

Welcome to ACE Engineering Academy - online live class

Subject: **Computer Organization and Architecture**

Faculty: **Y.V. Ramaiah**

9866339106

Subject

Computer organization & Architecture

Chapters (Topics)

- I. Computer Arithmetic ✓
- II. Memory Organization
- III. Secondary Memories
- IV. Basic processor organization and Design
- V. Pipeline organization
- VI. Control unit Design
- VII. IO Organization

Chapter 2 Memory Organization

- Introduction ✓
- Memory Basics ✓
- Memory Classification ✓
- Memory Size Expansion ✓
- Primary Memory
- Secondary Memory ✓
- ROM and its design ✓
- RAM and its design
- Memory Hierarchy
- Cache Memory
- Mapping Techniques
- Different misses occurred in cache
- Different block replacement techniques
- Tag directory design
- Associative Memory

Q. A processor can support a maximum memory of 4GB,

where the memory is word-addressable (a word consists of two bytes). The size of the address bus of the processor is at least 31 bits.

2016

Max. size of mem. to be supported = $4G \times 8$ bits.
 $= \frac{32}{2} \times 8$ bits.
word size = 2 bytes.

But memory is word Addressable.

No. of words in memory = $\frac{4G}{2B}$ words.
 $= \frac{2G}{2B} = 2^{31}$ words.

Q. The main disadvantage of DRAM over SRAM is ____.

(a) High package density *advantage*

(b) Costly *X*

(c) External memory refresh logic is required *✓*

(d) High power consumption *X*

Q. A function table is required in very large numbers. The memory most suitable for this purpose would be

(a) ROM *✓*

(b) RAM

(c) EPROM

(d) EEPROM

Q. When the power supply of a ROM is switched off, its contents

- (a) Become all zero's
- (b) Become all one's
- (c) Remain same
- (d) Are unpredictable

Byte Addressable Memory:-



In this, each memory location is used to store only one byte
Let word size = 32 bit (4 Bytes)
then Byte Addressable Memory requires 4 memory locations to store one word



Word Addressable Memory:-

Each memory location is used to store one word (whatever the word size)

Let word size = 4 Bytes, only one memory location (Register) is used to store 4 Bytes.

→ Electronic memories are designed with electronic components like diode, BJT/FET and IC, to operate electronic memory, electro-mechanical devices (like Head and motor) are not required.

Electronic memory is faster than magnetic and optical memory.

CPU directly communicates with Electronic memory, hence all Electronic Memories are known as Primary Memory.

Generally Electronic Memory Speed
almost equal to the processor speed.



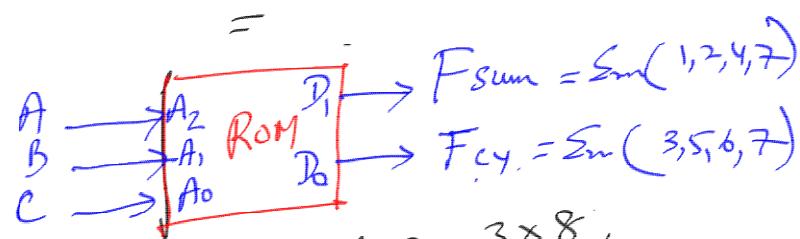
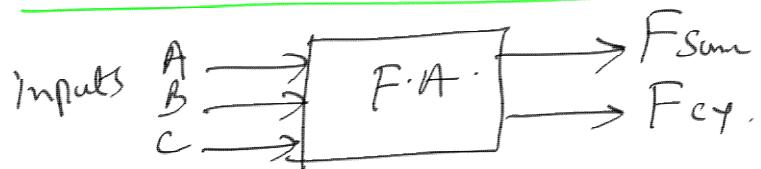
Secondary Memories:- All Magnetic
and optical memories are operated with
Head and motor, these devices are
operated in m-s. time only.
DMA Controller / IO processor is used
to transfer the Data b/w Secondary
memories and main mem.

Electronic Rom design:- Rom is the
the combination of Decoder and OR gates.
One OR gate is used to deliver only one data
bit. i.e. No. of OR gates required = no. of Data bits.
To design $2^n \times m$ bit size Rom;
size of Decoder needed $(n * 2^n)$ size Decoder
and no. of OR gates required = m .

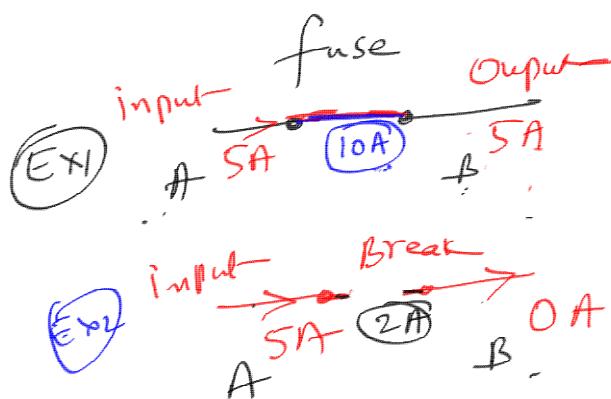
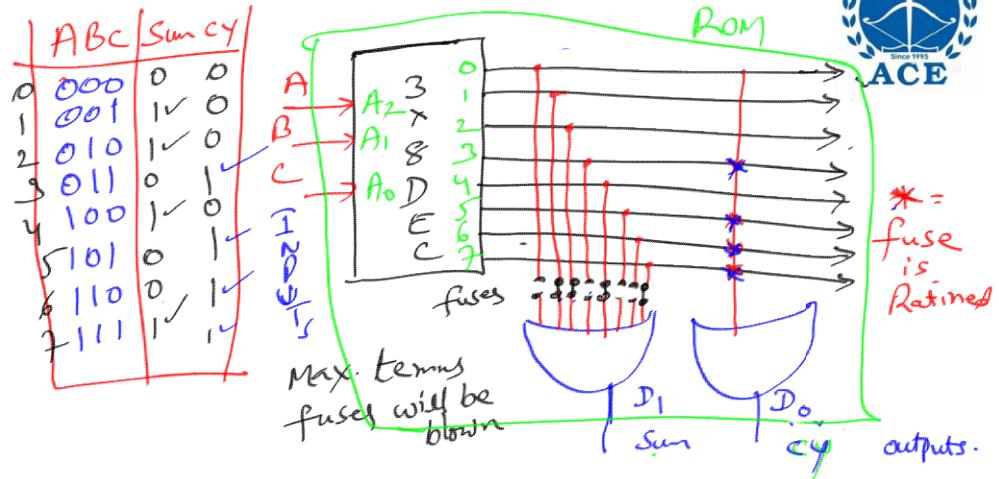


- Generally ROMS are used to implement function tables.
- In this process, functional inputs are applied to the ROM address, and functional outputs are taken from Data lines.

Implement Full Adder with ROM.



Size of the Decoder = 3×8 ,
No. of OR gates = 2



ROM is programmed according to the Table.

Programming is the process of blowing the fuses by sending higher current.



Special Applications of ROM

- ROMs are used for storing the result of arithmetical operations.
- In this process, input operands are applied to Address and their result is taken from Data.
Ex: Calculator



Inputs Result -

Ex: $4 \times 5 = 20$

$100101 : 10100_2$

Address Stored
at Data

Calculator brings pre-defined values.



Design a ROM for 2 bit multiplication

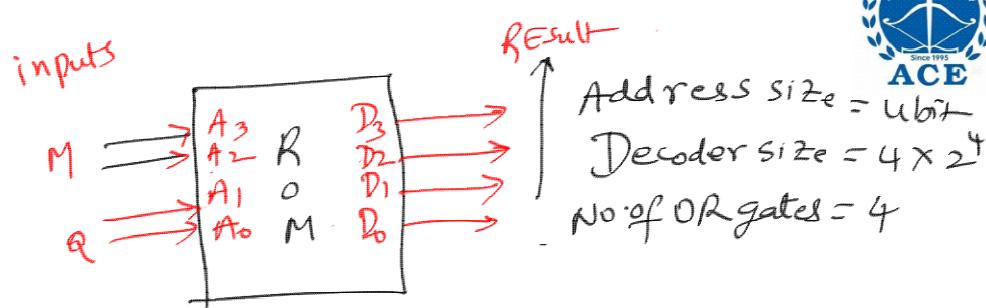
Each input operand size = 2 bits.

$$M = \text{multiplicand} \quad M * Q = \text{Product}$$

(2) (2) (4) bit

Q = multiplier

Address
After multiplying 2 no. of n bit data;
the longest result size = $2n$ bits.





M	Q	A ₃ A ₂	A ₁ A ₀	F	D ₃ D ₂ D ₁ D ₀	Result
0	0	00	00	0	0000	00000
1	0	00	01	0	0000	00000
2	0	00	10	0	0000	00000
3	0	00	11	0	0000	00000
4	0	01	00	0	0000	00000
5	0	01	01	0	0001	00001
6	0	01	10	0	0000	00000
7	0	01	11	0	0011	00011
8	1	10	00	0	0000	00000
9	1	10	01	0	0000	00000
10	1	10	10	0	0100	00100
11	1	10	11	0	0100	00100
12	1	11	00	0	0000	00000
13	1	11	01	0	0000	00000
14	1	11	10	0	0001	00001
15	1	11	11	1	0001	00011

$$F_{D_3}(A_3A_2A_1A_0) = \Sigma m(15)$$

$$F_{D_2}(A_3A_2A_1A_0) = \Sigma m(10, 11, 14)$$

$$F_{D_1}(A_3A_2A_1A_0) = \Sigma m(6, 7, 9, 11, 13, 14)$$

$$F_{D_0}(A_3A_2A_1A_0) = \Sigma m(5, 7, 13, 15)$$

Size of Decoder
needed = 4×2^4



ROM design for n bit multiplication

(M) multiplicand = n bit one Operand size = n bit
 (A) multiplier = n bit Address Size = 2n bit
 Max. no. of Addressable Regs = 2^{2n}
 Longest Result size = $2n$ bits
 No. of OR gates required = No. of Data lines
 $= 2^n$
 Size of Decoder = $2^n \times 2^{2n}$



ROM size Required for
'n' bit multiplication:-

Address
size = 2^n bit

$\frac{2^n}{2} * 2^n$ bits.

Longest Result
size = 2^n bit
(Data)

NAT
Q) The size of ROM required
for 8 bit multiplication is 128 kilobytes.

$$\begin{array}{l} \text{16} \\ 2 * 16 \text{ bit} \\ 64K * 2B = 128 \text{ KBytes} \end{array}$$

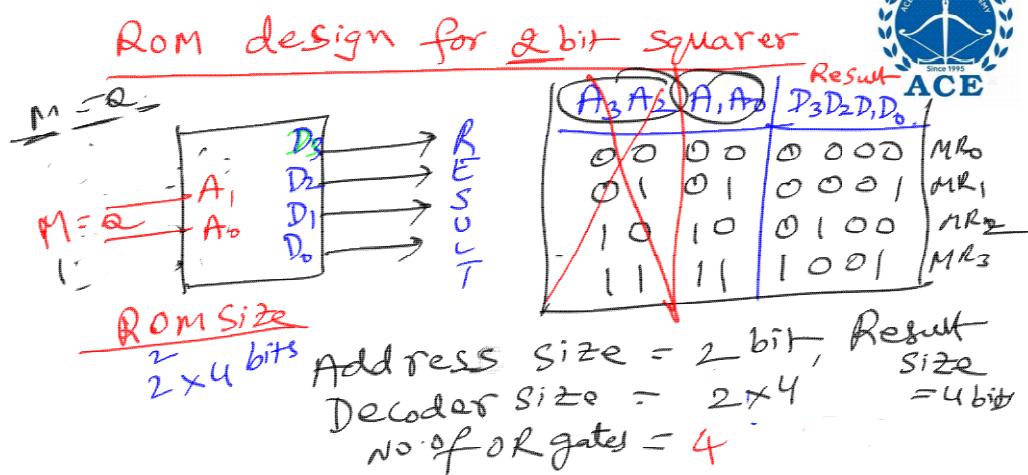
$1K = 2^{10}$



ROM design for squarer operation

Both multiplicand and multiplier are same

Let multiplicand size is 'n' bits, the longest result size = $2n$ bits





For ' n ' bit Square operation
Rom size Required

Operand size = n bit = Address size.
Longest Result size = 2^n bits = Data Bus size

$$\begin{array}{l} \text{No. of OR gates} = 2^n \\ \text{size of Decoder} \\ n \times 2^n \end{array}$$

$$\begin{array}{l} \text{Rom size} \\ 2^n * 2^n \text{ bits} \end{array}$$

NAT

Q) The size of the Rom required for 6 bit Square is 768 bits

$$\begin{aligned} & 2^6 * 12 \text{ bits} \\ & = 64 * 12 \text{ bits} \\ & = 768 \end{aligned}$$



2013 GATE

a) Size of ROM required for 4 bit
Square is — Bytes.

$$2^4 \times 8 \text{ bit} = 16 \text{ Bytes}$$

Design of Electronic RAM

- All Electronic RAMS are Volatile Memories
 - Electronic RAM can be designed in 2 ways :
 - (i) Dynamic RAM
 - (ii) Static RAM.
- These can be generally designed with BJT / FETs
- BJT = Bipolar Junction Transistor
- FET = Field Effect Transistor



→ BJT is faster and costlier than FET.

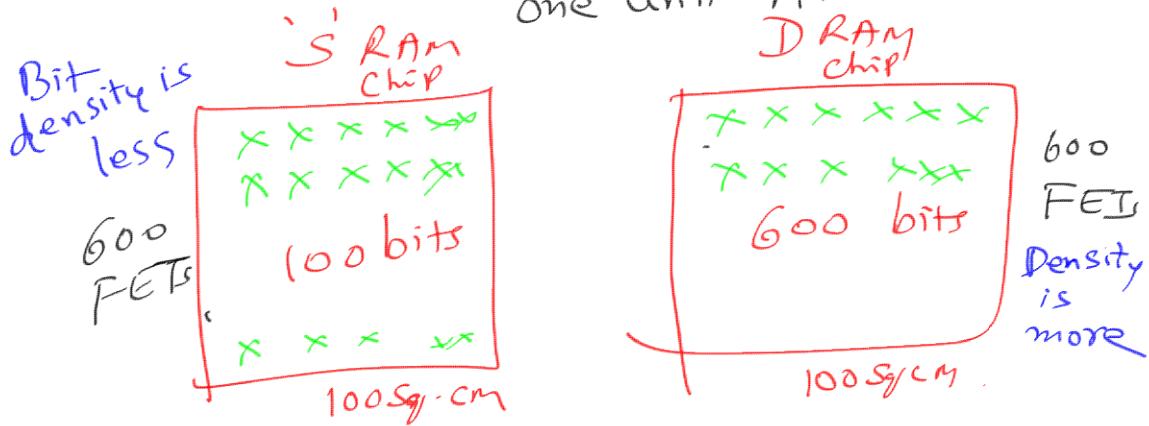
→ The memories in Superfast Computers are designed with BJT and the memories in General purpose Systems are designed with FETs.



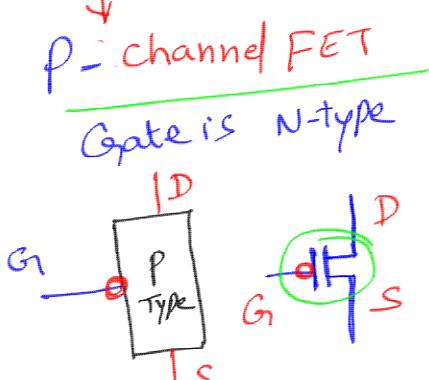
SNO	parameter	'D' RAM	'S' RAM
1	No of FETs needed to store one bit	One	Six
2	Cost /bit	Cheaper	Costliest
3	Capacitor	is needed	is not Needed
4	Refresh logic circuit	is needed	is not Needed
5	Power Consumption	Smaller	Very High
6	Bit Density.	More	Less
7	Speed	Moderate	7 to 10 times to 'D' RAM
8	Application in Digital Computer	Main Mem ER 4GB, 8GB	Cache Mem



Bit density = No. of bits in one unit Area

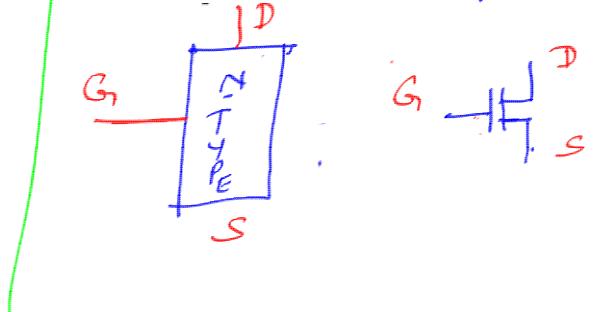


FET operation



N-channel FET

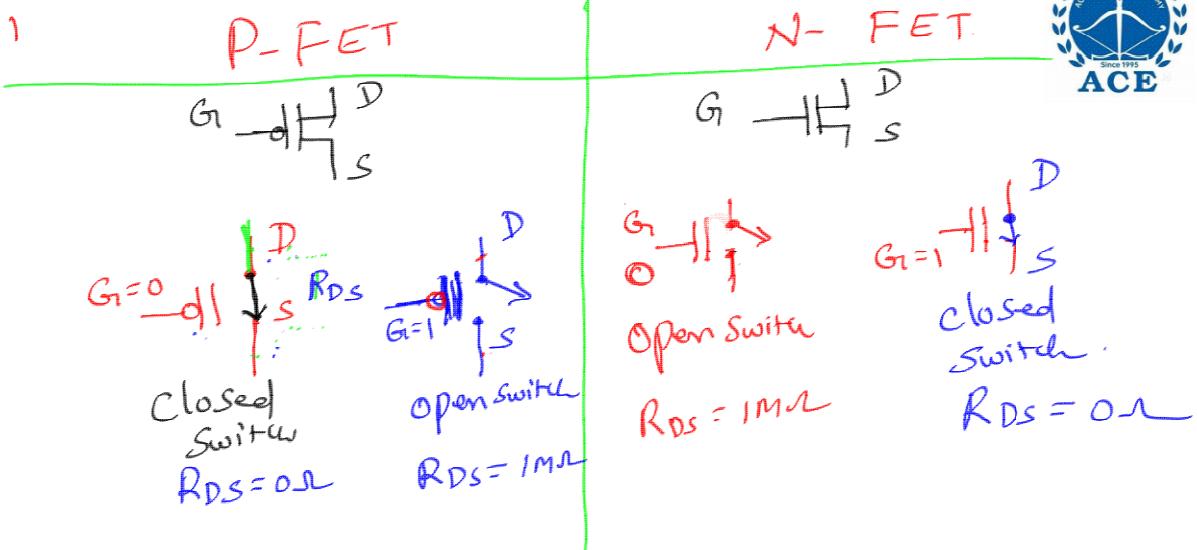
Gate is P-Type





FET is having 3 terminals
GATE is control input

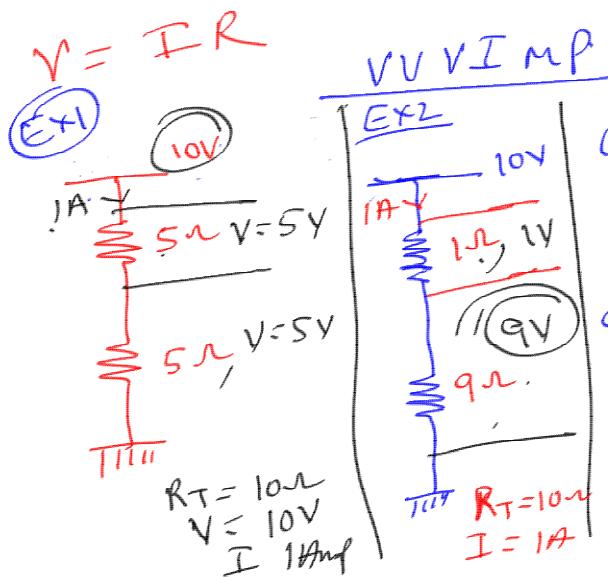
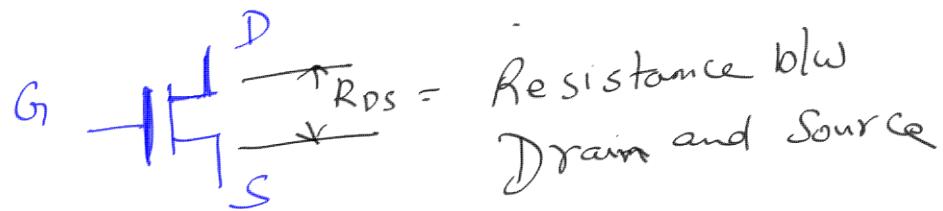
~~Drain~~
Source





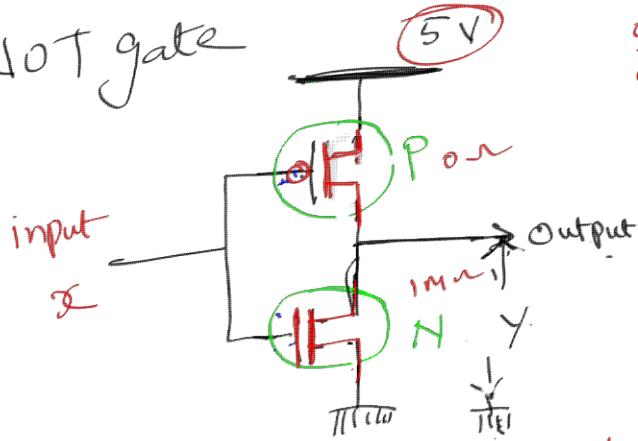
bit₀ = 0 to 0.6 V

bit₁ = 2.4 V to 5 V



When 2 Resistors are connected in Series then max Voltage is available across High Resistor Value.

NOT gate



5V = bit 1
0V = bit 0



(x)	P	N	(y)
0	ON	OFF	1
1	OFF	ON	0

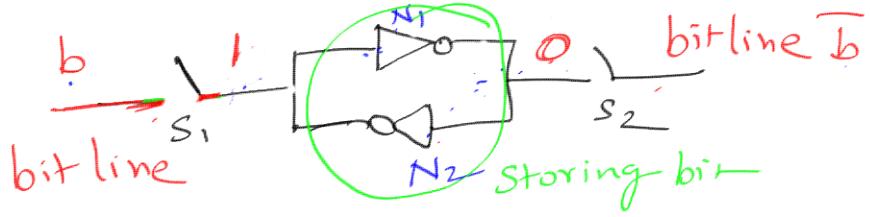
x	y
0	1
1	0

NOT gate

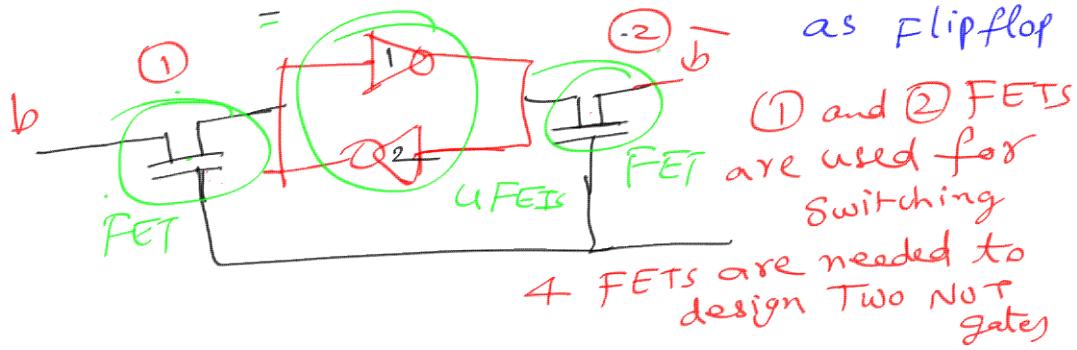
'S' RAM design for storing one bit

→ It requires 6 FETS; in these 6 FETS,
2 FETS are used for switching
and 4 FETS are used for
bit voltage storing





S' RAM is
also known
as flip flop



FET and BJT are used as
High Speed Electronic
switch.



Welcome to ACE Engineering Academy - online live class

Subject: **Computer Organization and Architecture**

Faculty: **Y.V. Ramaiah**

9866339106

Subject

Computer organization & Architecture

Chapters (Topics)

I. Computer Arithmetic ✓

II. Memory Organization

III. Secondary Memories

IV. Basic processor organization and Design

V. Pipeline organization

VI. Control unit Design

VII. IO Organization

Chapter 2 Memory Organization

- Introduction ✓
- Memory Basics ✓
- Memory Classification ✓
- Memory Size Expansion ✓
- Primary Memory
- Secondary Memory ✓
- ROM and its design ✓
- **RAM and its design** ✓
- Memory Hierarchy ✓
- Cache Memory
- Mapping Techniques
- Different misses occurred in cache
- Different block replacement techniques
- Tag directory design
- Associative Memory

Q. A 32-bit wide main memory unit with a capacity of 1 GB is built using $256M \times 4$ - bit DRAM chips. The number of rows of memory cells in the DRAM chip is 2^{14} . The time taken to perform one refresh operation is 50 nanoseconds. The refresh period is 2 milliseconds. The percentage (rounded to the closest integer) of the time available for performing the memory read/write operations in the main memory unit is $\frac{59}{60}$.

$$\begin{aligned} \text{Chip refresh time} &\approx 16\text{K} \\ &= 16384 \times 50\text{ns} \\ &= 0.8192 \text{ ms} \end{aligned}$$

$R/W \text{ time} = \frac{1}{2} - 0.8192 \text{ ms}$

$$\therefore \text{time}_{R/W} = \frac{1.1808 \text{ ns}}{2 \text{ ns}} = 59.04 \approx 59$$

$\approx 0.8 \text{ ms}, \approx 60\%$

Q. Assume that for a certain processor, a read request takes 50 nanoseconds on a cache miss and 5 nanoseconds on a cache hit.

Suppose while running a program, it was observed that 80% of the processor's read requests result in a cache hit. The average read access time in nanoseconds is 14.

$$\begin{aligned}
 T_{cm} &= 5 \text{ ns} & H &= 0.8 & \text{AMAT} \\
 T_{mh} &= 50 \text{ ns} \\
 n &= 50, \quad \text{Hit time} = \frac{(0 \times 5)}{200 \text{ ns}} = (0.8 \times 5) + \\
 \text{miss time} &= \frac{50 \text{ ns}}{50} = (0.2 \times 50) \text{ ns} \\
 & \quad \text{Total} = 700 \text{ ns} = 14 \text{ ns}
 \end{aligned}$$

$$50 \text{ instructions} = 700 \text{ ns}$$

$$\begin{aligned}
 \text{AMAT} &= \frac{700 \text{ ns}}{50} \\
 &= 14
 \end{aligned}$$





Q2
GATE

Q. A cache memory that has a hit rate of 0.8 has an access latency 10 ns and miss penalty 100 ns. An optimization is done on the cache to reduce the miss rate. However, the optimization results in an increase of cache access latency to 15 ns, whereas the miss penalty is not affected.

The minimum hit rate (rounded off to two decimal places) needed after the optimization such that it should not increase the average memory access time is _____.

$$H = 0.8, T_{cm} = 10 \text{ ns}$$
$$T_{mm} = 100$$

$$AMAT = (0.8 \times 10 \text{ ns}) + (0.2 \times 100) \text{ ns} = 28 \text{ ns}$$



After optimization

$$T_{cm} = 15 \text{ ns}$$

$$T_{mm} = 100 \text{ ns}$$

$$AMAT \leq 28 \text{ ns}$$

$$H = ?$$

$$AMAT \leq H_{cm} * 15 \text{ ns} + (1 - H_{cm}) * 100 \text{ ns}$$

$$28 \leq$$

$$28 \leq 15 H_{cm} + 100 - 100 H_{cm}$$

$$+72 \leq +85 H_{cm}$$

$$85 H_{cm} \geq 72$$

$$H_{cm} \geq \frac{72}{85} = 0.847$$

$$= 0.85 \checkmark$$

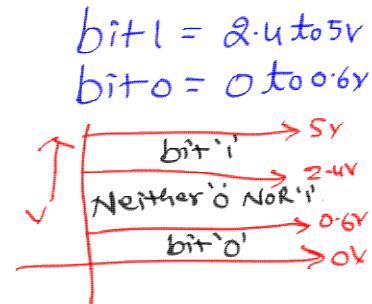
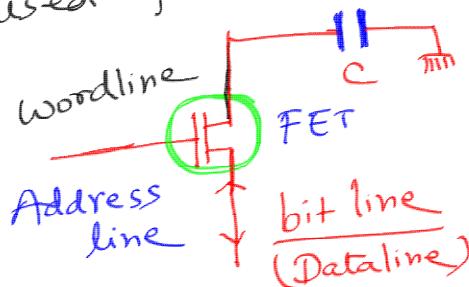


MSQ Q) In the give below, with reference to static RAM for storing one bit, the correct statement(s) is/are

- x a) Bit density is more than DRAM
 - ✓ b) Total 6 FETs BJTs are used
 - ✓ c) Power consumption is High Compared to DRAM
 - ✓ d) Memory Refresh logic circuit is not needed.
- b, c and d

One bit storage in DRAM

- It requires one FET and one capacitor
- FET is used for switching and Capacitor is used for bit storing





- Due to the capacitor nature; the bit voltage across the capacitor gets discharge gradually; when it falls down to 2.4 Volts then bit '1' becomes bit '0' which may cause lot of disturbance in the program.
- To prevent this discharge; one separate memory refresh logic circuit is used for each DRAM chip; this circuit refreshes all the capacitors which are used to store bit; will be recharged.

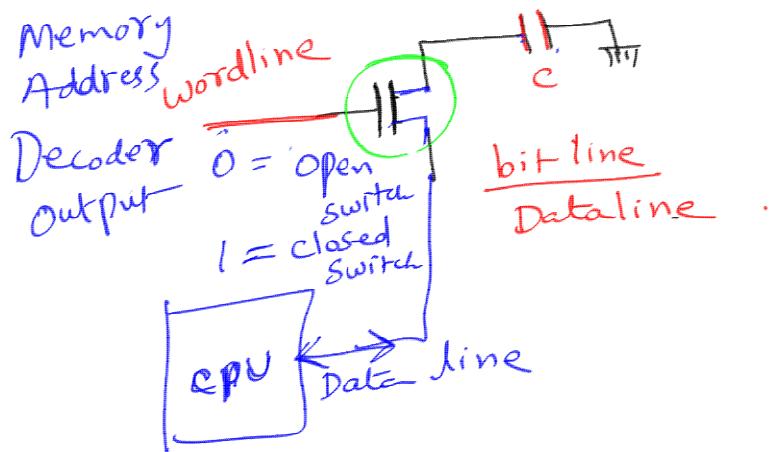


- One refresh operation refreshes all the capacitors available in one Row.
- Since each DRAM chip is having separate Refresh circuit; Refreshing operation is performed in parallel; one chip refresh time is sufficient to refresh all the chips available in the system.
- Total no of Refresh operations needed to refresh one chip = No of Rows
- Chip refresh time = No of Rows * ^{one} Refresh time



(Ex) In a $1K \times 16$ bit 'D' RAM chip, one Refresh operation takes 10 n.s. , total amount of time required to Refresh the entire chip is $10 \cdot 24$ micro second.

$$\begin{aligned} \text{No. of Rows} &= 1K = 2^{10} = 10^3 \\ \text{Chip refresh time} &\approx 10^3 \times 10 \times 10^{-9} \text{ sec} \\ &\approx 10 \mu\text{sec} \\ &= 1024 \times 10 \times 10^{-9} \text{ sec} \\ &= 10.24 \mu\text{sec} \end{aligned}$$





→ while performing READ/WRITE operation for certain Row, all the Capacitors will be automatically refreshed

→ Refresh time = $\frac{\text{Amount of time needed}}{\text{Row/Chip}}$ for only Refresh operation

Refresh period of Row :- Refresh time + RD/WR time of a Row

Refresh period of a chip :- Refresh time of the given chip + RD/WR time of the chip

(Ex)

one Refresh operation takes 100ns in 4Kx8 bit DRAM chip, Chip refresh period is 1200 micro second; the % of time used to perform RD/WR

$$\text{No. of Rows} = 4K = 4 \times 10^3 = 4096$$

chip is _____ Chip refresh time = $4096 \times 100 \text{ ns}$
 $= 409.6 \mu\text{sec}$

$$\text{chip RD/WR} = 1200 - 409.6 \mu\text{sec}$$

$$\% \text{ of time} = \frac{790.4 \mu\text{sec}}{1200 \mu\text{sec}} \times 10 = 65.8666 = \underline{\underline{65.87}}$$



Approximate Value

$$\text{Chip Refresh time} = \frac{4000 \times 100ns}{4} = 100\mu\text{sec}$$
$$\text{Chip Read time} \approx 80\mu\text{sec}$$

$$\% \text{ of time} = \frac{80\mu\text{sec}}{1200\mu\text{sec}} \times 100$$

$$\approx 66.666 \approx \underline{\underline{66.67}}$$

✓ Mark will be awarded; for the value
65 to 67.

Memory Hierarchy



- It is the combination of all different Memories used in the system from costliest to cheapest i.e. (fastest to slowest)
- Faster Memories are costlier and slower Memories are cheaper.
- In a System, fastest memory (Cache) is used to store current executing file and slower Memories are used for storing back-up files.

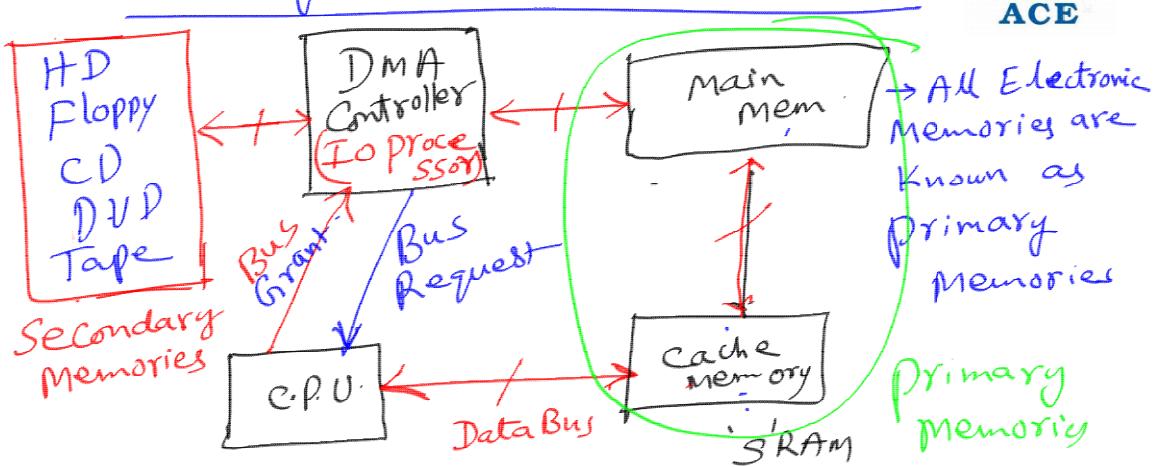


VUV Imp

In order to control the overall system design cost, we use smaller size Cache memory and larger size Secondary Memories.



Memory Hierarchy Basic diagram





- During program execution, the CPU first visits cache memory, let target word is available in C.M; the CPU directly reads that word and that operation is known as Cache Hit.
- when the target word is not available in Cache the operation is known as Cache Miss; then CPU visits main memory



- Let target word is available in M.M then CPU transfers (Maps) the Associative block from M.M to C.M then only CPU reads the target word.
- Let the Target word is not available in M.M, then O.S. takes responsibility to initialize the DMA Controller for transferring the target words from Secondary Mem. to M.M.



- Generally buses are under the control of CPU; for asking the Bus Control, DMAc sends Bus request Signal to the CPU.
- After Receiving Bus request Signal, the CPU first finished the present operation and then handovers all the Buses to DMAc by sending Bus Grant Signal.



→ word access time in Secondary Memories is in millisecond (because of R/w Head and motor); hence CPU can't communicate with these memories.
To transfer data b/w Main Memory and Secondary Memory; DMA Controller (I/O processor) is used.



- In Memory Hierarchy; the fastest memory is electrically closer to the C.P.U. and Cache occupies top position in the memory Hierarchy.
- Main Memory occupies central position and Secondary Memory occupies least position in memory Hierarchy.



- After Receiving Bus Grant Signal, the DMA starts data transfer operation.
- During DMA operation, the CPU neither READ nor writes.
- After Completion of DMA operation, all the buses will be automatically relinquished to C.P.U.



block = Group of words .

mapping:- It is the process of transferring the block of words from Main memory to Cache memory when miss has occurred in Cache



Hit:- When the Searched word is available in Searched memory, the operation is known as Hit

Otherwise Miss .

Hit Rate is directly proportional to the size of the space .

→ Always Highest level memory provides 100% hit rate .



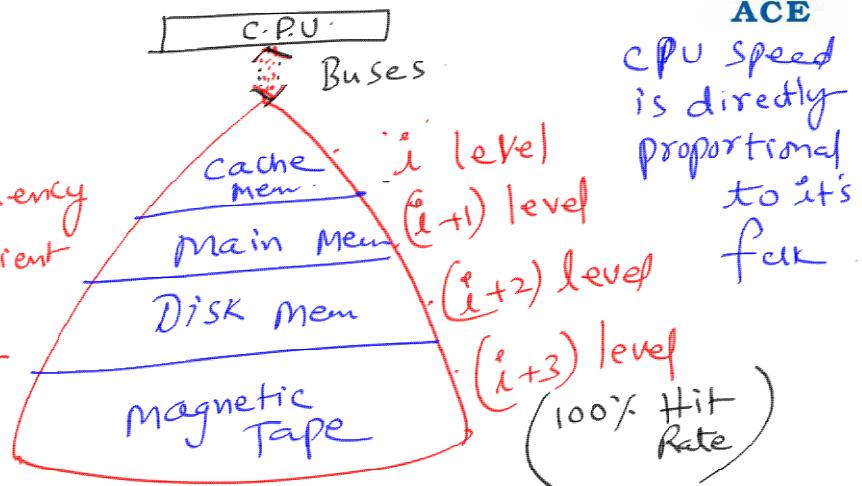
Hit Rate = $\frac{\text{No. of Hits}}{\text{Total no. of CPU References}}$

Generally Secondary Memory is Highest level memory - It gives 100% Hit Rate

Memory Hierarchy in Pyramid Structure

$$f_{clk} = \frac{1}{T}$$

High frequency
CPU is efficient
than low
frequency
CPU



CPU speed
is directly
proportional
to its
 f_{clk}



f = Speed

S = size in
no. of words

H = Hit Rate

C = Cost / bit

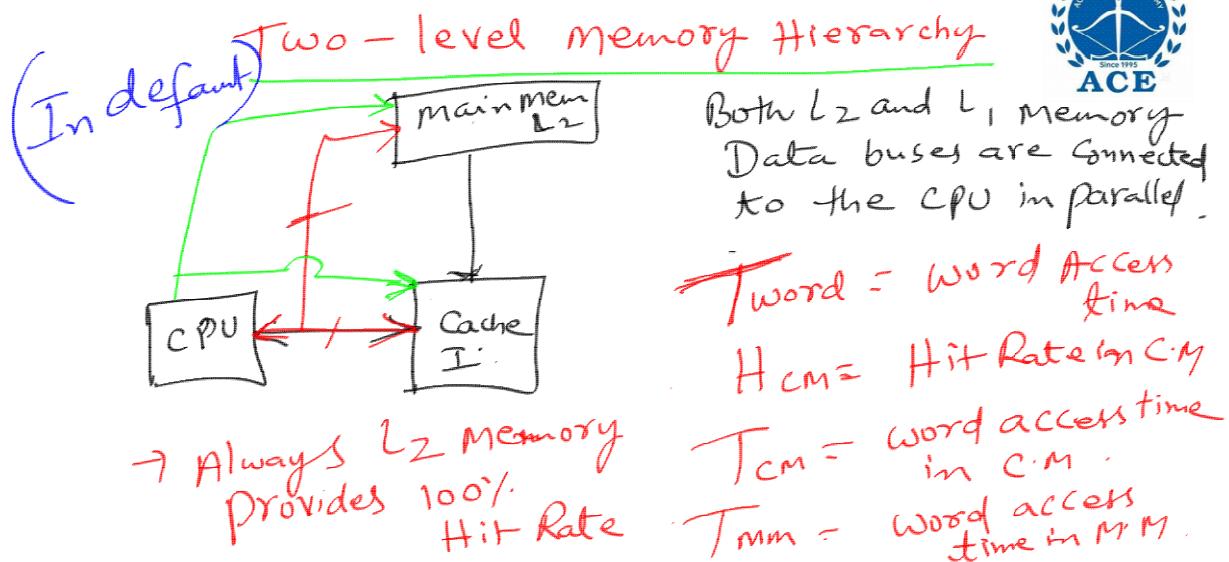
Note:- Highest level memory produces
100% Hit Rate.

$C_i > C_{i+1} > C_{i+2} > C_{i+3}$

$f_i > f_{i+1} > f_{i+2} > f_{i+3}$

$S_i < S_{i+1} < S_{i+2} < S_{i+3}$

$H_i < H_{i+1} < H_{i+2} < H_{i+3}$





$$\text{Let } T_{cm} = 1 \text{ ns}$$

$$T_{mm} = 10 \text{ ns}$$

$$H_{cm} = 90\%$$

Average Memory Access time

$$= \frac{\text{Average word Access time}}{\text{Access time}}$$

$$\text{let } n = 100$$

n - length of the program

$$H = 90, M = 10$$

$$\begin{aligned} \text{Total Cache access time} \\ = 90 \times 1 \text{ ns} = 90 \text{ ns} \end{aligned}$$

Total M.M. access time

$$= 10 \times 10 \text{ ns} = 100 \text{ ns}$$

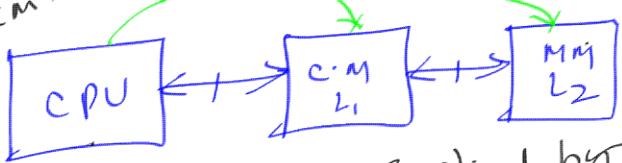
$$T_{100} = (90 + 10) \text{ ns} = 190 \text{ ns}$$

$$A.M.A.T = \frac{190 \text{ ns}}{100} = 1.9 \text{ ns}$$

$$\begin{aligned} A.M.A.T &= (H_{cm} * T_{cm}) + (1 - H_{cm}) * T_{mm} \\ &= (0.9 * 1 \text{ ns}) + (0.1 * 10 \text{ ns}) \\ &= \underline{\underline{1.9 \text{ ns}}} \end{aligned}$$



$T_{CM} = 1 \text{ ns}$ Two level memory hierarchy
 $T_{MM} = 10 \text{ ns}$ (Serial)
 $H_{CM} = 90\%$



words are supplied by
L₁ Mem. only

$$\text{AMAT} = \left(H_{CM} * T_{CM} \right) + \left(1 - H_{CM} \right) * \left(T_{MM} + T_{CM} \right)$$

Total C.M. access
time = $90 * 1 \text{ ns} = 90 \text{ ns}$

Total M.M. access
time = $10 * (10 + 1) = 110 \text{ ns}$

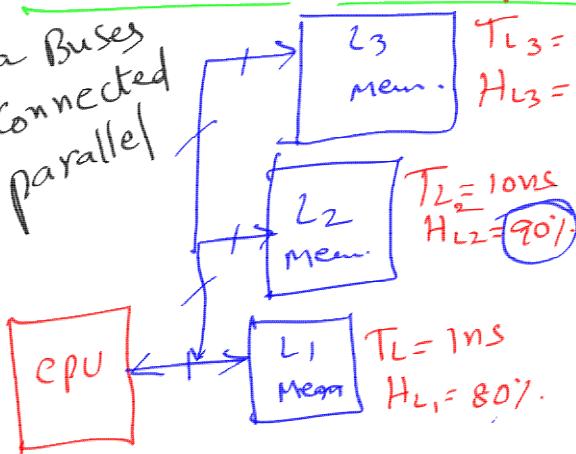
$$T_{100} = 200 \text{ ns}$$

$$\text{AMAT} = 2 \text{ ns}$$



Three level memory hierarchy

All
Data Buses
are Connected
in parallel



$n=100$

$\text{Total L}_1 \text{ time} = 8.0 \times 1 \text{ ns} = (8 \text{ ns})$

$\text{Total L}_2 \text{ time} = 18 \times 10 \text{ ns} = 180 \text{ ns}$

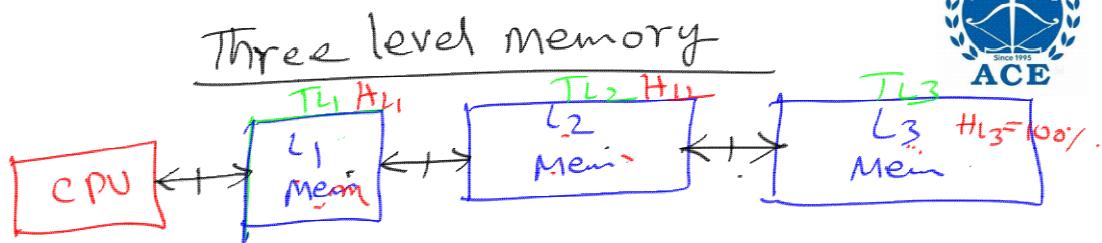
$\text{Total L}_3 \text{ time} = 2 \times 100 \text{ ns}$

$T_{100} = (80 + 180 + 200) \text{ ns} = 460 \text{ ns}$

$\text{AMAT} = 4.6 \text{ ns}$



$$\begin{aligned}
 \text{AMAT} &= (H_{L1} * T_{L1}) + (1 - H_{L1}) * H_{L2} * T_{L2} \\
 &\quad + (1 - H_{L1}) * (1 - H_{L2}) * T_{L3} \\
 &= (0.8 * 1\text{ns}) + (0.2 * (0.9) * 10\text{ns}) \\
 &\quad + 0.2 * 0.1 * 100\text{ns} \\
 &= 0.8\text{ns} + 1.8\text{ns} + 2\text{ns} = \underline{4.6\text{ns}}
 \end{aligned}$$



only L₁ supplies the words to CPU

$$\begin{aligned}
 \text{AMAT} &= (H_{L1} * T_{L1}) + (1 - H_{L1}) * H_{L2} * (T_{L2} + T_{L1}) \\
 &\quad + (1 - H_{L1}) * (1 - H_{L2}) * (T_{L3} + T_{L2} + T_{L1})
 \end{aligned}$$



Mapping

- mapping is the process of transferring the block of words from M·M to CM when miss has occurred in Cache only
- Block is the group of words.
- Block transfer is preferred for reducing the future miss penalties

Different mapping techniques

- (i) Direct mapping (Cheapest)
- (ii) Block Set Associative mapping
- (iii) Fully Associative mapping
(Costliest)



Main Memory = physical Memory.

MMW = No. of words in main memory.

CMW = No. of words in Cache Memory.

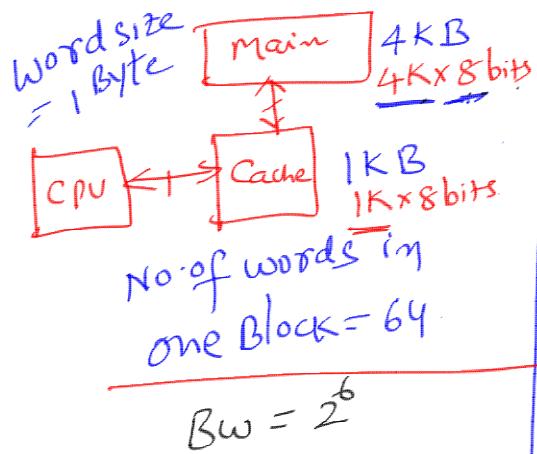
BW = No. of words in one Block.

MB = No. of blocks in Main memory.

No. of blocks in Cache memory -

CB
 CL

Sometimes one Cache block is
also known as one Cache LINE



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



MA = PA = main memory Address Size

(12) bits = Physical Address Size

$$CMW = 1K = 2^{10}$$

CA = 10 bits

Cache memory word Address size

$$MM \text{ Capacity} = 4K \times 8$$

$$CM \text{ Capacity} = 1K \times 8$$

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

$$CMW = 1K = 2^{10}$$

$$PA = 12, CA = 10$$

$$TAG = \underline{(PA - CA)}$$

2 bit

$$BW = 64 = 2^6$$

$$CB = \underline{CL} = \frac{CMW}{BW}$$

$$= \frac{2^{10}}{2^6} = \underline{\underline{2^4}}$$



CL_0
 k_0
 CL_{15}

$$CMW = CL \cdot BW$$

$$MB = \frac{MMW}{BW}$$

$$= \frac{2^{12}}{2^6} = \underline{\underline{64}}$$

MB_0 to MB_{63}