address bus.

Ans. a No. of chips required = $\frac{1024 \times 8}{1024 \times 1} = 8$ chips.

\therefore 8 chips of 1024×1 are required in parallel to provide a memory capacity of 1024 bytes.

b. No. of chips required = $\frac{16 \times 1024 \times 8}{1024 \times 1}$
 $= 128$ chips are required.

Total no of lines required $16K = 2^{14}$.

14 address lines are required

10 lines are used to specify chip address.

4 lines are used to specify chip number.

Q. Extend the memory system of Fig(12-4) (512x8 ROM, 128x8 RAM chips are available) to 4096 bytes of RAM and 4096 bytes of ROM. List the memory address map and indicate what size decoders are needed.

Sol: No of RAM chips needed = $\frac{4096 \times 8}{128 \times 8} = 32$ RAM chips

No of ROM chips needed = $\frac{4096 \times 8}{512 \times 8} = 8$ ROM chips

Total Memory = $4096 + 4096 = 8192 = 2^{13}$

∴ Total 13 Address lines are required

- 1 is used to select the RAM or ROM.

- 12 lines of RAM are used as

7 lines are common to all RAM.

5 lines are used to select the RAM chip
+ (5×32) decoder is used.

- 12 lines are used for ROM as

9 lines are common to all ROM.

3 lines are used to select the ROM chip
+ (3×8) decoder is used.

Memory address map:

Component	Hexadecimal	Add. in	Address lines											
			16	15	14	13	12	11	10	9	8	7	6	5
RAM	0000-0FFF	- - - 0	5x32 decoder	X X X X X X X X X X X X X X X X										
ROM	1000-1FFF	- - - 1	3x8 decoder	X X X X X X X X X X X X X X X X										

- Q. A computer employs RAM chips of 256×8 and ROM chips of 1024×8 . The computer system needs 2K bytes of RAM, 4K bytes of ROM and 4 interface units, each with 4 registers. A memory mapped I/O configuration is used. The two higher order bits of the address bus are assigned 00 for RAM, 01 for ROM and 10 for interface registers.
- (a) How many RAM and ROM chips are needed
 (b) Draw a memory-address map for the system
 (c) Give the address range in hexadecimal for RAM, ROM and interface.

④ No. of RAM chips needed = $\frac{2 \times 1024 \times 8}{256 \times 8} = 8$ RAM chips
 of 256×8 .
 are required.

No. of ROM chips needed = $\frac{4 \times 1024 \times 8}{1024 \times 8} = 4$ ROM chips of
 1024×8 are required.

⑤ No. of address lines required for

RAM = $2^{10} = 2^8$
 → 8 lines are common to all chips.
 → 3 lines are used for chip selection
 + (3×8) decoder is used

ROM = $4096 = 2^{12}$
 → 10 lines are common to all chips
 → 2 lines are used for chip selection
 + (2×4) decoder is used

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

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

RAM 0000-0FFF 0 0 0 0 0 $\xrightarrow{3 \times 8 \text{ decoder}}$ X X X X X X X X X X X X

ROM 4000-4FFF 0 1 0 0 $\xrightarrow{2 \times 4 \text{ decoder}}$ X X X X X X X X X X X X

Interface 8000-800F 1 0 0 0 0 0 0 0 0 0 X X X X

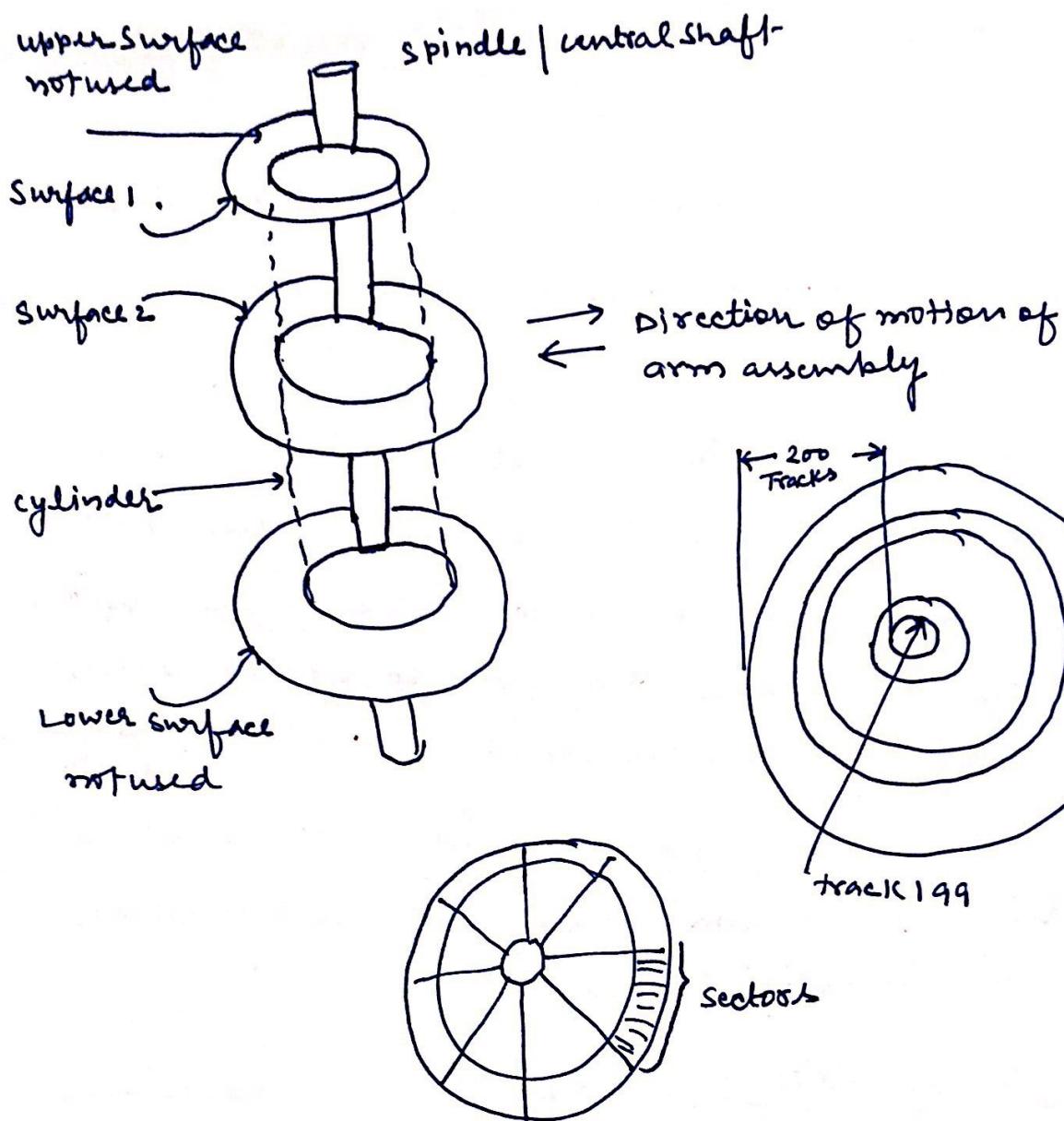
LECTURE

Magnetic Hard Disk : Magnetic disks are thin, circular metal plates coated on back the sides with a thin film of magnetic material.

- A set of such magnetic disks/ plates are fixed to a spindle one below the other to make up a disk ~~space~~ pack.
- The disk pack is sealed and mounted on a disk drive.
- The disk drive consists of a motor to rotate the disk pack about its axis at a speed of about 3600 rev/min.
- The drive also has a set of magnetic heads mounted on ~~axis~~ arms. The arm assembly is capable of moving in & out in radial direction.

Storage of Information : Information is recorded on both the surfaces of each disk plate, as it rotates about axis. However, the upper surface of the top plate & lower surface of the bottom plate are not used.

- Each disk consists of
- A set of corresponding tracks in all surfaces of a disk pack is called a cylinder.



Surfaces \therefore it has 18 tracks/cylinder

Each track is divided into sectors.

The total no. of bytes that can be stored in a disk pack is = no. of cylinders \times tracks/cylinder \times sectors/track
 \times Bytes per sector.

Eg: Disk pack has 10 disk plates each having 200 tracks.



18 recording surfaces \Rightarrow 18 tracks/cylinder

Suppose, there are 40 sectors/track & each sector can store 256 bytes.

200 tracks/plate \Rightarrow 200 cylinders

$$\therefore \text{Capacity of disk pack} = 200 \times 18 \times 40 \times 256$$

$$= 31864000 \text{ bytes}$$

$$= 36.864 \text{ Mbytes}$$

- A set of disk drives are connected to a disk controller. The disk controller accepts commands from the computer & positions the read-write heads of the specified disk for reading or writing.

- In order to read from or write on a disk, the computer must specify the drive no., cylinder no., surface no. & the sector no.. The drive no. should be specified because a controller normally controls more than one drive.

When a read-write command is received by the disk controller, the controller first positions the arm assembly so that the read-write head reaches the specified cylinder. The time taken to reach the specified cylinder is known as the seek-time (T_s). The seek time, the rate at which information is read from the disk

$$\therefore \text{transfer rate} = 10260 \times \frac{3600}{60} \text{ bytes/sec.}$$
$$= 615600 \text{ bytes/sec.}$$

$$\text{Time to read one sector is } \frac{256}{615600} = 4.138 \text{ msec.}$$

\therefore Avg time to access a sector is $20 + 8.3$ sec.

Compact Disk - Read only Memory (CD-ROM)

Also known as a Laser disk. It is a shiny metal like disk whose diameter is 5.25". It can store around 650 megabytes (equivalent to 2,50,000 pages) info. In CD-ROM is written by creating pits on the disk surface by shining a laser beam.

As the disk rotates the laser beam traces out a continuous spiral. The sharply focused beam creates a circular pit of around 5×10^{-4} mm whenever a 1 is to be written & no pit (also called a land) if a zero is to be written.

LECTURE

Content: Cache performance Issue..

Cache performance Issue :

Cache Hit : When CPU needs to read data/instruction from memory, first cache is examined, If the word is found in cache, it is called cache hit.

Cache Hit Time : In case of a cache hit, the required word is read from the cache, This time to access data from cache is cache hit time.

Cache miss : If the requested data is not found in cache memory, then a cache miss is said to occur.

Cache miss Time Penalty : In case of cache miss, the time required to fetch the required block from main memory (to deliver to CPU) is called cache miss time penalty

Hit Ratio (h):

$$= \frac{\text{No. of Hits}}{\text{No. of Cache Hits} + \text{No. of Misses}}$$

Miss Ratio ($1 - h$):

$$= \frac{\text{No. of cache Misses}}{\text{No. of cache Hits} + \text{No. of cache Misses}}$$

Average Memory Access Time:

$$t_A = h t_{c'} + (1-h) L t_c + t_m$$

↓ ↓ ↓ ↓

Average Memory Access Time Miss Ratio Cache memory access time
+ Main memory access time

↓

Hit Ratio

↓

Cache Memory hit Access Time

Suppose, the ratio of main memory access time to cache access time is γ , then efficiency (η) of a system that employs a cache is

$$\eta = \frac{t_c}{t_A}$$

$$\eta = \frac{t_c}{ht_c + (1-h)(t_c + t_m)}$$

$$= \frac{t_c}{t_c(h + (1-h)\left(1 + \frac{t_m}{t_c}\right))}$$

$$= \frac{t_c}{t_c(h + (1-h)(1+\gamma))}$$

$$= \frac{1}{h + (1-h)(1+\gamma)}$$

$$= \frac{1}{\cancel{h} + 1 - \cancel{h} + \gamma(1-h)}$$

$$\eta = \frac{1}{1 + \gamma(1-h)}$$

η is max when $h=1$. i.e all reference are from a cache.

Q. Calculate t_A (average access time), γ (ratio of main memory access time to cache memory access time) and λ (efficiency) of a memory system whose parameters are as indicated:

$$t_c = 160 \text{ ns}$$

$$t_m = 960 \text{ ns}$$

$$\lambda = 0.90$$

Sol:

$$\begin{aligned} t_A &= \lambda t_c + (1-\lambda) (t_c + t_m) \\ &= 0.90 \times 160 + (1-0.90) \times (160 + 960) \\ &= 144 + (0.1) \times 1120 \\ &= 144 + 112 \end{aligned}$$

$$t_A = 256 \text{ ns}$$

$$\gamma = \frac{t_m}{t_c} = \frac{960}{160} = 6, \text{ so } \boxed{\gamma = 6}$$

$$\lambda = \frac{1}{1 + \gamma(1 - \lambda)}$$

$$= \frac{1}{1 + 6(1 - 0.90)}$$

$$= \frac{1}{1.6} = 0.625$$

$$\boxed{\lambda = 0.625}$$

The access time of a cache memory is 50 ns and that of the main memory is 500 ns. It is estimated that 80% of the main memory requests are for read operation and the remaining are for write. The hit ratio for read access only is 0.9 and a write through policy is used.

- (a) what is the average access time of the system considering only memory read cycles?
- (b) what is the average access time of the system for both read and write requests?
- (c) what is the hit ratio taking into consideration the write cycle?

Let:

$$\begin{aligned}
 @ \text{average access time} &= t_A = h \cdot t_C + (1 - h) \cdot (t_C + t_M) \\
 &= 0.9 \times 50 + (1 - 0.9) \times (50 + 500) \\
 &= 45 + 55 \\
 t_A &= 100 \text{ ns}
 \end{aligned}$$

$P_R \rightarrow \%$ of read operation

$$\begin{aligned}
 @ \text{Average access time} &= P_R \cdot \text{avg read time} + (1 - P_R) \cdot t_M \\
 &= 0.8 \times 50 + (1 - 0.8) \times 500
 \end{aligned}$$

$hr \rightarrow$ hit ratio for read cycle
 $hw \rightarrow$ hit ratio for write cycle

③ Hit Ratio = $Pr \times hr + (1 - Pr) \times hw$

$$= 0.8 \times 0.9 + (1 - 0.8) \times 0$$
$$= 0.72$$

Virtual Memory : In a memory hierarchy system, programs and data are first stored in auxiliary memory. Portion of a program or data are brought into main memory as they are needed by the CPU.

Virtual memory is a concept used in some large computer systems that permit the user to construct programs as though a large memory space were available, equal to the totality of auxiliary memory. Each address that is referenced by the CPU goes through an address mapping from the so-called virtual address to a physical address in main memory. Virtual memory is used to give programmers the illusion that they have a very large memory at their disposal, even though the computer actually has a relatively small memory.

A virtual memory system provides a mechanism for translating program-generated addresses into correct main memory locations.

This is done dynamically, while programs are being executed in the CPU - the translation or mapping is handled automatically by the hardware by means of mapping Table.

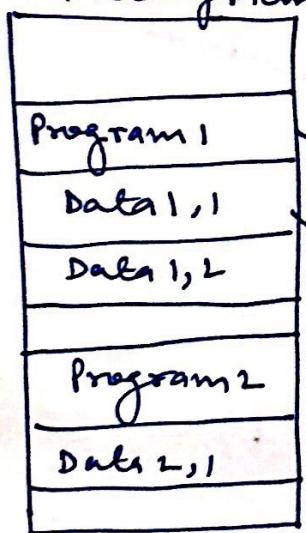
Address Space and Memory Space:

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

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

Thus the address space is the set of addresses generated by programs as they reference instructions and data; the memory space consists of the actual main memory locations directly addressable for processing. In most computers, the address and memory spaces are identical. The memory space is allowed to be larger than the address space in computers with virtual memory.

Auxiliary Memory



Address space
 $N = 1024K = 2^{20}$

Main Memory



Memory Space = M

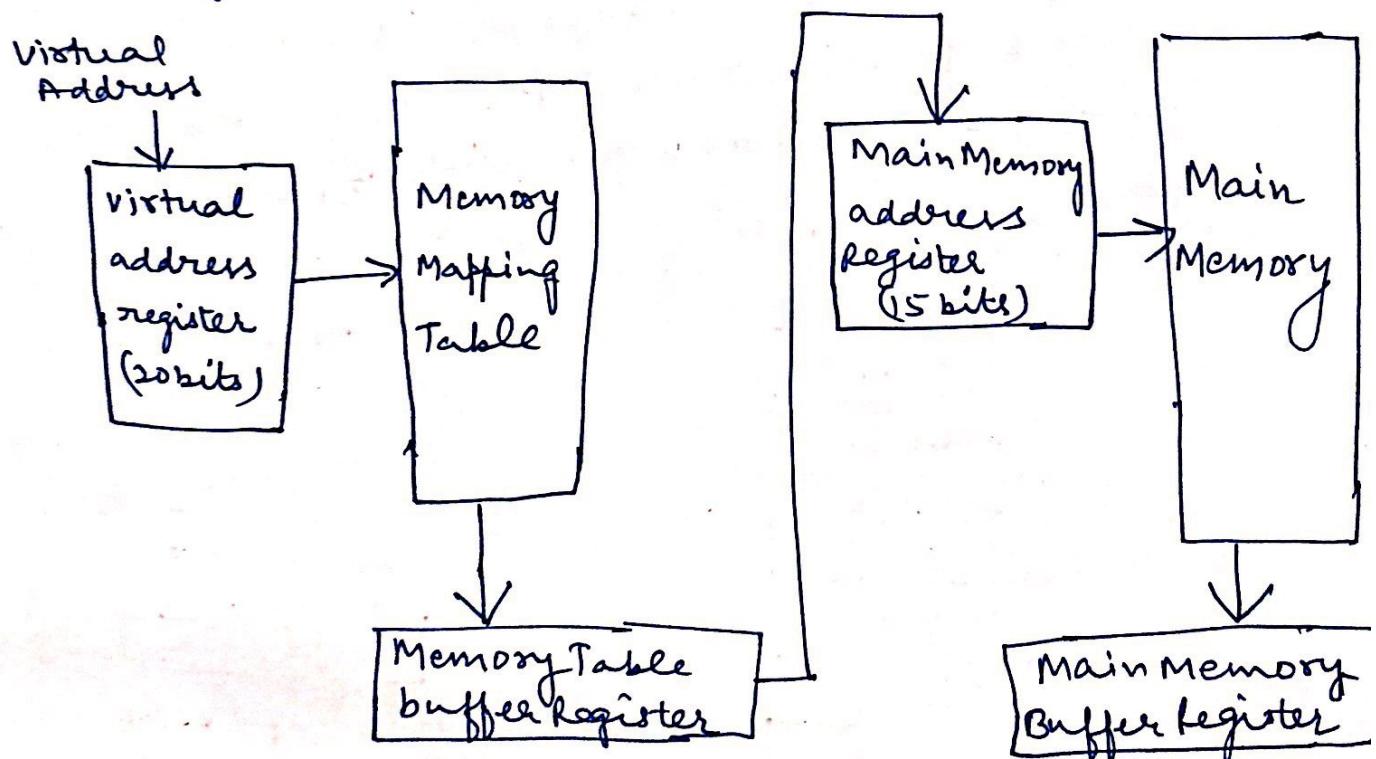
$$M = 32K = 2^{15}$$

Relation b/w address and Memory space in a virtual Memory system

Mapping Table: A table is needed, as shown in figure, to map a virtual address of 20 bits to a physical address of 15 bits. The mapping is ~~dynamic~~^{by dynamic operation}, which means that every address is translated immediately as a word is referenced by CPU.

The mapping table may be stored in a separate memory. In the first case, an additional memory unit is required as well as one extra memory access time.

In second case, the table takes space from main memory and two accesses to memory are required with the program running at ~~half speed~~.

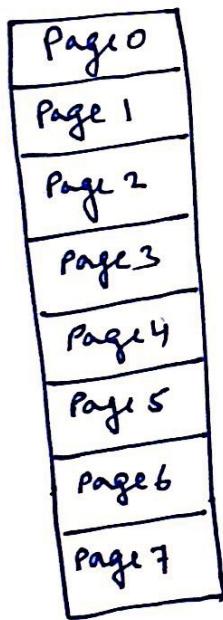


Memory table for mapping a virtual Address

Address Mapping using Pages

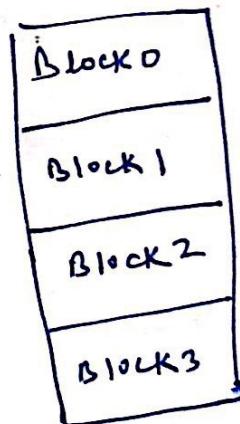
The table implementation of the address mapping is simplified if the information in the address space and the memory space are each divided into groups of fixed size. The physical memory is broken down into groups of equal size called blocks.

Consider a computer with an address space of 8K and a memory space of 4K. If we split each into group of 1K words we obtain eight pages and four blocks.



Address space

$$N = 8K = 2^{13}$$

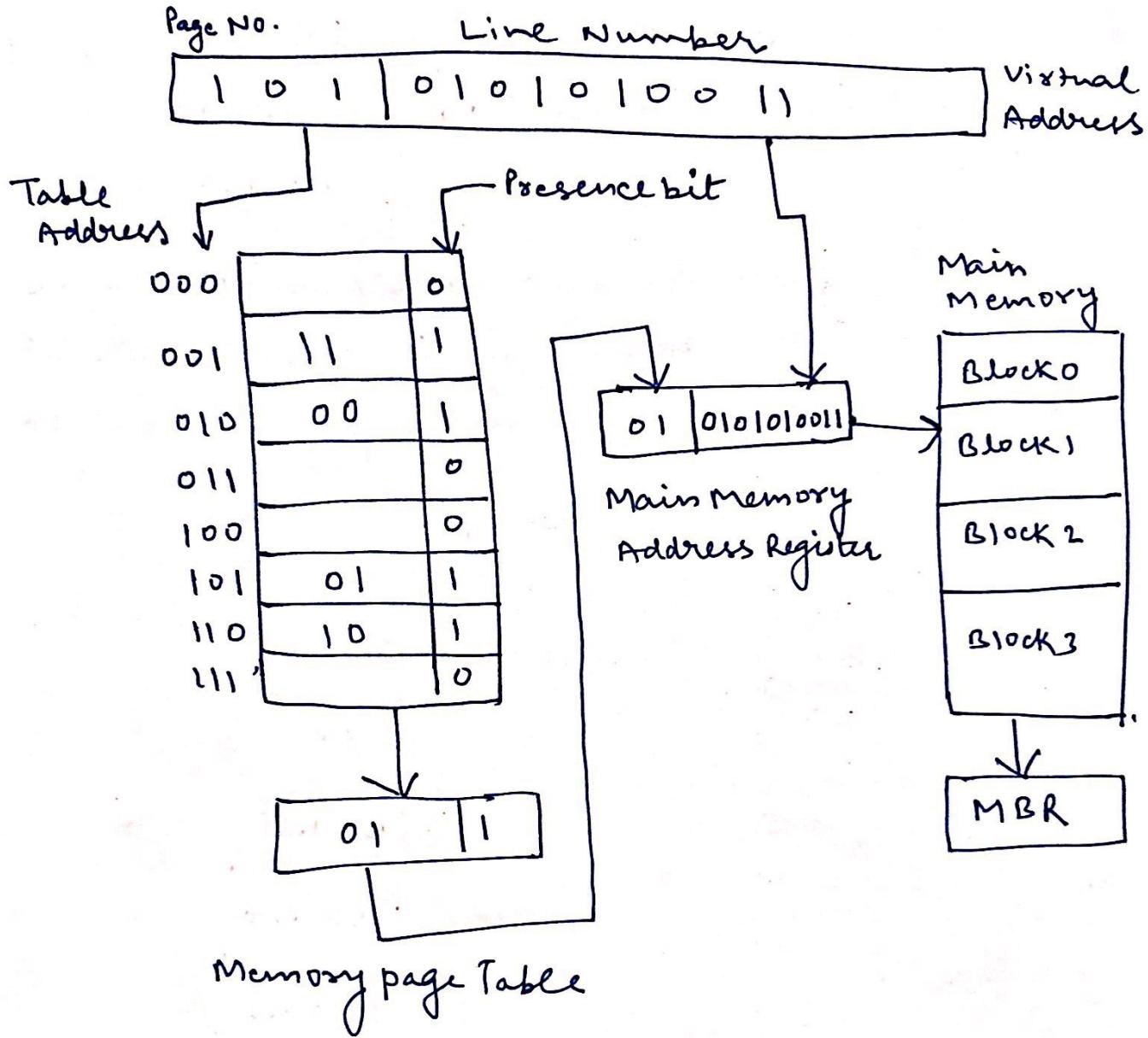


Memory space

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

Address space and memory space split into groups of 1K words

The organization of memory mapping table in a paged system is shown in figure



Memory table in a paged system

Page Replacement :

The virtual memory system is a combination of hardware and software techniques. The memory management software system handles all the software operations for the efficient utilization of memory space.

It must decide :

- 1- Which page in main memory ought to be removed to make room for a new page.
- 2- When a new page is to be transferred from auxiliary memory to main memory.
- 3- Where the page is to be placed in main memory.

The hardware mapping mechanism and the memory management software together constitute the architecture of a virtual memory.

Page fault:

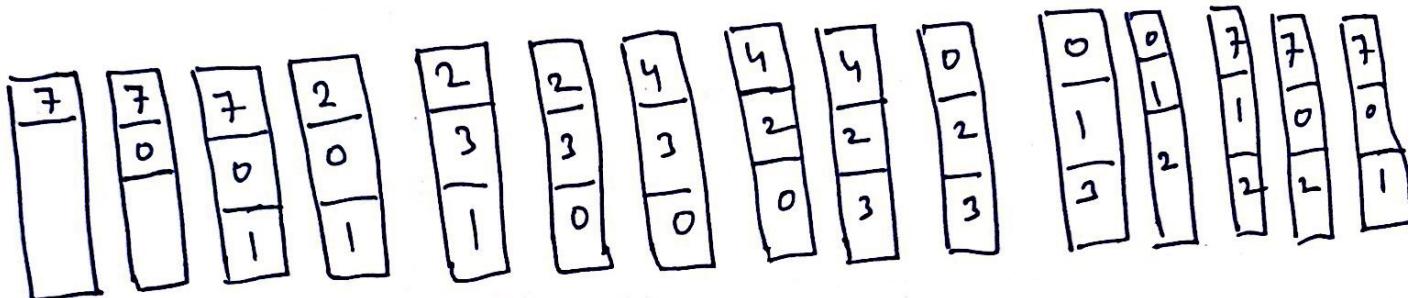
When a program starts execution, one or more pages are transferred into main memory and the page table is set to indicate their position. The program is executed from main memory until it attempts to reference a page that is still in auxiliary memory. This condition is called page fault.

When page fault occurs, the execution of the present program is suspended until the required page is brought into main memory.

There are two replacement algorithms used:

FIFO (First-In First-Out):

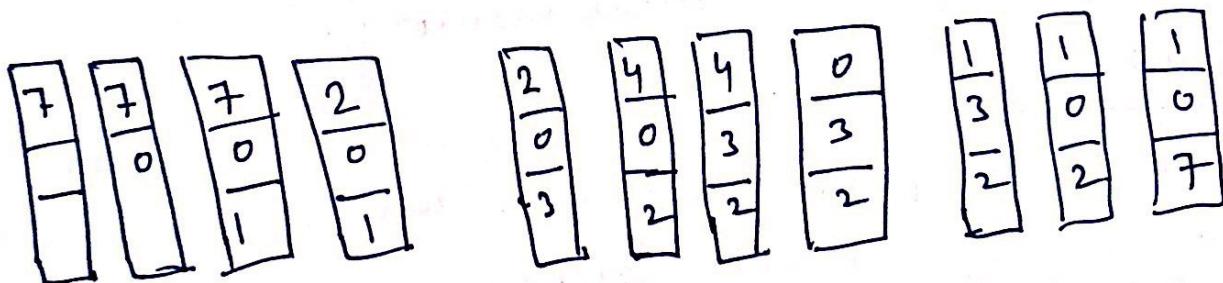
7 0 1 2 0 3 0 4 2 3 0 3 2 1 2 0 1 7 0 1



FIFO Page Replacement Algorithm

LRU (Least-Recently-used):

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



LRU page Replacement Algorithm