

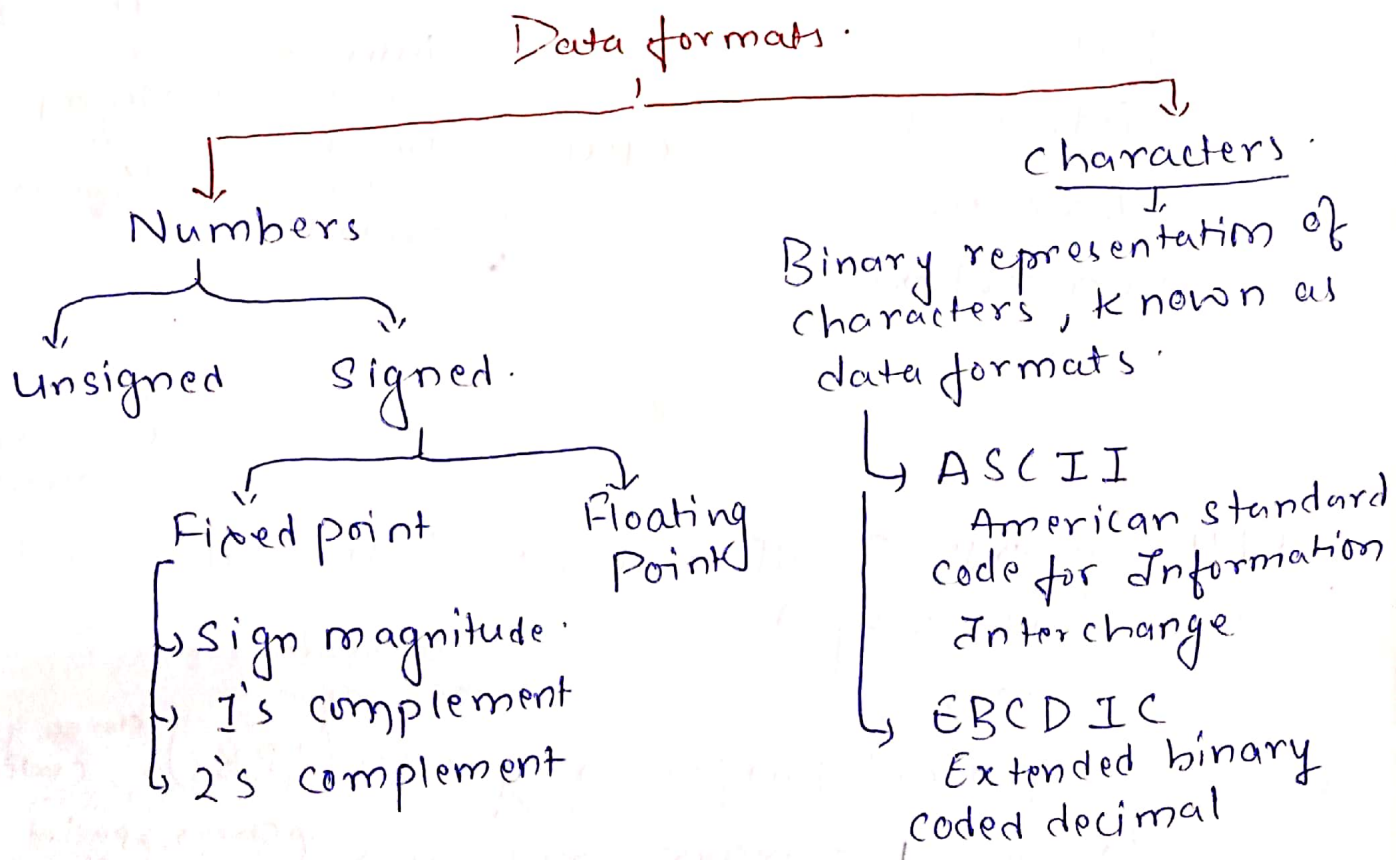
Computer Organization and Architecture.

Computer Architecture:- Conceptual design & fundamental operational structure of a computer is called computer Architecture.

Computer Organization:- It deals with physical devices and their interconnection with a perspective of improving the performance.

Things Included:

- | | |
|------------------------------|------------------------------|
| Computer Architecture | Computer Organization |
| ↳ CPU design | ↳ I/O organization |
| ↳ Instructions | ↳ Memory organization |
| ↳ Addressing modes | ↳ performance |
| ↳ Data formats | ↳ Instruction activity |
| | ↳ Pipeline |



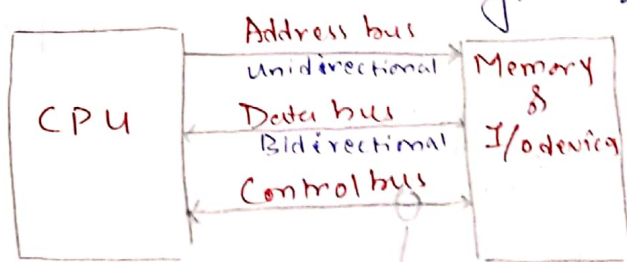
Computer's components:-

- ↳ CPU
 - ↳ control unit
 - ↳ ~~ALU~~ ALU (Arithmetic Logic unit)
 - ↳ Memory
 - ↳ Primary or main memory
 - ↳ Secondary or auxiliary memory
 - ↳ I/O devices
 - ↳ Input devices
 - ↳ Output devices
- Other components.
- System bus.
 - CPU register.

System bus: ~~connection~~ collection of communication line within the computer system.

In the system we have three types of bus.

- ↳ Address bus
- ↳ Data bus
- ↳ control bus



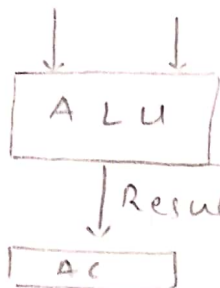
Individual lines are unidirectional

CPU registers: Small memories inside the CPU are called registers

Types:

- ↳ General purpose registers (GPRs) (stores general contents)
- ↳ Special purpose registers (stores special contents)
 - ↳ Accumulator
 - ↳ Status or flag register
 - ↳ Program counter
 - ↳ Address register
 - ↳ Instruction Registers
 - ↳ Data register
 - ↳ Stack pointer

→ **Accumulator**:- It is used to store the result of ALU and sometimes one of the input operand for ALU also.

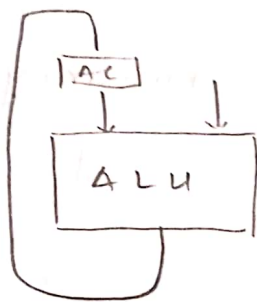


* If both the inputs are taken from Registers they are called "Register based architecture"

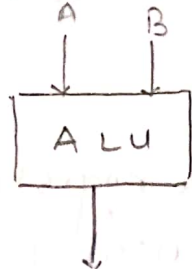
* If both the inputs should be taken from stack it is called "Stack based architecture"

* If out the input is coming from either Register or memory simultaneously it is called complex-architecture

* If out of two inputs one is necessarily coming from accumulator it is called accumulator based architecture.



32 bits 32 bits

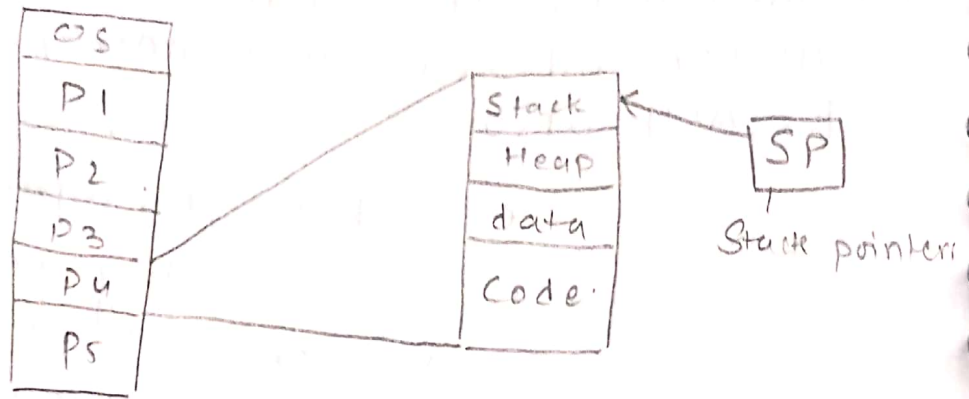


32-bits architecture

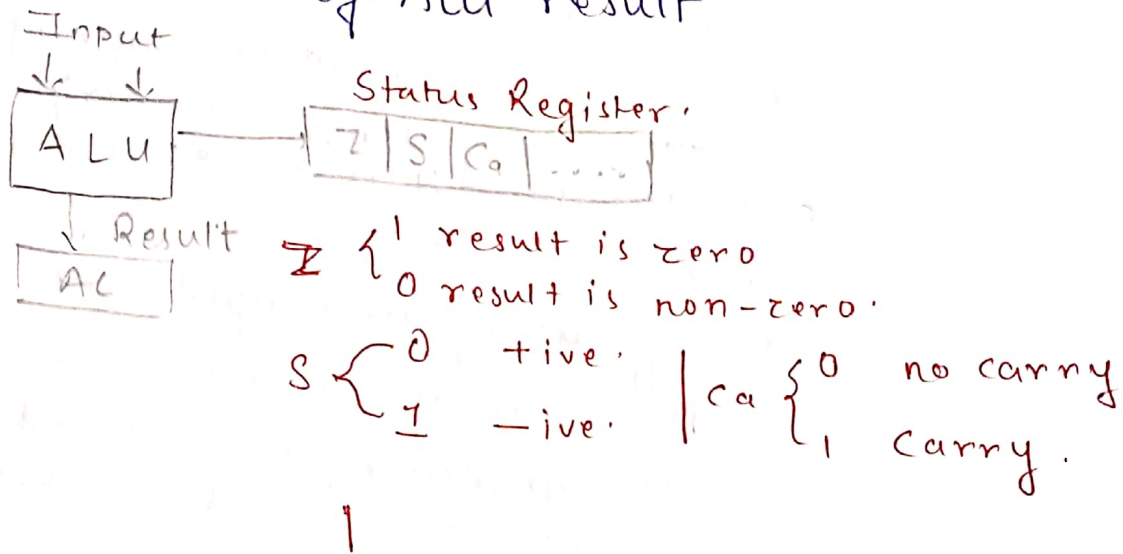
or
32-bits CPU

↓
This will give 1 word size

- **Program Counter**:- It stores the address of next instruction to be executed.
- **Instruction Register**:- It is used to store the current instruction which CPU executes; It is used to store the instruction before it goes for decode.
- **Stack pointer**:- It is used to store the address of the top of the process stack.

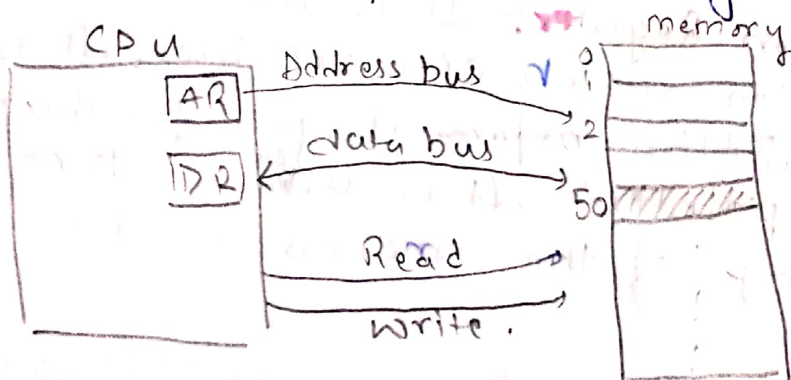


→ **Flag or Status register:** It is used to store the status of ALU result



→ **Memory Address Reg:-** It is used to send address to memory (MAR)

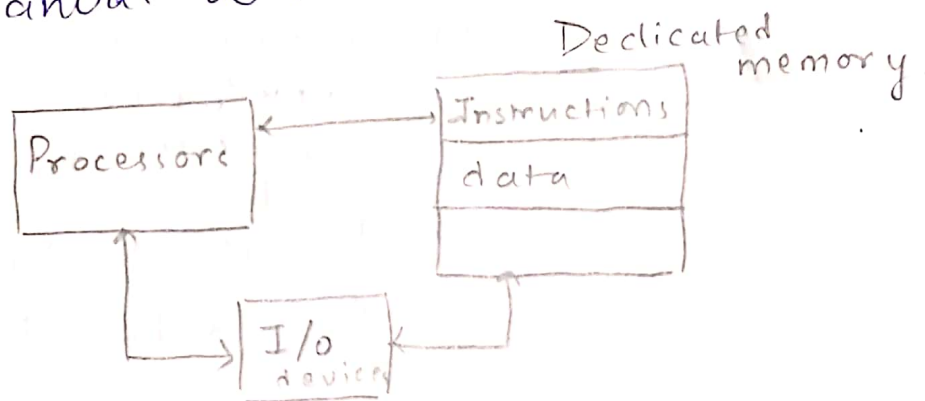
→ **Data register or memory Data register or memory buffer register:-** It is used to send data to memory (memory write) and to receive data from memory (memory read)



* Von neumann architecture:-

There is a dedicated memory connected to a processor in which instructions and data are stored.

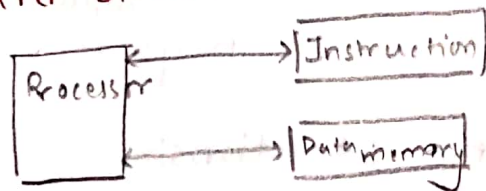
The processor can fetch the instruction & data from this memory, can execute the instructions using the data. & can copy the result in the memory without any manual work.



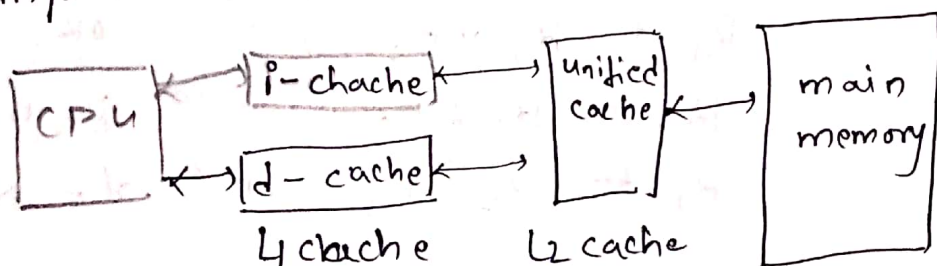
Bottle neck of vonneumann's model:-

The instructions and data both cannot be fetched from memory simultaneously. Hence to solve this problem scientist gave another architecture 'Harvard Architecture'.

Harvard architecture.



Improvised Harvard architecture



Stored program architecture.

⇒ Micro-operation :-

The operations performed on the data stored in the registers

Types of microoperation.

↳ Register transfer.

↳ Register to register $(R_1 \leftarrow R_2)$
content transfer $(R_2 \rightarrow R_1)$

↳ Memory access.

↳ Read: $Reg \leftarrow Memory$

↳ Write: $DR \leftarrow M[Address]$

$Memory \leftarrow Reg$ or $DR \leftarrow M[100H]$
or $M[Address] \leftarrow DR$ or $DR \leftarrow M[AR]$

↳ Arithmetic micro-operation :-

$$R_1 \leftarrow R_2 + R_3 \quad R_1 \leftarrow R_1 - 1$$

$$R_1 \leftarrow R_2 - R_3 \quad R_1 \leftarrow \overline{R_2} \text{ (1's complement)}$$

$$R_1 \leftarrow R_1 + 1 \quad R_1 \leftarrow \overline{R_2} + 1 \text{ (2's complement)}$$

$$R_3 \leftarrow R_1 + \overline{R_2} + 1 \text{ (Addition with 2's complement)}$$

2's complement means its negative representation

↳ Logical micro-operations:-

$$R_1 \leftarrow R_2 \wedge R_3 \quad \text{Logical AND}$$

$$R_1 \leftarrow R_2 \vee R_3 \quad \text{" OR}$$

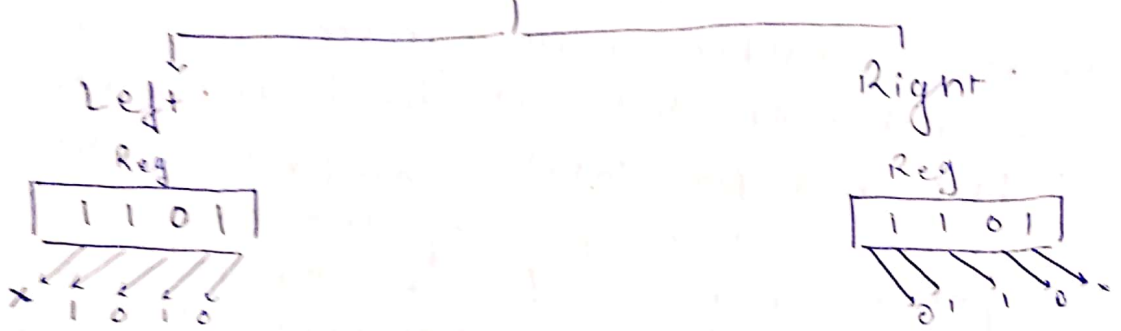
$$R_1 \leftarrow R_2 \oplus R_3 \quad \text{" XOR}$$

$$R_1 \leftarrow \overline{R_2 \oplus R_3} \quad \text{" NOR}$$

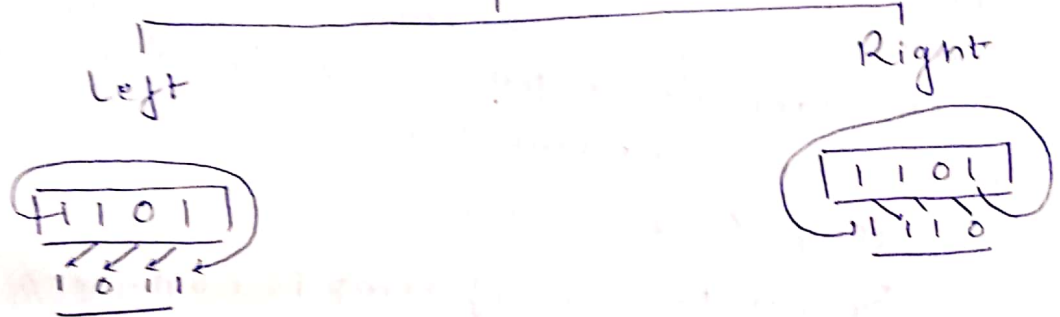
Shift micro-operations

- ↳ Logical shift
- ↳ Circular shift (Rotation)
- ↳ Arithmetic

Logical shift

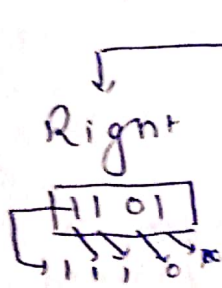


Circular shift



Arithmetic shift: the shift operation is performed on arithmetic numbers (signed no.), after the shift operation sign of the no. should remain same.

Arithmetic

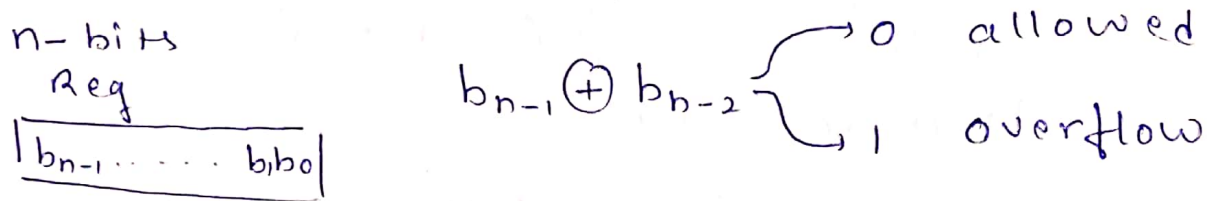


Left :- It is same as logical left shift but it is performed only after CPU ensures that sign will not change after the operation.

eg:- $\begin{array}{r} 1100 \\ \swarrow \searrow \\ 1000 \end{array}$

eg:- $\begin{array}{r} 1010 \\ \downarrow \\ \text{not-allowed} \\ \downarrow \\ \text{arithmetic-left} \\ \downarrow \\ \text{shift overflow} \end{array}$

to ensure the operation is allowed or not CPU uses \oplus XOR operation on the MSB



Question :- Consider a new instruction i.e "branch-on-bit-set" (bbs), the instruction "bbs reg, pos, label" jumps to label if bit in position pos of the register operand reg is '1'. a register is 32 bits wide and bits are numbered 0 \rightarrow 31, bit in position '0' been least significant. consider the following imulation of this instruction on a processor that does not have bbs implemented.

$$\text{temp} \leftarrow \text{reg} \& \text{mask}$$

Branch to label if temp is non-zero

The variable temp is a temporary register. for correct imulation the variable "mask" must be generated by .

(a) $\text{mask} \leftarrow 0x1 \ll \text{pos}$

(b) $\text{mask} \leftarrow 0xFFFFFFFF \gg \text{pos}$

(c) $\text{mask} \leftarrow \text{pos}$

(d) $\text{mask} \leftarrow 0xF$

$x \ll y$ means shifting x in left direction y times.

Question

Pg-42

$$R_1 \leftarrow M[3000]$$

$$\text{loop: } R_2 \leftarrow M[R_3]$$

$$R_2 \leftarrow R_1 + R_2$$

$$M[R_3] \leftarrow R_2$$

$$R_3 \leftarrow R_3 + 1$$

$$R_1 \leftarrow R_1 - 1$$

Branch on nonzero
Halt

$$R_1 = \boxed{\begin{matrix} 109 \\ 107 \end{matrix}}$$

$$R_2 = \boxed{\begin{matrix} 15010 \\ 109 \\ 100 \end{matrix}}$$

$$R_3 = \boxed{\begin{matrix} 2000 \\ 2001 \\ 2002 \end{matrix}}$$

1000	Prog. Instru
2000	110
2001	109
2002	108
2003	107
2004	106
2005	105
2006	104
2007	103
2008	102
2009	101
2010	100
	⋮
3000	10

Q3 memory references

$$R_1 \leftarrow M[3000] \text{ — 1 times = 1 bit}$$

$$\left. \begin{matrix} R_2 \leftarrow M[R_3] \\ M[R_3] \leftarrow R_2 \end{matrix} \right\} \text{ — 10 times = } 10 \times 2 = 20 \text{ bit}$$

$$20 + 1 \text{ bit} = 21 \text{ bit}$$

Q4 Content of memory location

$$2010 = '100'$$

Q5 If CPU receives interruption

during the execution of any instruction then CPU completes the execution of current instruction, then CPU stores the address of next instruction onto the stack as return address and then CPU goes for interrupt service.

$$1000 \quad R_1 \leftarrow M[3000]$$

$$1008 \quad R_2 \leftarrow M[R_3]$$

$$1012 \quad R_2 \leftarrow R_1 + R_2$$

$$1016 \quad M[R_3] \leftarrow R_2$$

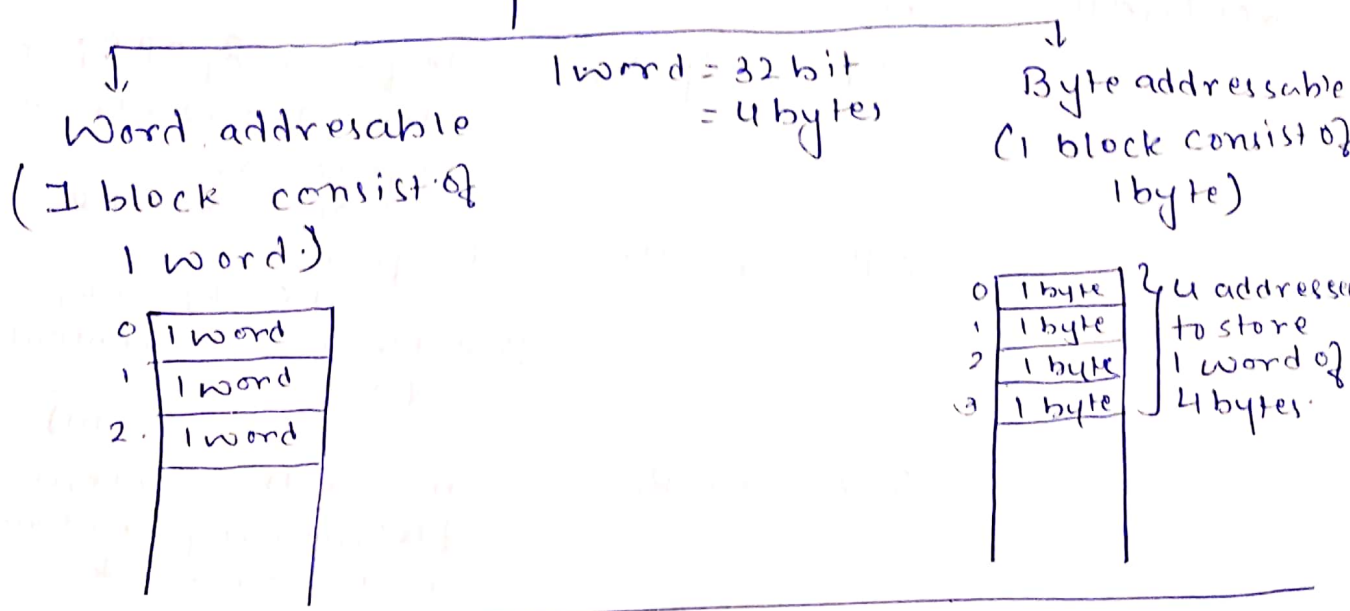
$$1020 \quad R_3 \leftarrow R_3 + 1$$

$$1024 \quad R_1 \leftarrow R_1 - 1$$

$$1028 \quad \text{Branch}$$

As byte addressable uses 4 blocks for 1 word. hence the answer will be $1000 + 8 + 4 + 4 + 4 + 4 = 1024$

Memory



Q1
P42

1 word = 8 bit = 1 byte

$A_0 \Rightarrow$	00100011	C ₀
	1C ₀ 0010001	1
	1C ₀ 0010000	1
	11C ₀ 001000	0
	011C ₀ 0010	0
	0011C ₀ 001	0
	00011C ₀ 00	1
	100011C ₀ 0	0
	0100011C ₀	0

Carry

B → 0XZ.3
C → ~~XXXXXXXXXX~~0

Q₁
(b) no. of 1 bits in A₀

Q₂
(a) RRC A, #1:
Right rotate A through carry by one bit.

Question
2007 GATE

In a simplified computer the instructions are:

OP R_j, R_i — Performs $R_j \text{ OP } R_i$ & stores the result in Register R_i
 $R_j \text{ OP } R_i \rightarrow R_i$

OP m, R_i — Performs $\text{val OP } R_i \rightarrow R_i$
[val denotes the content of memory location m]

MOV m, R_i — moves the content of ~~the~~ memory location 'm' to register R_i

MOV R_i, m — moves the content of R_i to memory location m

The computer has only two registers and 'OP' is either ~~add~~ ADD or SUB. Consider the following basic block.

$$t_1 = a + b$$

$$t_2 = c + d$$

$$t_3 = e - t_2$$

$$t_4 = t_1 - t_3$$

assume that all operands are initially in memory.

The final value of computation should be

in memory what is the minimum no. of MOV instruction in the CODE generated for the basic block.

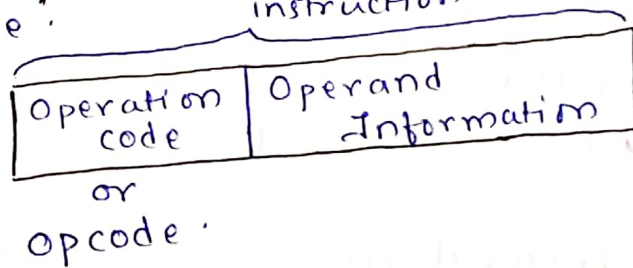
Answer:-

<u>MOV b, R₁</u>	$b \rightarrow R_1$
ADD a, R ₁	$a + R_1 \rightarrow R_1$
<u>MOV D, R₂</u>	$D \rightarrow R_2$
ADD c, R ₂	$c + R_2 \rightarrow R_2$
SUB e, R ₂	$e - R_2 \rightarrow R_2$
SUB R ₁ , R ₂	$R_1 - R_2 \rightarrow R_2$
<u>MOV R₂, x</u>	$R_2 \rightarrow x$

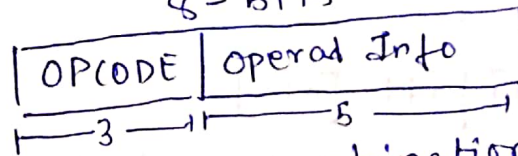
⇒ Instructions:-

definition: a group of bits which instructs computer to perform some operation.

"Instructions are binary statements which is produced by compiler in a machine language or byte code or binary code".



eg:- For a system
8-bits



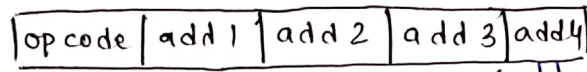
opcode combination = $2^3 = 8$

∴ Max no. of operations possible = 8
no. of instruction (different type of instructions) supported by the system = 8

If 12 opcode combination is required then bits of opcode will be = 4

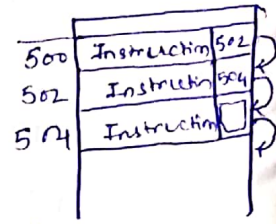
⇒ 4-address instruction:-

Within an instruction maximum 21 address can be specified



operands

↓
address of next instruction.



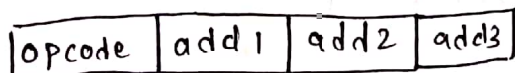
This type of instructions were used in the type of CPU that does not have program counter register and hence the last address is used to point to next instruction. ~~These type~~ and the rest all 3 addresses could be memory or register operand.

⇒ Disadvantages

- Large size instruction.
- Due to this more memory is required to store the program.
- If the program is relocated then the 4th address part of the instructions should be updated accordingly and updation of instruction is very costly operation.

⇒ 3-address instruction.

maximum three addresses can be specified within an instruction.



Program counters used in such types of computer this stores the address of next instruction.

Hence the disadvantages of above cases are resolved and this turns into advantages of 3-address instructions.

- lesser size instruction
- no updation of addresses or relocation

ex:- $\frac{10101}{\text{ADD}} \quad \frac{10}{R_2} \quad \frac{11}{R_3} \quad \frac{01}{R_1} \quad R_2 \leftarrow R_3 + R_1$

* All modern days computer uses 3-address instruction.

→ 2-address instructions

Maximum two addresses can be specified in an instruction

example :- $\frac{10101}{\text{opcode}} \quad \frac{10}{R_2} \quad \frac{11}{R_3} \quad R_2 \leftarrow R_2 + R_3$

one operand is used as source and destination both.

Disadvantages:

- The value of common operand is updated by the result of the operation
- More no. of instructions for a program as compared to 3-address instruction format.

example :- $x = (A+B) * (C+D)$

3 address

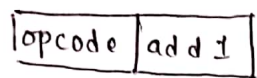
$R_1 \leftarrow A + B$
 $R_2 \leftarrow C + D$
 $x \leftarrow R_1 * R_2$

2 address instruction:

• $R_1 \leftarrow A$
 • $R_1 \leftarrow R_1 + B$
 • $R_2 \leftarrow C$
 • $R_2 \leftarrow R_2 + D$
 • $R_1 \leftarrow R_1 * R_2$
 • $x \leftarrow R_1$

→ 1 - address instruction

only one address can be specified in an instruction.



eg:- $\frac{10101}{\text{opcode}} \quad \frac{11}{R_2} \quad AC \leftarrow AC + R_2.$

accumulator is used as 2nd operand implicitly (within the instruction we don't need to mention it)

~~eg~~ disadvantages:

→ more no. of instructions for a program as compared to two-address instruction format.

eg:- $x = (A+B) * (C+D)$

- $AC \leftarrow C$
- $AC \leftarrow AC + D$
- $R_1 \leftarrow AC$
- $AC \leftarrow A$
- $AC \leftarrow AC + B$
- $AC \leftarrow AC * R_1$
- $x \leftarrow AC$

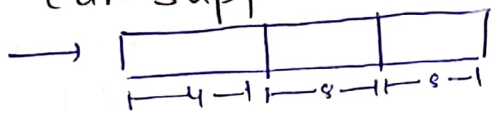
→ 0 - address instruction

no any address is specified within an instruction

Such type of instruction are used in stack-based architecture.

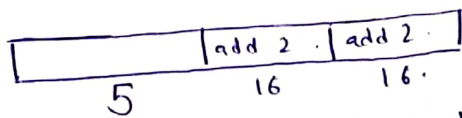
eg:- $!0101$ 2 operands are taken from stack

question a digital computer has 8-bit addresses and 20-bit instructions, minimum and maximum how many 2-address instructions the computer can support.



maximum instruction = $2^4 = 16$
 minimum " = 1

question-a computer supports 25 2-address instructions if the address size is 16-bit then the size of an instruction in bits is?



opcode bits = $\lceil \log_2(\text{no. of instruction supported}) \rceil$

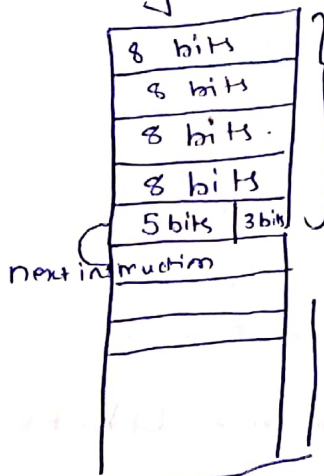
opcode bits = $\lceil \log_2(25) \rceil = 5$

Size of instruction = $16 + 16 + 5 = 37$ bits.

→ In above question if the instructions are stored in the memory in byte align fashion then the amount of memory (in bytes) to store the program that consists 100 instructions is?

byte aligned fashion.

⇒ 100 instructions will take 500 bytes of instruction
 $37 \text{ bits} \cong 5 \text{ byte}$

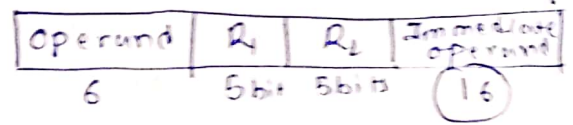


ISA :- Instruction set architecture
 collection of all the instructions supported by a processor (CPU)
 that collection is called Instruction set architecture

Question 17:
Page 44

no. of instructions = 40.
opcode - bits = 6 - bits

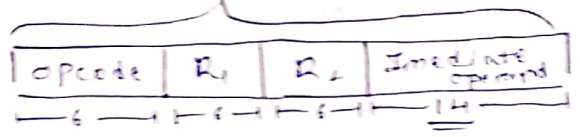
We have 24 registers hence to store these many registers we will require 5 bits ($R_0 - R_{23}$) i.e. 5 bit numbering is required to store 24 registers



Question 15 32 bit architecture means ALU.
Page 44 take 32 bit at once
and 1 word = 32 bits

Instruction length = 1 words

no. of register bits (operand bits) = 64 = 6 - bits
op code bits = 6 - bits

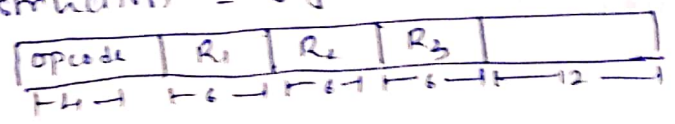


14 bits $\Rightarrow 2^{14} - 1$ values = 16383 values

Question 18
Page 44.

Given :- 64 registers.
register reference = 6.
Instruction set = 12.
i.e. 12 instructions are possible.
 \therefore instruction bit = 4 bits

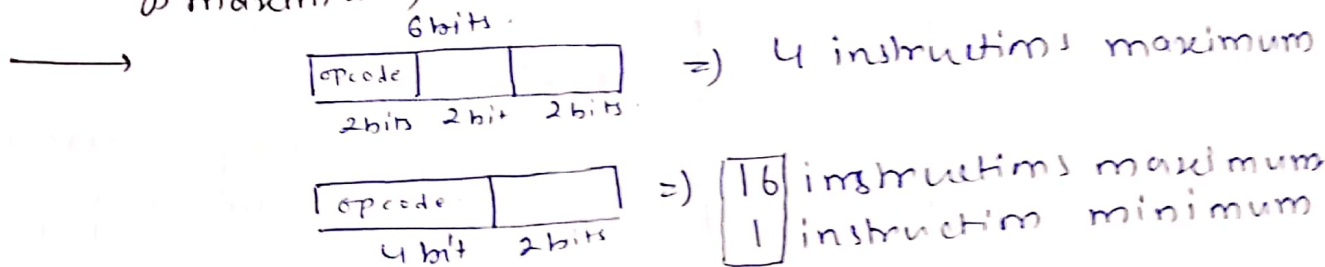
Instruction = 5 fields



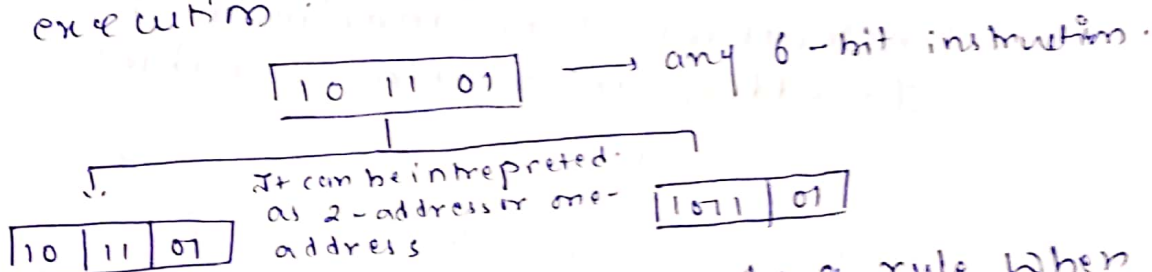
= 34 bits
= 5 bytes for an instruction

5 bytes for 1 instruction
500 bytes for 100 instructions

Question a digital computer has 6 bit instruction and 2 bit addresses, it supports 1 address & 2 address instructions both. If there are 3 : 2 address instructions then how many '1' address instructions can be formulated (minimum & maximum)



Considering an instruction comes to CPU for execution:



hence we have to apply a rule when the interpreter will go for one-address instruction or two address instruction. hence we have to apply some restriction on lower bit opcode.

total no. of instructions possible for (2 address) = $2^{\text{bit}} = 2^2 = 4$
 (00, 01, 10, 11)

total no. of used instructions = $\frac{3}{(00, 01, 10)}$

total unused instruction = $\frac{1}{(11)}$
 ↓
 assumed

Hence compiler will check for the first two bits and will go for one address instructions when it encounters '11'

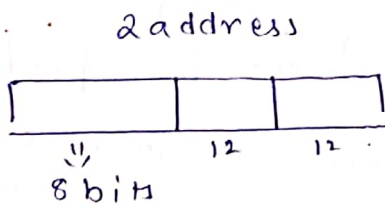
Hence possible 1 address instruction =

maximum
 Hence possible 1 address instructions = 4
 minimum = 1

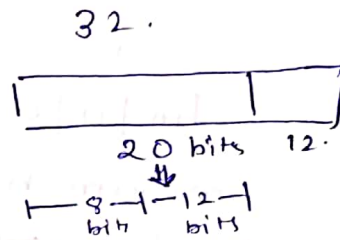
1100
1101
1110
1111
4

2-address instructions	Unused OpCodes	1-address instructions
4	0	$0 * 2^2 = 0 \Rightarrow$ only 2-address instructions supported.
3	1	$1 * 2^2 = 4$
2	2	$2 * 2^2 = 8$
1	3	$3 * 2^2 = 12$
0	4	$4 * 2^2 = 16 \Rightarrow$ only 1-address instruction supported.

question a computer has 32 bit instructions & 12 bit addresses if there are 254 2-address instructions then maximum how many 1-address instructions possible.



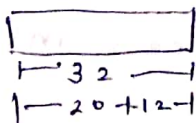
$$\begin{aligned} \text{max opcodes} &= 2^8 = 256 \\ \text{used opcodes} &= 254 \\ \hline \text{unused opcodes} &= 2 \end{aligned}$$



$$\begin{aligned} \text{max} &= 2 * 2^{12} = 2^{13} \\ &= \underline{\underline{8192}} \end{aligned}$$

If ~~3 types of~~ there are 8000 1 address instructions are used then maximum how many 0-address are possible.

$$\begin{aligned} \text{Max 1-address possible} &= 8192 \\ \text{used} &= 8000 \\ \hline &= 192 \end{aligned}$$

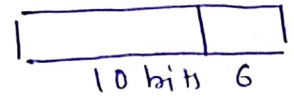
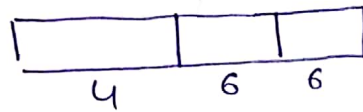


$$\Rightarrow 192 * 2^{12} = \text{maximum possible 0-address instructions}$$

question 16
P 44.

length of instructions = 16 bit.
address length = 6 bit.

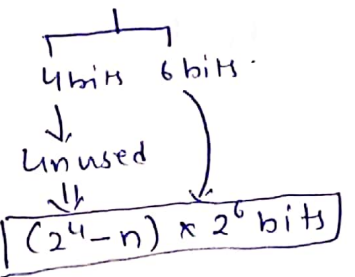
Supports 2 address & 1 address



max opcodes = 2^4 .

used opcodes = n

$$\frac{\text{Unused} = 2^4 - n}{}$$



gate

A processor has 16 integer register (R_0, R_1, \dots, R_{15}) and 64 floating point registers (F_0, F_1, \dots, F_{63}) it uses a two byte instruction format. There are 4 categories of instruction: type 1, type 2, type 3, & type 4. Type 1 category consist of 4 instructions each with three integer register operands (3Rs). Type 2 category consist of 8 instructions each with 2 floating point register operands (2Fs). Type 3 category consists of 14 instructions each with 1 integer register operand & 1 floating point register operand (1R + 1F). Type 4 category consist of 'N' no. of instructions, each with a floating point register operand. The maximum value of N is?

⇒ Type 1



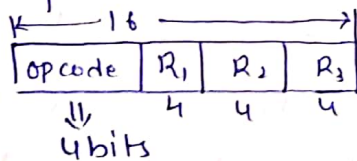
$(R_0, R_1, \dots, R_{15}) \Rightarrow \text{Integer} = 4 \text{ bits}$

$(F_0, F_1, \dots, F_{63}) \Rightarrow F = 6 \text{ bits}$

2 bytes instruction

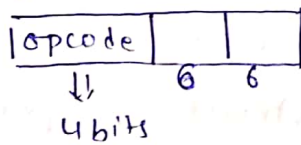
Type 1	Type 2	Type 3	Type 4
4	8	14	N
3R	2F	(1R+1F)	(1F)

Type 1



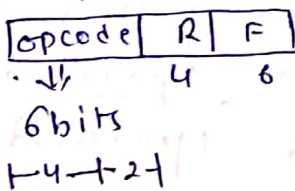
$$\begin{aligned} \text{max} &= 2^4 = 16 \\ \text{Used} &= 4 \\ \hline \text{Unused} &= 12 \end{aligned}$$

Type 2



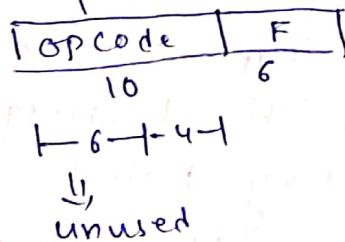
$$\begin{aligned} \text{max} &= 12 \times 2^0 = 12 \\ \text{Used} &= 8 \\ \hline \text{Unused} &= 4 \end{aligned}$$

Type 3



$$\begin{aligned} \text{max} &= 4 * 2^2 = 16 \\ \text{Used} &= 14 \\ \hline \text{Unused} &= 2 \end{aligned}$$

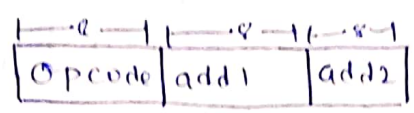
Type 4



$$2 * 2^4 = \underline{32} = N$$

question a computer supports 3 byte instructions and 8 bit addresses there are maximum 1024 '1' address instructions. then how many 2 address instructions are supported by the system?

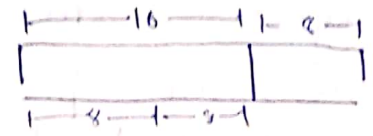
→ Assume n 2-add instruction supported



max = $2^8 = 256$

used = n

unused = $256 - n$

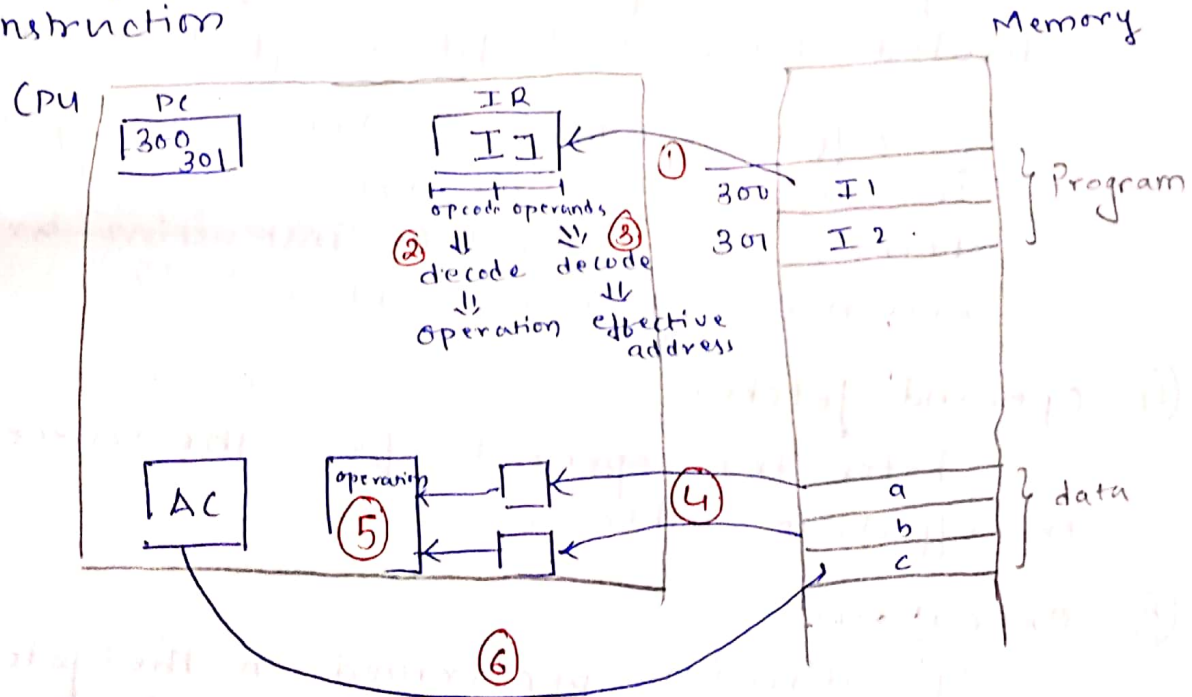


$(256 - n) \times 2^4 = 1024$

$n = 252$

* Instruction cycle :-

System performs 6 phases to execute an instruction



1. Instruction fetch
2. Instruction decode
3. Effective address calculation
4. Operand fetch

① Instruction fetch

Using 'PC' value next instruction is copied from memory to 'IR' Instruction register. during Instruction fetch only the value of 'PC' is incremented by the size of instruction.

② Instruction decode.

opcode part of instruction is decoded by CPU to obtain operation.

③ Effective address calculation

Operand Information part of instruction is decoded by CPU to obtain effective address

=> effective address :- Address of operand in computation type instruction & target address in branch type instruction are known as effective address.

④ Operand fetch.

Fetch the operands from the source using the effective addresses.

⑤ execution:

Operation is performed on the fetched operand and result is generated.

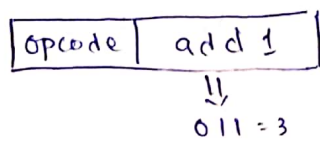
⑥ Write back

The generated result is copied back to the destination

Note: all type of instructions do not require all '6' phases

⇒ Addressing modes:-

It specifies how and from where the operand is obtained using the address field of instruction. They are used in 3rd step of instruction cycle i.e. effective address calculation phase.



add 1 could be anything.
#3 (no. three)

R3 (Register no. 3)

M[3] (Memory address 3)

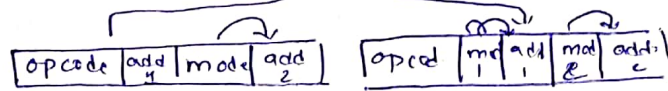
therefore to satisfy ambiguities the 'mode' is introduced.



00 ⇒ operand

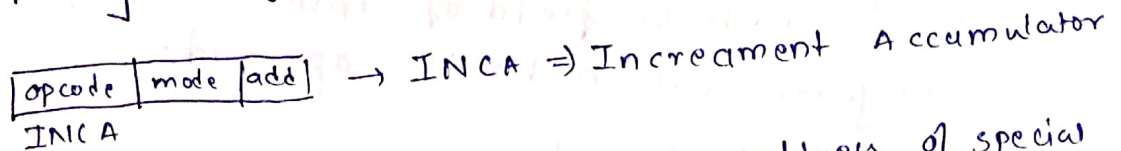
01 ⇒ Reg.

10 ⇒ Mem. address



→ Types of addressing modes

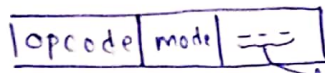
→ Implied mode: The opcode definition itself specifies the operand hence no need to specify any explicit address for the operand.



There is no need to specify the address of special registers explicitly hence the operation's opcode will not specify it

→ Immediate mode: The address part of instruction specifies operand value

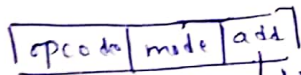
This mode is used to initialize register with constant value.



directly value is mentioned as operand instead of address.

→ Direct mode: (absolute mode):

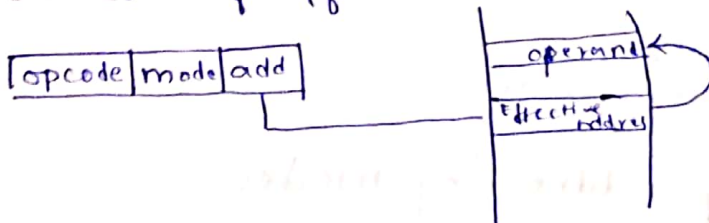
The address field of instruction specifies the effective address



memory address of operand (Effective Address)

→ Indirect mode:

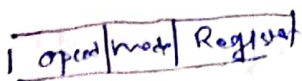
The address field of instruction specifies address of effective address



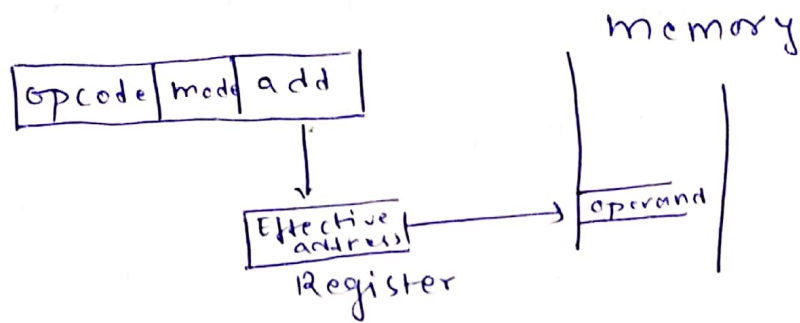
This mode is used to implement pointers

→ Register mode:

address field of instructions specifies a register which holds operand.



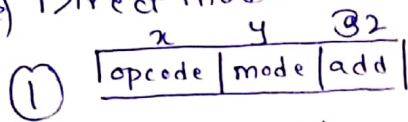
→ Register indirect mode:-
Address field of instruction specifies a register which holds effective address



This mode is used to shorten the instruction length.

consider a system with 4 GB Ram & 64 Regs
Mem add = 32-bits
Reg. reference = 6-bits

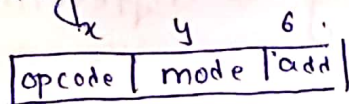
1) Direct mode



Length of instruction = $x + y + 32$

time to get operand
① = 1 mem. access time.

2) Register indirect



Length of instruction = $x + y + 6$

② = 1 Reg access time + 1 mem access time
 \approx 1 mem access time

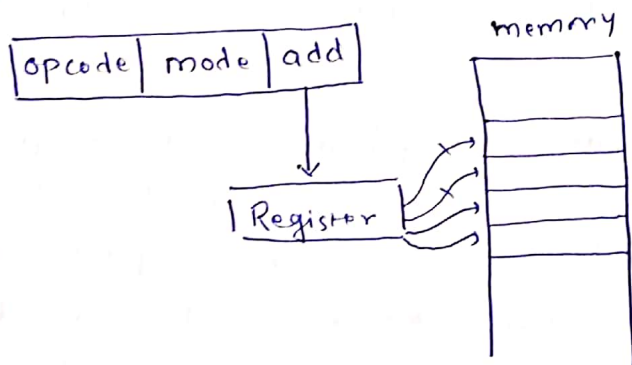
So within same amount of time Register Indirect mode gives better performance by reducing the length of instruction.

Auto/

→ Auto increment / Auto decrement mode:-

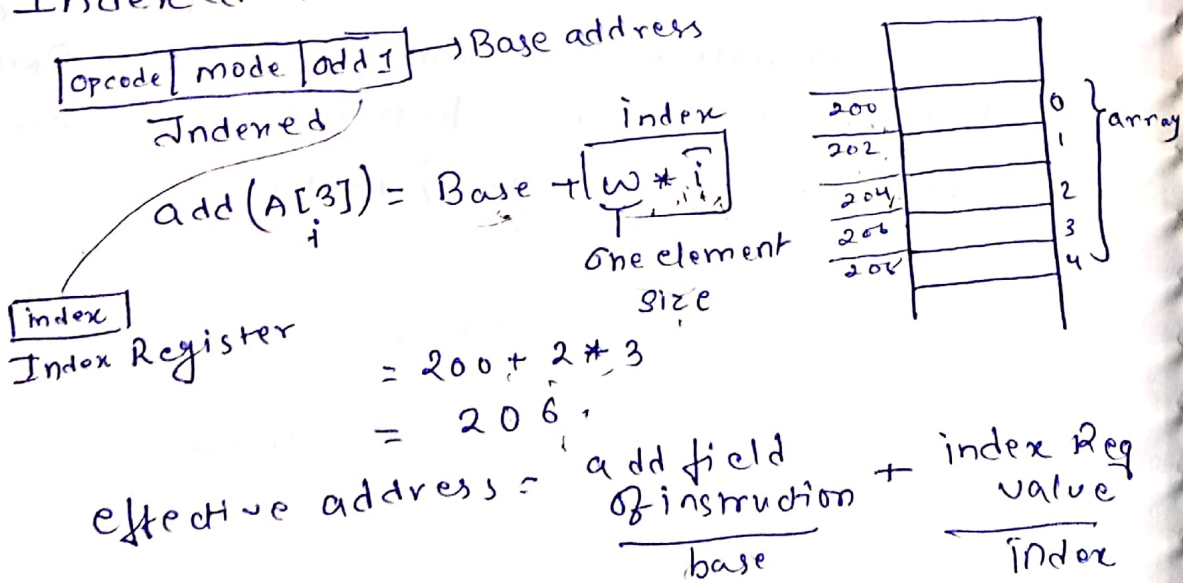
It is a variant of Register indirect mode in which the content of Register (Effective address) is automatically incremented or decremented. This mode is use to access a table of content (array) sequentially

Autoincrement → post increment
 Auto decrement → pre decrement



Here same instruction is used 'n' no. of times into a loop instead of creating n different instructions and string it.

→ Indexed mode or Indexed Reg. mode:-

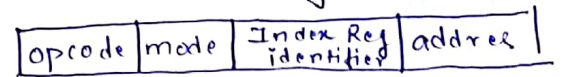
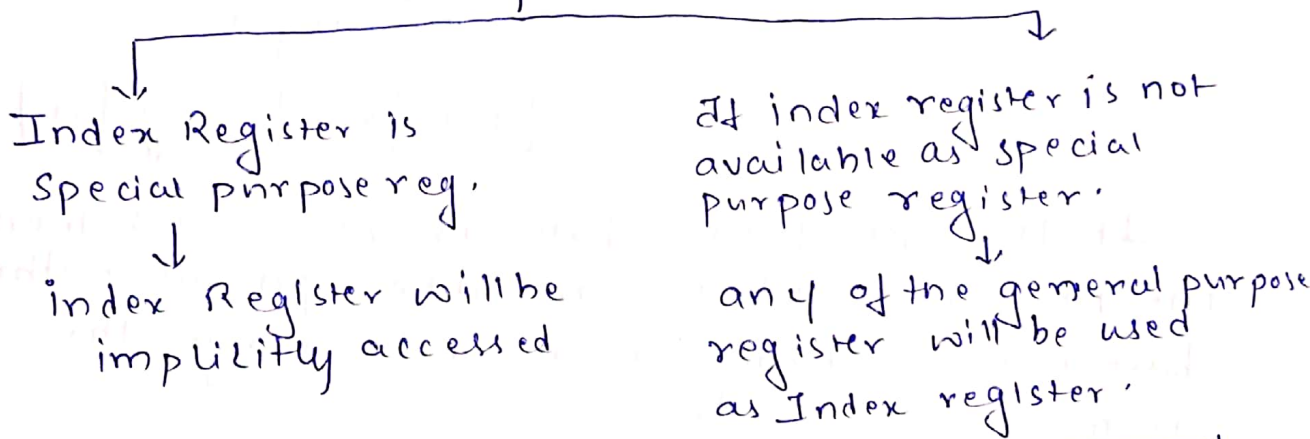


Address of element z of array = Base + index

Instruction to calculate index is used before the instruction of access array element to calculate the index value for further usage

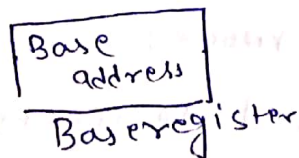
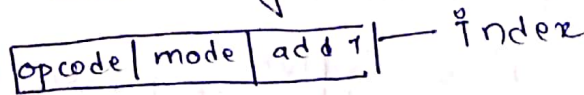
Index register is a special register which may or may not be available in CPU of different type.

Implementation



If the program data (array) is relocated then the new base address should be updated in the address part of the instruction but updation of instruction is a costly operation hence to solve this we have another mode

→ Base register mode:



$$\text{effective address} = \frac{\text{add. field of instruction}}{\text{Index}} + \frac{\text{Base Reg Value}}{\text{Base add.}}$$

Here the index is kept inside the instruction and the base address is stored in base register hence on relocation only value of register is updated that reduces the complexity of relocation

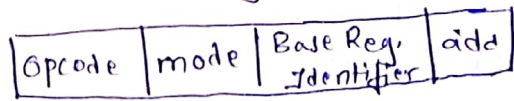
Implementation

Base register as special purpose register.

Base Register is implicitly accessed.

Base register is not available as special purpose register.

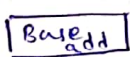
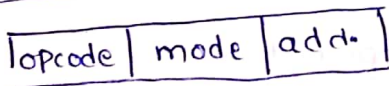
any of the general purpose register will be used as Base register



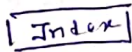
If the program data is relocated then the new base address should be updated in the base register hence no need to change the instruction or code.

Instruction to calculate the index value (with i) will be runned in compile time.

→ Base + index register mode:-



Base Register



Index Reg.

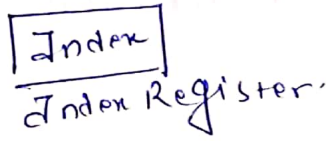
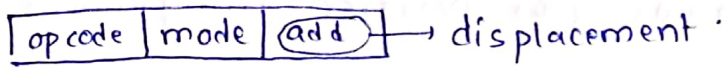
Base address and Index, both are stored into the registers. and effective address is passed to the instruction

$$\text{Effective add} = \text{Base Register Value} + \text{Index Reg. Value}$$

→ Base + Index + displacement mode:-

This mode is used to access array of records (structures)

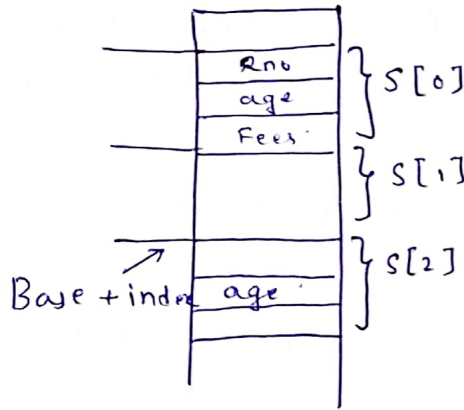
displacement field will be used to enter into specific element of structure.



$$\text{Effective add} = \text{Base Reg. value} + \text{index Reg. value} + \text{add. field value of instr. displacement}$$

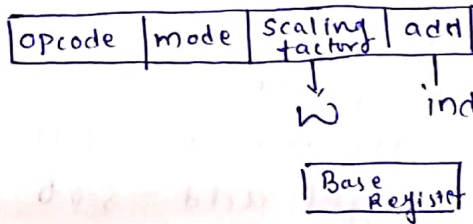
Structure

- Rno
- age
- fee.



Base + index will give address of S[2] but to access inside it displacement is going to use.

→ Scaled mode:



Here the instruction to calculate the index is not excluded externally but 'scaling factor' helps to get the value of 'w' and no. of instructions in Base or index mode is reduced.

This can be used with both Base & index register mode.

$$\text{Effective add} = \left(\begin{matrix} \text{Base} \\ \text{Reg} \\ \text{value} \\ \text{Base} \end{matrix} \right) + \left(\begin{matrix} \text{Scaling} \\ \text{factor} \\ w \end{matrix} \right) * \left(\begin{matrix} \text{add field} \\ \text{value of} \\ \text{inst?} \\ i \end{matrix} \right)$$

PC Relative mode (or position independent mode)

This mode is used for branch type of instructions

↳ Branch condition :-

Assume that CPU is currently executing instruction I₂.

PC = 502.

I₂ is a branch instruction

memory	
500	I ₁
501	I ₂
502	I ₃
503	I ₄
504	I ₅
505	I ₆
506	I ₇
	!

} Program

condition = false
↓
Branch is not taken

↓
Next instruction ⇒ next in the sequence.
⇒ I₃

No change in PC explicitly

condition = true
↓
Branch is taken

↓
Next instruction ⇒ Target instruction. (Assume it is I₇)
⇒ Target add = 506

PC = Target add.

Target add = PC + Relative location to skip.

So to calculate the target address we use.
(PC-value will give the value to jump in branch)

op code	mode	address
---------	------	---------

PC Relative

↳ Relative location to skip

Effective add = PC value + add field value of instruction

When the branch instruction is encountered the pc-relative mode will be activated only when the condition is true. and condition is checked in execution phase

When the program is relocated the address part of instruction will not be affected as it will contain only the no. of skips. Hence it is also called (position independent instruction)

*Relative no. of skips is also known as 'offset'
 → Forward branching ⇒ offset ⇒ +ve value
 → Backward branching ⇒ offset ⇒ -ve value

Example.

PC = 200
 202

Reg = 400
 399

XR = 150

XR = Index Reg.

200	opcode	mode
201	Address = 500	
202	Next instruction	
	⋮	
399	850	
400	760	
	⋮	
500	800	
	⋮	
600	900	
	⋮	
702	Target instruction	
	⋮	
800	350	

memory.
 ←

PTO

mode.	Effective add.	Operand.
1) Immediate mode	201	500
2) Direct mode	500	800
3) Indirect mode	800	350
4) Reg. mode.	Register address is not available.	400
5) Reg. Indirect mode.	400	700
6) Auto decrement	399	850
7) Index Reg mode	$500 + \frac{100}{\times R} = 600$	900
8) PC - Relative	$202 + 500 = 700$	

Register address is never taken.

For PC - Relative mode :-

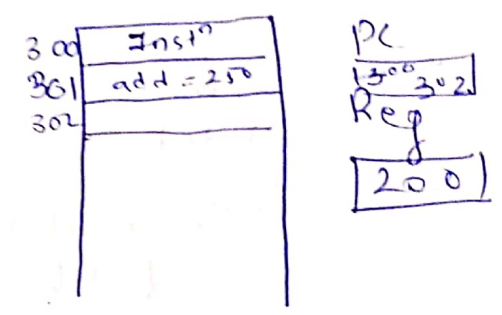
~~in~~ The target address is calculated in effective address phase.

and in execution phase the branch is checked for the condition and if condition is false 'PC' is given next instruction address but if condition is true then the 'PC' will be given address i.e. target address

question an instruction is stored at location 300 with its address field at location 301. the address field has the value 250. A processor register contains the no. 200. Evaluate the effective address if the addressing mode is

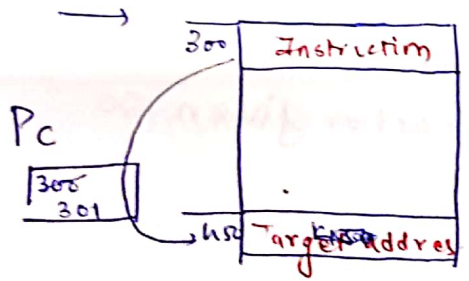
- 1) direct mode.
- 2) Immediate mode
- 3) PC - relative
- 4) Register indirect

mode	Effective add.	operand.
→ direct mode	250	
→ Immediate	301	
→ Pc-relative	200	
→ Reg-indirect	$302 + 250 = 552$	



Question a relative branch mode type instruction is stored in memory at address 300. The branch is made to an address 450.

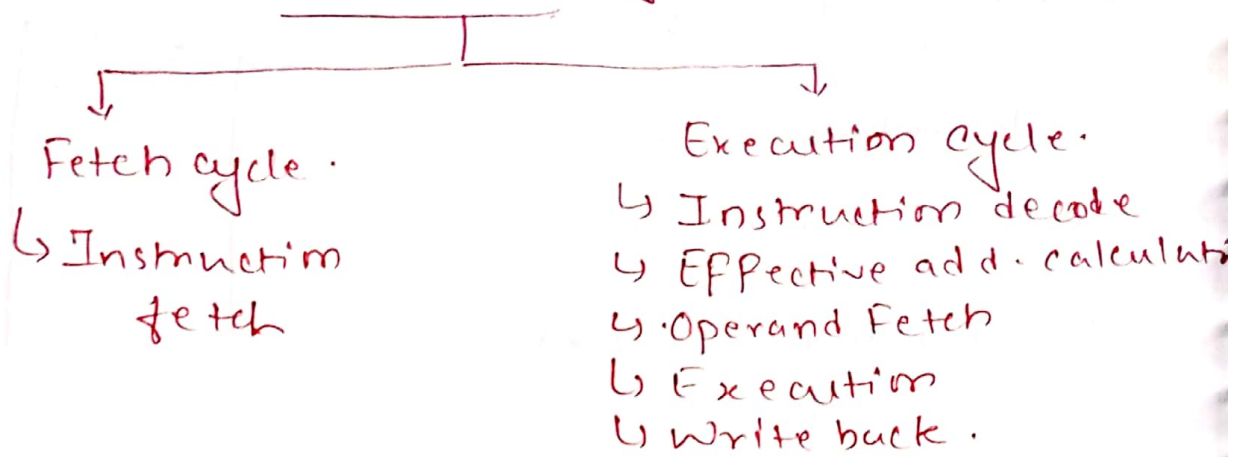
- (1) What should be the value of relative address field of the instruction
- 2) Determine the value of program counter before instruction fetch, After instruction fetch & after execution of this instruction



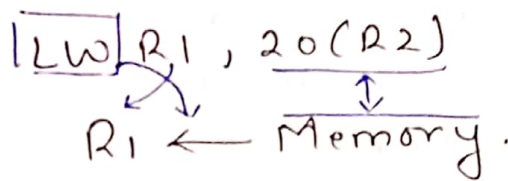
1) Relative address field.
 Target add = PC + Relative add.
 $450 = 301 + \text{Relative add.}$
 Relative address = 149

- 2) Before instruction fetch
 $PC = 300$
- After instruction fetch
 $PC = 301$
- After Execution
 $PC = 450$

Instruction cycle.



Page 44.
Question 14



LW \Rightarrow operation. (Copy memory to $R1$)

Effective add = $20 + R2$.

LW \rightarrow Load Word.

$R2 \rightarrow$ General purpose register.

hence its not PC-relative.

\rightarrow There is no scaling factor given so not scaling mode.

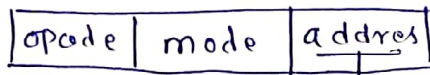
Hence 'D' is correct

- Q-7
- Rom is Read only memory
 - PC points to next instruction
 - ✓ Stack is LIFO
 - not all instruction affects the flag.

Q8
Pg 43

Implementation of Branch instruction

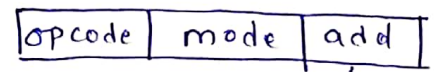
Position dependent



Target addr.

This will give problem when there is a need to relocate program.

Position independent



relative location to skip

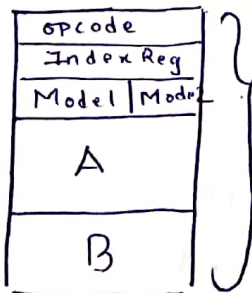
=> Relative mode is also known as Position independent mode.

Q9 In Pc relative mode relocation of Instruction takes place without any change in code
 -> Base register mode relocation of Data takes place.
 Hence Both are true here.

Q10

instruction => 3 words,

ADD A[R0] @B



3 word Instruction

Execution cycle does not include the instruction fetch.

-> A & R0 are already in the CPU

So in A[R0] '1' memory cycle is accessed

@B is indirect mode then it will take

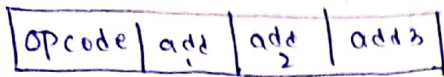
'2' memory access

Hence '3' memory access for address effective address calculation and operation fetch & '1' memory access for write back.

Total

$$1 + 3 = 4$$

Q12

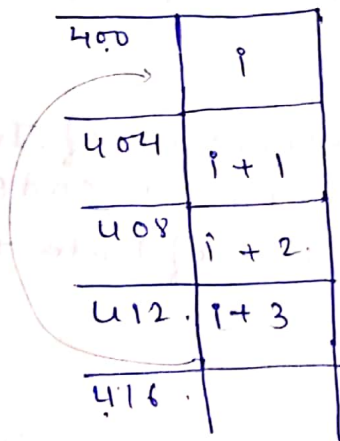


An implied accumulator register is a special register i.e we can use it explicitly through its address.

Q19

Instruction size = 4B

Offset = relative location. are stored in Bytes out of all the instruction the instrⁿ with offset is branch instruction.



CPU executing i+3.
hence PC = 416.
Target add = 400
PC + offset = 400
416 + offset = 400
offset = -16

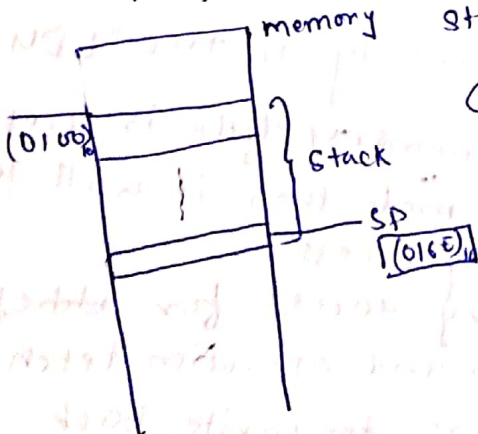
IF the target of the branch instruction is 'i' that means target address = 400

~~Q20~~

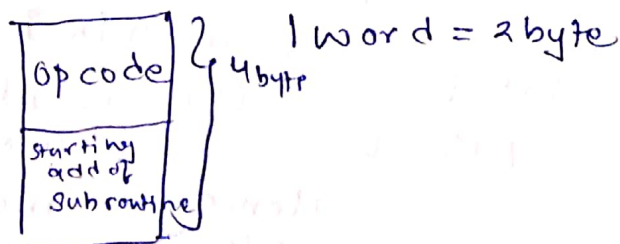
~~Q21~~

Q13

Byte addressable memory
PC, PSW. all register are of 2 Bytes
stack grows upward.



CALL instruction → 2 words

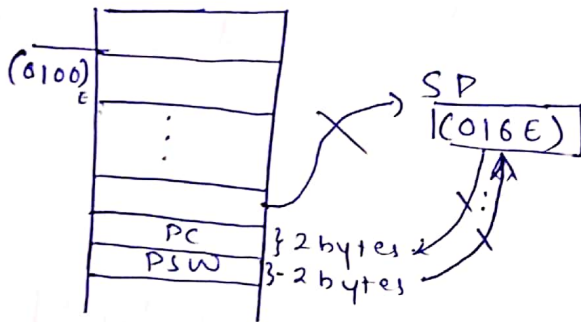


CALL

- ↳ store current PC value in stack
- ↳ " PSW value in stack
- ↳ Load starting address of subroutine in PC

PC = $(5FA0)_{16}$ Before fetch of call instruction.

Hence $(5FA0)_{16}$ is address of call instruction



$016E$
+ 2 After PC

 0170 in stack
+ 2

 0172 After PSW
 in stack.

SP points on the top of stack and call function pushes the two value in stack.

question:- In the above question what will be the value of PC before fetch of call instruction, After fetch of call instruction, and after the execution phase of call instruction.

Before fetch $\Rightarrow (5FA0)_{16}$.

After Fetch $\Rightarrow 5FA0 + 4 \text{ bytes} = 5FA4$

After Execution \Rightarrow Starting address of subroutine. one instruction is 4 byte long.

Branching

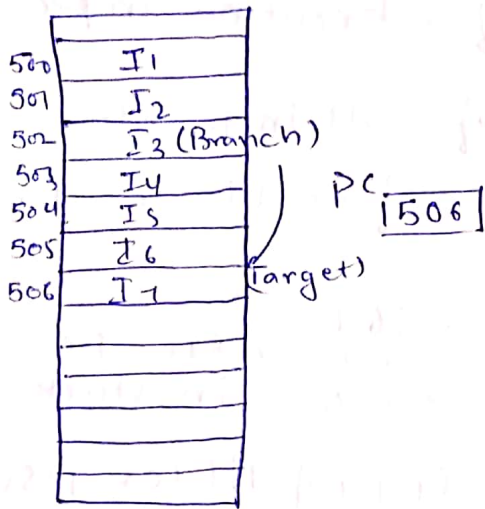
Jump

→ PC value updated by Target

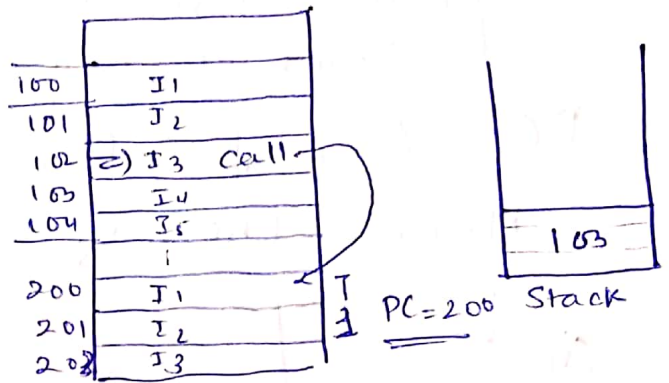
Function call

→ PC value + other regs are stored on stack
→ PC = target (starting add of function)

Jump



Function call



Page-43
Que 11

$R_1 \rightarrow M[100]$

$M[100] \rightarrow R_2$

$M[100] \rightarrow R_3$

\Rightarrow

$R_1 \rightarrow M[100]$

$R_1 \rightarrow R_2$

$R_1 \rightarrow R_3$

⇒ CPU Design ⇐

CPU cycle time:- The time in which CPU can perform one micro-operation.

CPI:- (Cycles per instruction): No. of cycles required to execute an instruction

$$\Rightarrow \text{Execution time} = \text{no. of instruction} * \text{CPI}_{\text{avg}} * \text{CPU cycle time}$$

$$\Rightarrow \text{clock rate of CPU} = 1 / \text{cycle time.}$$

$$\Rightarrow \text{Execution time} = \frac{\text{no. of instruction} * \text{CPI}_{\text{avg}}}{\text{clock rate.}}$$

Instruction Type	No. of Instruction (n_i)	CPI	CPI
Load & store	30	6	180
ALU	40	5	200
Branch	10	4	40
Other	20	5	100
	<u>100</u>		<u>520</u>

$$\text{CPI}_{\text{avg}} = \frac{520}{100} = \underline{\underline{5.2}}$$

$$\text{CPI}_{\text{avg}} = \frac{\sum_{i=1}^n \text{CPI}_i * n_i}{\sum_{i=1}^n n_i}$$

$\text{CPI}_i \rightarrow$ CPI for instruction of type i

$n_i \rightarrow$ no. of instructions of type i

Performance of CPU is given by MIPS
(million instructions per second)

$$\text{MIPS} = \frac{\text{no. of instructions}}{\text{Execution time} \times 10^6}$$

$$\text{MIPS} = \frac{\text{Clock rate}}{\text{CPI}_{\text{Avg}} \times 10^6}$$

question
gate.

Consider two processors P_1 & P_2 executing the same instruction set. Assume that under identical conditions, for the same input, a program running on P_2 takes 25% less time but incurs 20% more CPI more as compared to program running on P_1 . If the clock frequency of P_1 is 1 GHz then clock frequency of P_2 (in GHz) is?

	P_1	P_2	No. of instructions $n_1 = n_2$.
Execution Time	t_1	$t_2 = 0.75 t_1$	
CPI	C_1	$C_2 = 1.2 C_1$	
Freq.	1 GHz	?	

$$\text{Execution time} = \frac{\text{No. of inst.} \times \text{CPI}}{\text{Frequency}}$$

$$\text{No. of inst.} = \frac{\text{Ex time} \times \text{freq}}{\text{CPI}}$$

$$\frac{t_1 \times F_1}{C_1} = \frac{t_2 \times F_2}{C_2} \Rightarrow \frac{t_1 \times 1}{C_1} = \frac{0.75 t_1 \times F_1}{1.2 C_1}$$

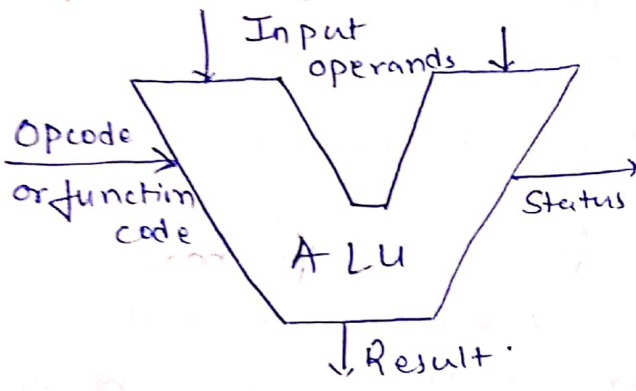
$$\Rightarrow F_1 = 1.0 \text{ GHz.}$$

Note:

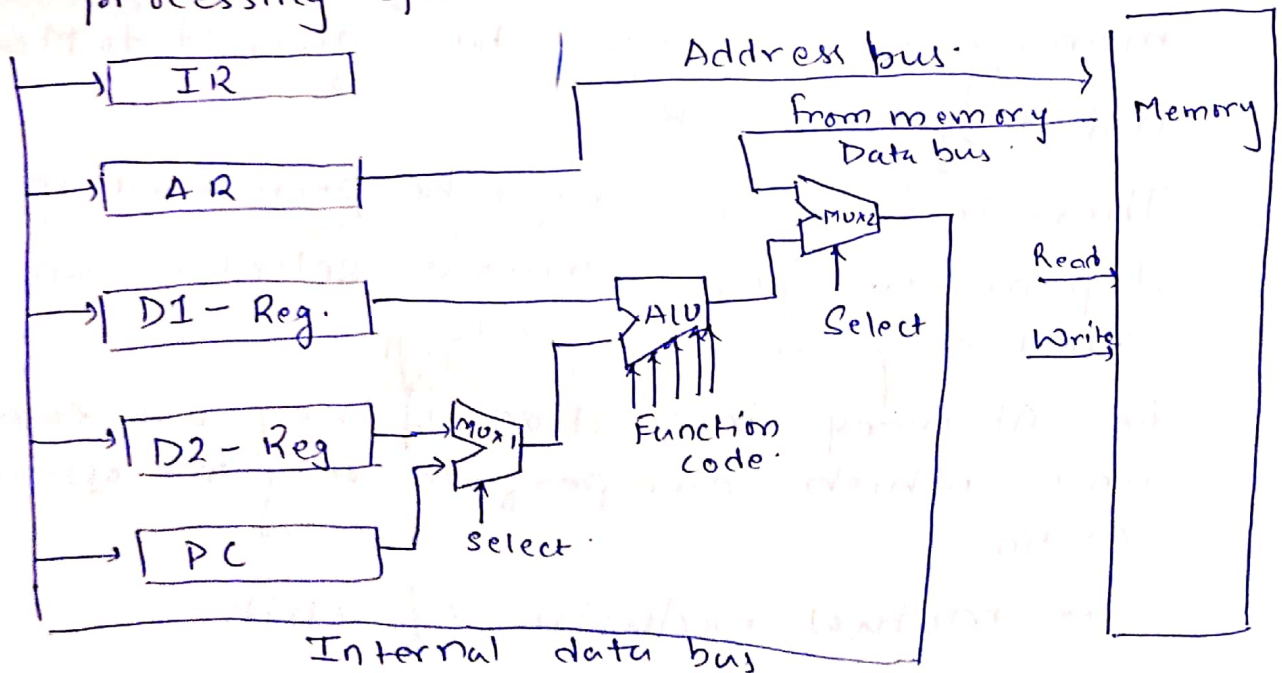
MIPS count does not provide always the correct performance of the CPU. because MIPS does not consider the complexity of instruction.

* FLOPS — Floating point operations per second. It is another way to measure performance of CPU i.e performance measure metric.

⇒ ALU (Arithmetic Logic unit)



⇒ Data path: collection of functional units such as ALU, multiplexer etc to perform data processing operations.



Instruction Fetch: $IR \leftarrow M[PC]$

But this instruction cannot be executed directly there are lot of processes included to execute this instruction.

This instruction will be divided as:

$AR \leftarrow PC$

$IR \leftarrow M[AR], PC \leftarrow PC + 1$

i.e. First PC value is supposed to be added to the AR and the $M[AR]$ to IR

this happen in many steps further.

Firstly PC value will be forwarded to MUX1 with select = 1 and then it moves further to ALU, ALU uses some function code to move $A + B$ further such that $A * B \rightarrow B$ and then it moves further to MUX2 with select 1 which further passes to Internal data bus and through that data bus it gets copied into AR and further AR moves address through address bus to memory and memory bus gives it to MUX2 and it gets to 'IR'

These instructions cannot be performed in 1 step as everytime mux's selector can select only one output to flow.

i.e. at every instruction if they use common unit which can perform only one operation at time.

i.e. mutual exclusion of units

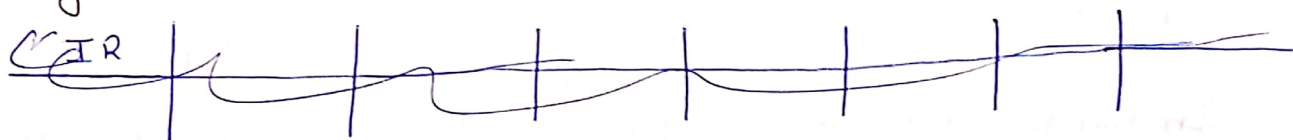
CPU can perform only those instructions at a time which do not require common unit. Control unit will generate signals.

* Control unit:

It generates control signals and sends those signals to the units of the system and units perform their respective operation.

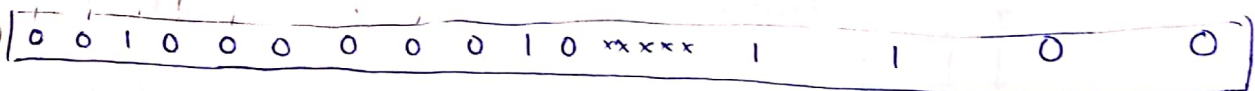
→ Control variable: The appropriate name of a control signal.

→ Control word: Set of all control signals generated by control unit at a time.

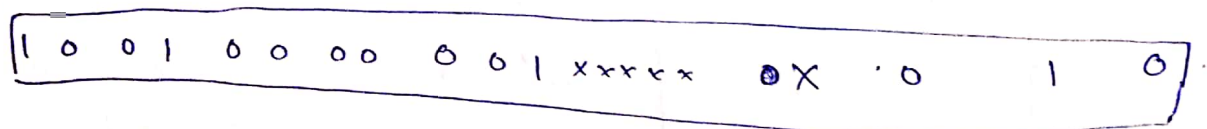


IR		AR		D1		D2		PC			Mux1		Mux2		memory	
In	out	In	out	In	out	In	out	In	out	Inc			Select	Select	Read	write

→ Control word for (AR ← PC)



Control word for (IR ← M[AR], PC ← PC + 1)



→ Control unit organization:-

→ Q. Two types of control unit organization.

- 1) Hardwired control unit.
- 2) Microprogrammed control unit.

Hardwired control unit: The control logic is implemented using logic gates, decoder & other digital circuits (Hardwired CU.)

Advantages:

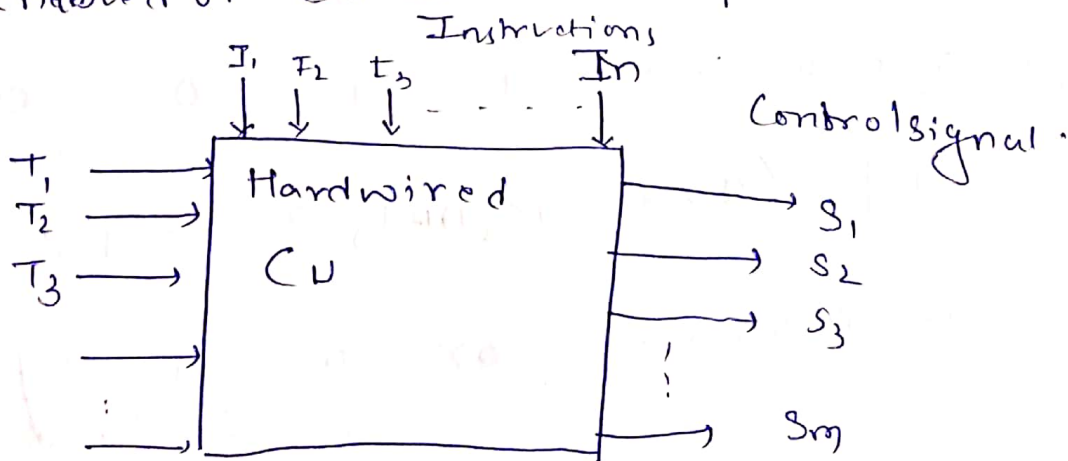
→ Faster mode of operation.

Disadvantage:

→ updation in control logic is difficult because rearranging wire among circuits is difficult.

In hardwired the logical gates are implemented in some circuit called control unit.

the instructions will come and for every instruction certain signals are supposed to be enabled or disabled at a particular set of time



So every instruction cycle needs different set of signals this set of signals will lead to microoperations. For every instruction the time will start again from t_1 .

	I_1	I_2	I_3	I_n
t_1	S_1, S_5, S_9			
t_2	S_2, S_3, S_4, S_{10}			
t_3			
\vdots				
t_n				

at one ~~epoch~~ microoperation time set of signals will generated. and for all instruction cycle time sequence and set of signals are formed. i.e. for every instruction again the time sequence is formed.

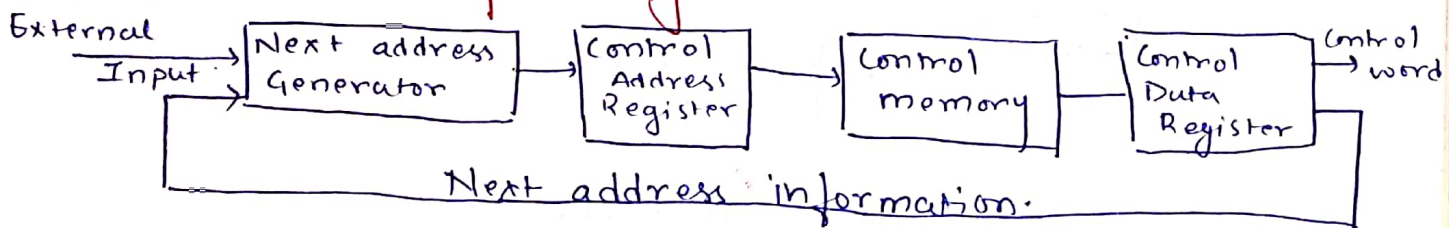
Microprogrammed control unit :-

all the possible control words are stored in a memory. Based on the requirement, the appropriate control word is fetched and signals are sent to corresponding units.

Advantage: updation in control logic is easy

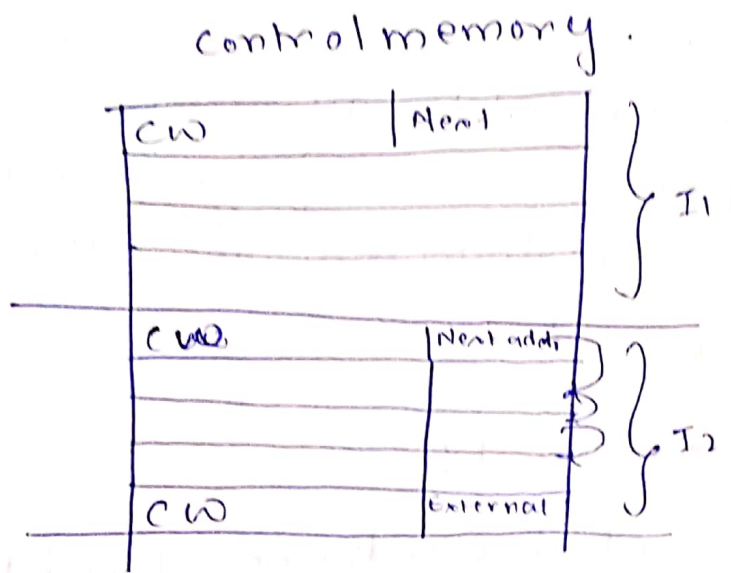
Disadvantage: It is slower than hardwired CU.

Control word sequencing :-

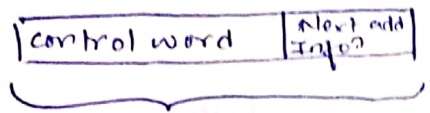


all the control word which are responsible to execute one type of instruction, are stored in control memory sequentially.

Inside control memory control word block is attached to another block which alerts whether to take the next address or to take stop the instruction cycle and wait for external input.

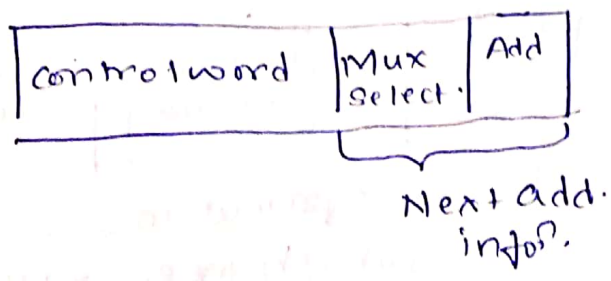


at every address in control memory the following information is stored



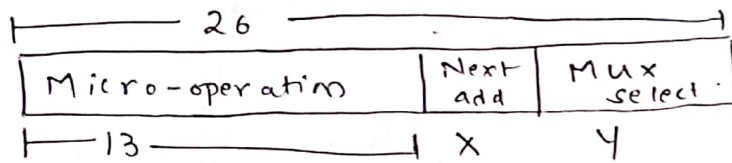
This is known as microinstructions.

Format of microinstruction is



Q4
P57

micro-instruction = 26 bit



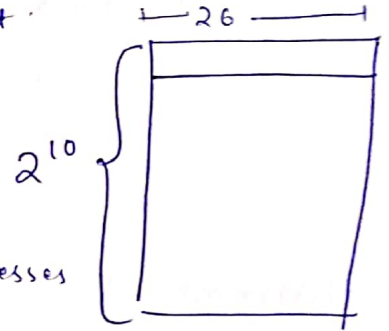
Mux needs 8 status bits that means
Select requires 3 bits input.

i.e. $y = 3$

$$13 + x + 3 = 26.$$

$\therefore x = 10$ i.e. 2^{10} addresses
can be stored

Now size of memory = $2^{10} \times 26$ bits
= 26 k bits.



1 word = 26 bits.

\therefore size of memory in words = 1024 words.

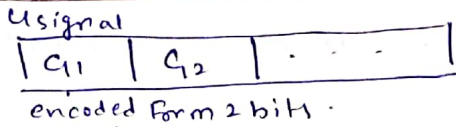
Types of microprogrammed CU.

↳ Horizontal.

→ 1 Bit for each control
signal in a control word
Hence larger sized
control word.

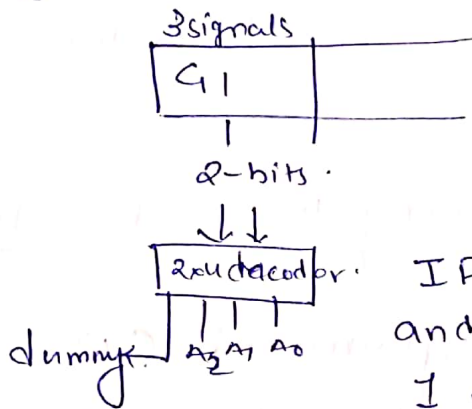
↳ Vertical.

The control signals are
divided into multiple
groups in such a way
that within each group
only one signal is enabled
at a time.
Each group information is
stored in encoded form
to decrease the size of
control word. The control
signals are generated
using the decoders.



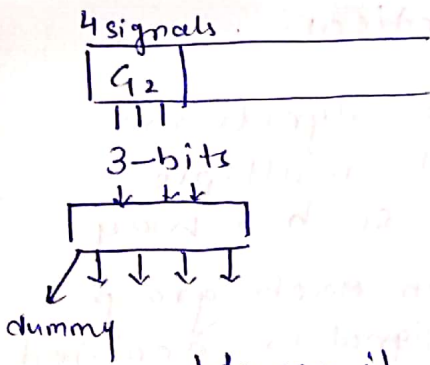
Horizontal is faster than vertical.
 among all these CUs.
 Hardwired is faster than Horizontal is
 faster than vertical.

→ atmost 1 signal enabled at a time within each group.



IF 3 signals are passed and atmost we require 1 at a time.

then encoder will need 2 bits and decoder will generate one extra signal for dummy.



Hence if 4 signals are to be generated at a time the encoder will require

$\lceil \log_2 5 \rceil$ bits i.e 3 bits to store

4 for 4 signals and 1 for dummy dummy says when to keep all the signal disabled

If any signal cannot be grouped then it is stored as horizontal way.

Qn. A microprogrammed controlled unit is required to generate a total of 25 control signals assume that during any microInstruction, atmost two control signals are active minimum no. of bits required in the control word to generate the required control signals will be ?

There are 25 signals. Without mutually exclusion 2 groups will be generated as to make atmost 2 pairs of signals to be active at a time.

25+1	25+1
------	------

$\therefore 5 \text{ bits} + 5 \text{ bits} = \underline{10 \text{ bits}}$

Qn. A CPU can support 16 distinct instructions. To execute each instruction a sequence of 64 micro-operations is required. the microprogrammed control unit stores each micro-instructions in three fields 1) Horizontal control word to generate 128 signals 2) Next address field 3) MUX select field to select 1 of 8 Mux inputs.

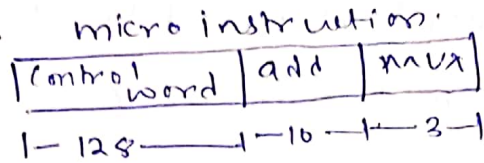
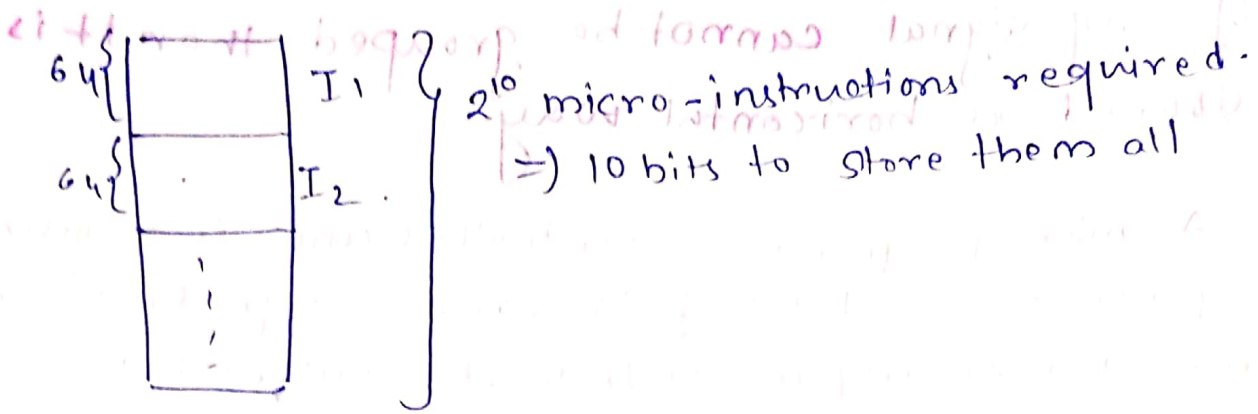
What is the size of control memory required for this control unit.

→ no. of inst? = 16

no. of micro-operation sequence per inst? = 64.

\therefore total microinstructions stored in memory = $16 * 64 = 1024 = 2^{10}$

\therefore control memory address size = 10 bit. to store 2^{10} instructions.

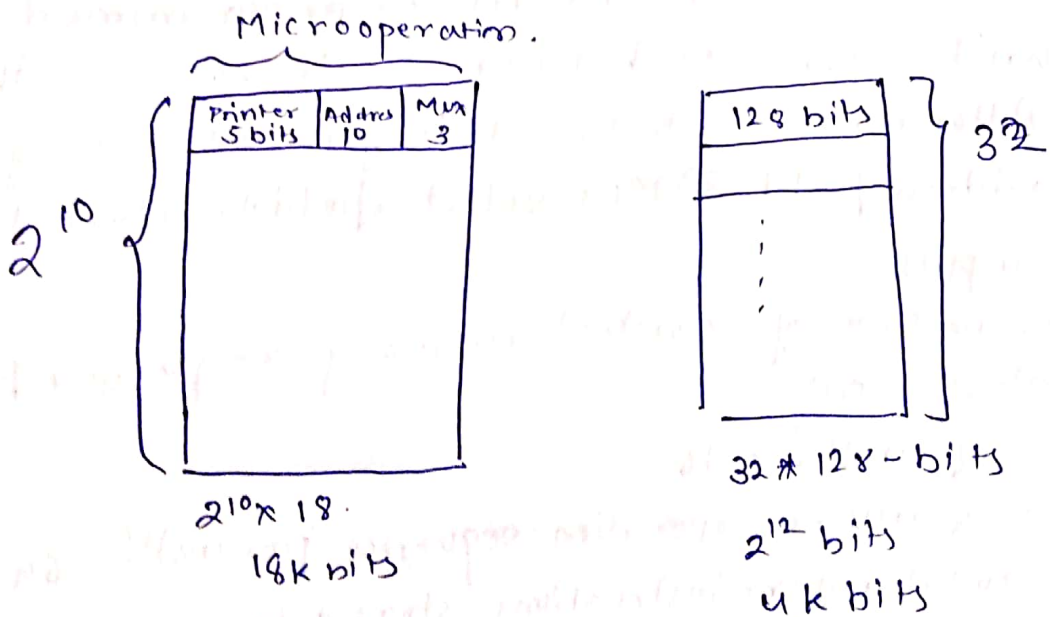


$$\begin{array}{r}
 = 128 \\
 + 10 \\
 + 3 \\
 \hline
 141 \text{ bits}
 \end{array}$$

\therefore Size of control memory = $2^{10} * 141$
 = 141 k bits

If we consider nano programmed control unit.
 then.

Unique no. of microoperations (control words) = 32.



22 k bits will only be required

RISC vs CISC

RISC

- Reduced Instⁿ. set
- Less no. of Instⁿ. supported.
- Simple instⁿ.
- Less no. of addressing modes
- Fixed length instⁿ.
- Hardwired control unit
- Reg-To-Reg arithmetic operations only.
Input output
- More no. of Registers

CISC.

- Complex instⁿ. set.
- more no. of instⁿ. supported.
- Complex instⁿ.
- more no. of addressing modes.
- variable length Instⁿ.
- microprogrammed CU also possible.
- Reg-to-memory & memory to-Reg also possible.
- Less no. of Registers.

Q2
[SI]

$q(2)$
 $P(51)$ | all are correct

$\frac{98}{P52}$ (d)

$\frac{9|3}{51}$ (d) Initiation of interrupt services.

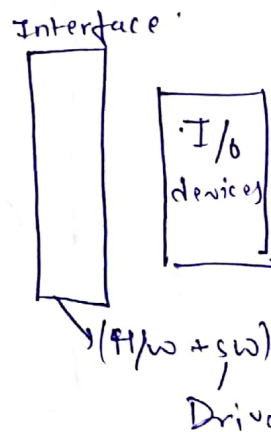
I/O organization

Peripheral device: a device which is connected to CPU externally apart from memory.

category of peripheral devices.

- ↳ Input device
 - ↳ Output device
 - ↳ Storage device
- (aka) I/O devices.

CPU



Need for interface:

- 1) Conversion of signals required.
- 2) Synchronization is required.
- 3) Data Format conversion is required.
- 4) All the i/o devices should be managed in such a way that each device should not disturb the operation of CPU & any other i/o device.

^{Primary} memory is not included in I/O devices because there is no require of any interface to communicate with CPU i.e. ~~why they are not incl~~ Peripheral devices

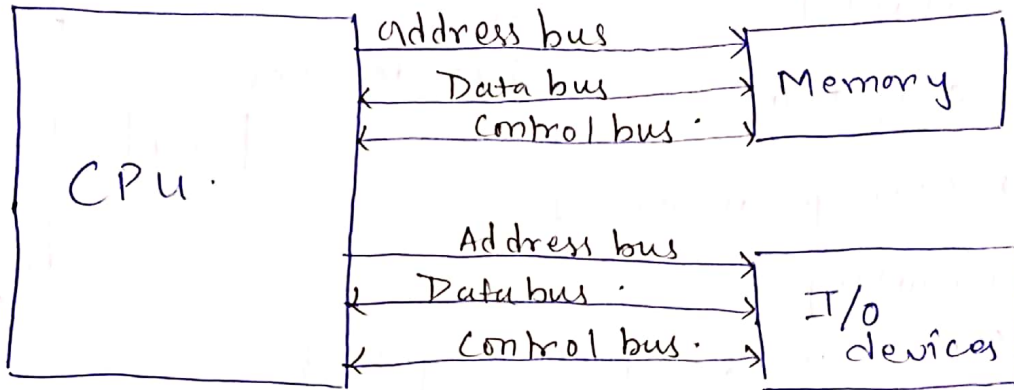
I/O processor = Interface + DMA Controller + I/O instruction execution.

There are three ways bet in which CPU can communicate with memory & I/O

* Communication of CPU with Memory & I/O

There are three ways in which CPU communicates with memory & I/O

1) Separate Buses For both memory & I/O



Address Bus in I/O is required to justify which I/O device is used at instance.

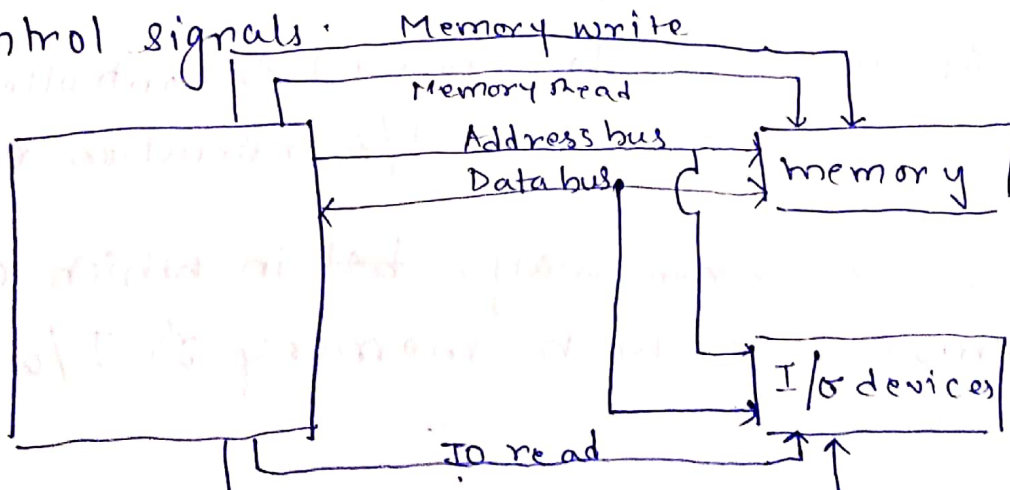
Advantages:

Easy to implement.

Disadvantage:

Costly.

2) Common address & Data buses with separate control signals.



Here control signals (Memory Read/write, I/O Read or Write) will suggest which address should be checked. and CPU will decide which operation to be done and control signal which which say what to access

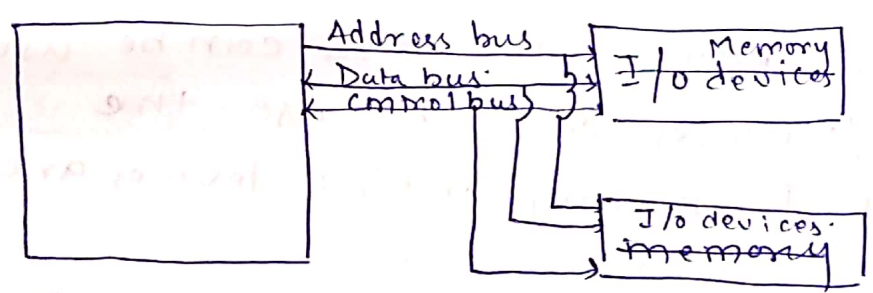
This type of communication is termed as (I/O mapped or (port mapped I/O)

Memory Read/write, I/O read or write is stored in a single group in vertical microprogramming i.e only one out of 4 will be active at time

Advantage:
No memory wastage

Disadvantage:
- Less no. of Instructions to access I/O devices.
- Less no. of Addressing modes to access I/O devices.

3) Common buses For both memory & I/O.



From all available addresses few are assigned to I/O devices, these addresses are stored in memory hence for every address which directs to I/O device memory is wasted.

eg: consider a system with 8B (Byte addressable) memory and 2 I/O devices d_1, d_2 .





assume

device	address
d1	010
d2	100

These two addresses are useless to memory which is memory wastage.

These address table will be stored in memory controller which decides the allocation of memory and assigns the control to memory or I/O device.

Port mapped I/O

→ Advantages:

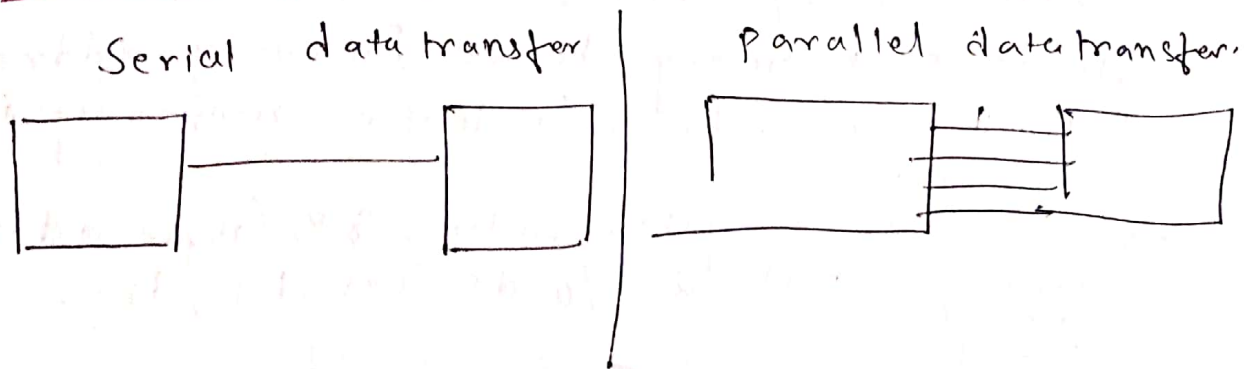
- more no. of Instructions to access the I/O devices
- more no. of addressing mode to access the I/O device

→ Disadvantage.

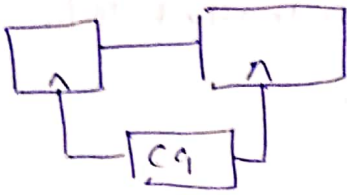
- a part of memory is wasted.

All the memory access instructions can be used to access I/O devices also because the control signals for memory and I/O devices are common.

* Serial asynchronous data transfer.

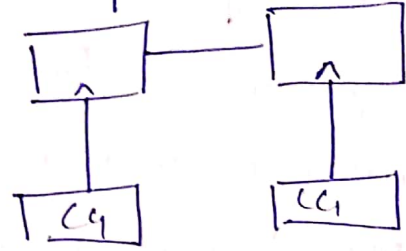


Synchronous



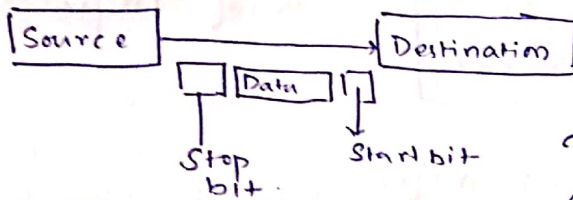
Here common clock generator is used & is tuned based on the slowest m/c which tends to lower performance

asynchronous



Here externally synchronisation is provided to reduce the data loss between m/c. but performance will be upgraded.

Synchronous Serial data transfer: In this a fixed amount of data is sent from source to destination which also takes '2' bits 1 for 'stop' & 1 for start bit this will synchronize the two machines. both sender & destination is aware with the time gap of data transfer.



assume.
data size = 8 bit.
a total size of $8 + 1 + 1 = 10$ bit.
are transferred to transfer a char of 8-bits.

* Modes of Data transfer:

⇒ Modes of data transfer:-

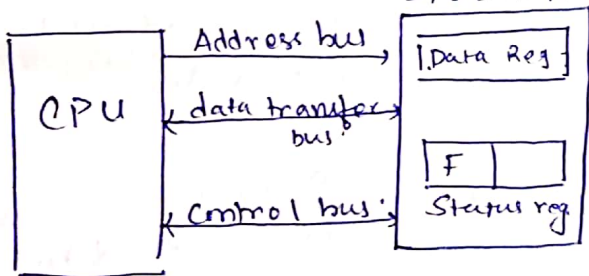
1. Programmed I/O or Program controlled I/O
2. Interrupt initiated / driven I/O
3. DMA (Direct memory Access).

→ Programmed I/O or program controlled I/O
There is no any provision (signal) through which I/O device can inform to CPU about data transfer.

→ In such case the device sets its own status and ~~wait~~ waits.

→ CPU runs a program periodically and checks the status of every device one by one.

→ If any device has its status set then CPU performs data transfer for that device.



Status register by default is of "1 Byte"

Total time required to transfer data = time to check the status of device
i.e. Time required to transfer status register (1B) from I/O device to CPU
+
time to transfer actual data.

disadvantage here is the wastage of CPU time. everytime it has to check for the status of I/O device.

→ Interrupt initiated I/O:-

I/O devices have a provision (interrupt signal) to inform CPU about data transfer.

— When CPU receives Interrupt from any device.

(i) CPU completes the execution of current instruction

(ii) PC value + other status values CPU stores onto the stack.

(iii) CPU takes a branch to service the interrupt.

(iv) after the service using the stack content CPU resumes previous process.

Interrupt Service Routine (ISR): a function execution of which services the interrupt.

Vector address: Starting address of ISR is known as vector address.

Interrupts are of two types based on vector address.

vector

Non-vector
or
scalar

In vectored Interrupt CPU will generate an Interrupt and it will also give vector address of ISR to resolve the interrupt service and service is done.

The device sends starting address of ISR (vector address) along with the interrupt signal. In such case CPU directly jumps to the provided address, executes ISR & services the interrupt.

→ Non vectored Interrupt.

The device sends only interrupt signal to CPU. In such case CPU first executes default service routine which helps the CPU to resolve the scenario and to obtain the actual ISR address.

→ Based on acceptability interrupts are of two types

Maskable

CPU can accept or reject

non maskable

always accepted
eg: Trap.

→ Based on source Interrupts are of two types

External (Hardware)

If any I/O device generates the interrupt.

Internal (Software)

During execution of any instruction if any unexpected error occurs then an internal interrupt is generated.

Internal interrupt are always non maskable.

eg for internal interrupt :

↳ $R_1 \leftarrow R_2 / R_3 ; R_3 = 0$

↳ Page fault .

↳ addressing error in segmentation .

↳ Exceptional handling in program .

If multiple devices generate the interrupt simultaneously then CPU will service the interrupt of highest priority device first .

Priority based interrupt handling .

Software solⁿ (Polling)

Hardware solⁿ

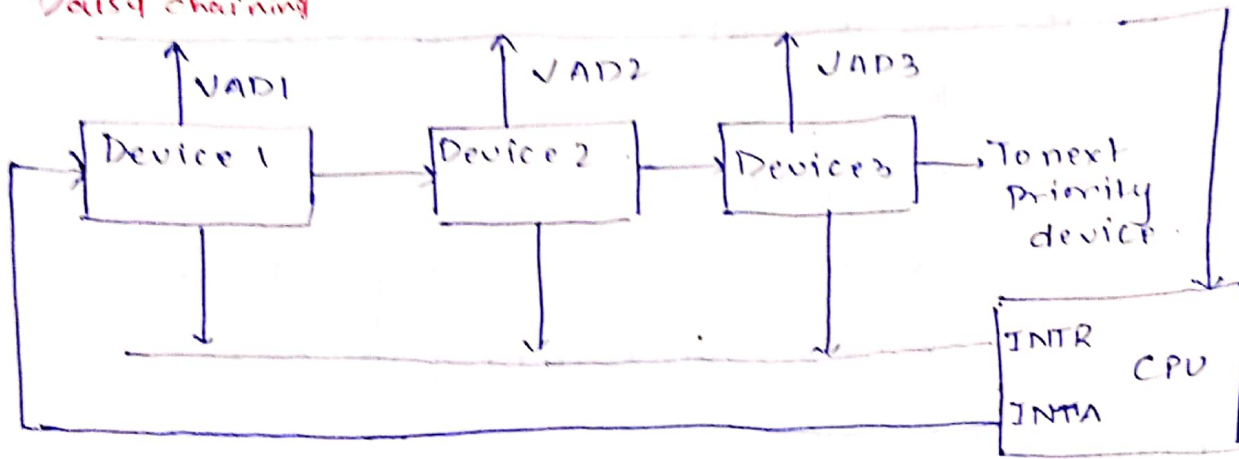
Serial
(Daisy chaining)

Parallel .

Polling: Whenever the interrupt occurs the program will check every devices one by one from which the interrupt is emitted then it will resolve the interrupt .

Polling will only be started when interrupt is encountered .

Daisy chaining



INTR: Interrupt Request.

INTA: Interrupt Acknowledgement.

VAD: Vector address.

Devices will create an interrupt requests and INTR will be set to '1' and CPU will generate an acknowledgement which will move to device 1 first and then after the Device 1 is done with interrupt service the CPU will check for further more interrupts available if yes then it will generate the acknowledgement and will send it to Device 1, Device 1 will pass it to another Device & so on. If the all the interrupts are handled then only the CPU will stop the process hence not all devices are scanned for interrupt only the one with interrupt.

In daisy chaining the device which electrically nearest to CPU will be having highest priority.

Else the generally the devices which is most frequently accessed by CPU will have highest priority i.e. HDD.

	Time.	Otherwise bits, bytes etc
k	10^3	2^{10}
M	10^6	2^{20}
G	10^9	2^{30}

$$\begin{aligned} \text{Time required in Interrupt I/O} &= \text{interrupt overhead time} + \text{data transfer time.} \\ &= 4 \mu\text{s} + 0 \\ &= 4 \mu\text{sec.} \end{aligned}$$

$$\text{Performance of Programmed I/O} = 1/100 \mu\text{sec.}$$

$$\text{Performance of Interrupt I/O} = 1/4 \mu\text{sec.}$$

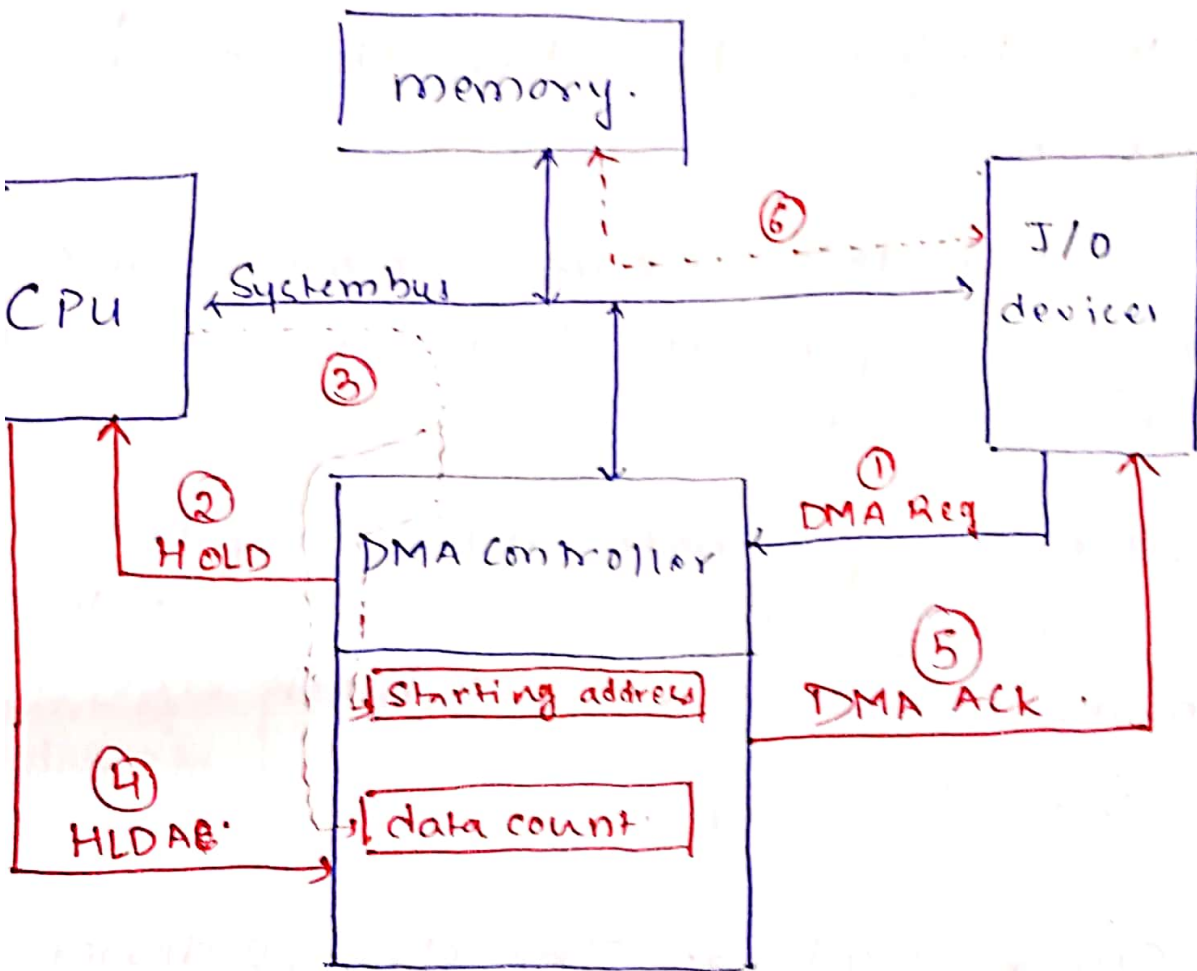
$$\text{Performance gain} = \frac{1/4}{1/100} = \frac{100}{4} = 25$$

$$\text{Performance gain} = \frac{\text{Slower technique time.}}{\text{Faster technique time.}}$$

* DMA :-

Direct memory access : It is a technique which enables direct data transfer between memory and I/O device without and Interferen of CPU.

To implement this technique ; A specific hardware is used known as DMAC (DMA Controller)



Starting address: Memory address starting from where the data transfer should be done.

data count: No. of bytes or words to be transferred.

AFTER 1B is transferred	AFTER 2B	AFTER 3B
-------------------------	----------	----------

Note: During DMA Transfer CPU cannot perform those operation which require system buses. which means the CPU will be blocked mostly. System bus is used by DMA controller.

* Modes of DMA Transfer:

Once the control of the buses is given to DMA controller how much data is transferred through DMA before CPU takes the control of the buses back.

↳ (i) Burst mode:- The entire amount of data is transferred before the CPU takes the control of the buses back.

↳ (ii) Block transfer mode:- A block of data is transferred before the CPU takes the control of the buses back. usually a block size is bet? [512 B to 1 kB]

↳ (iii) Cycle Stealing:- The slow IO device takes some time to prepare the data. during that time CPU keeps the control of the buses. Once the data is prepared CPU gives the control of the buses to ~~the~~ DMA controller for one cycle. in that cycle the prepared data is transferred to memory and then CPU takes the controls of Buses back.

Time to prepare the data = t_x .

time to transfer data to memory = t_y .

$$\% \text{ of time CPU blocked} = \frac{t_y}{t_x + t_y} \times 100\% \rightarrow \textcircled{1}$$

When I/O can start preparation.

of next data during transfer to memory.

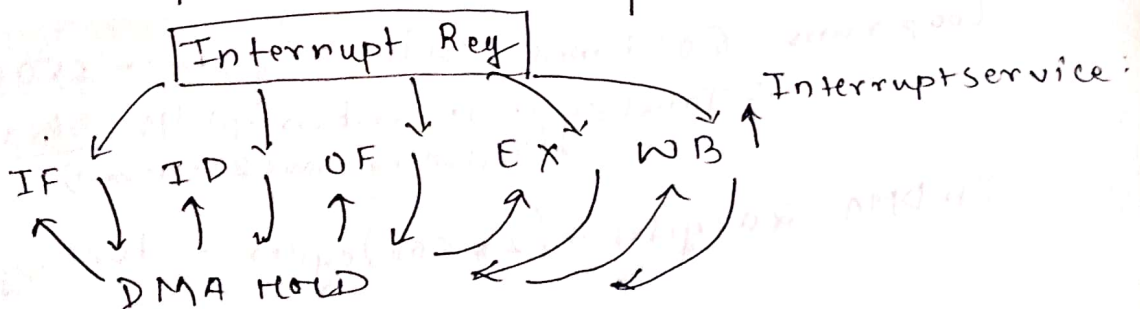
$$\% \text{ of time CPU blocked} = \frac{t_y}{t_x} \times 100\% \rightarrow \textcircled{2}$$

as t_y is included in t_x itself

IF t_x is larger then only use the second formula.

IF $t_x \geq 2t_y$. then only use second formula. else use the first one.

	Burst mode	Cycle stealing.
1) Amount of total data to be transferred.	Small	Large.
→ Amount of data transferred in one time	More.	less.
→ CPU waiting	More	less.
→ DMAC waiting	Less	more.



DMA controller is a special purpose processor which is used for data transfer between memory and I/O.

Q8: 64 Here 500B from I/O to memory using interrupt I/O which uses CPU. Hence it is not DMA. It is interrupt driven. program based I/O

1 cycle ← Initialize add. Reg.
 1 cycle ← Initialize Count = 500.
 2 ← Loop: load a byte from device.
 2 ← Store in the mem at add. given by add. reg.
 1 ← Inc add. reg.
 1 ← Dec the count.
 4 ← if count ≠ 0 goto Loop.

DMA
 500B from I/O to memory
 20 cycles for initialization.
 2 cycles/byte.

$$\text{Speed up} = \frac{\text{Slower time}}{\text{Faster time}}$$

$$= \frac{\text{Interrupt I/O time}}{\text{DMA time}}$$

Loop runs 500 times with 7 cycles = 3500
 ∴ Total cycle in Interrupt I/O = $\frac{3500}{2}$
 (Including above 2 statements)

In DMA 20 cycles + (2 * 500) cycles = 1020 cycles.

$$\therefore \text{speed up} = \frac{3502}{1020} = 3.4 \text{ (approx)}$$

Q2
68

Data count Reg = 16 bits

File 29,154 KB from disk to memory

Data count Reg says how many bits of data can be transferred at time.

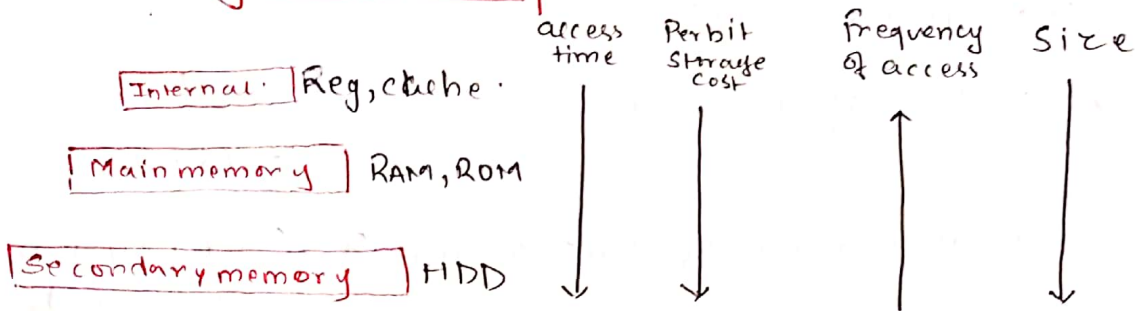
$$\begin{aligned} \text{Max value of data count register} &= 2^{16} - 1 \\ &= 65535 \end{aligned}$$

$$\begin{aligned} \therefore \text{minimum no. of time required} &= \frac{29154 \text{ KB}}{65535 \text{ B}} \\ &= \frac{29154 \times 2^{10} \text{ B}}{65535 \text{ B}} \end{aligned}$$

* Memory organization

a physical device which is use to store the contents in the system is known as memory.

→ Memory Hierarchy



→ arrow represents direction of increment.

Memory is divided among two functionality purpose

- ↳ (i) Maximize access speed.
- (ii) Minimize cost of bit per storage.

Memories	Size (in bits)	Per bit storage cost	Cost of memory
M_1	S_1	C_1	$S_1 C_1$
M_2	S_2	C_2	$S_2 C_2$
M_3	S_3	C_3	$S_3 C_3$
	$(S_1 + S_2 + S_3)$ bits		$S_1 C_1 + S_2 C_2 + S_3 C_3$ is Total cost of memory system

$$\therefore \text{average per bit storage cost} = \frac{S_1 C_1 + S_2 C_2 + S_3 C_3}{S_1 + S_2 + S_3}$$

$$= \frac{\sum_{i=1}^n S_i C_i}{\sum_{i=1}^n S_i}$$

S_i = size (in bits) of memory M_i
 C_i = per bit storage cost of memory M_i

memory representation:

The memory is represented by

= no. of cells * 1 cell capacity.

= no. of memory address * bits per address

no. of cells | 4 | 8 | 16 | 2^n | ∞

4 GB memory \rightarrow $\underbrace{4\text{ GB}}_{\text{no. of cells}} \times \underbrace{1\text{ B}}_{\text{size of each cell}}$

$$= 2^{32} \times 1\text{ B}$$

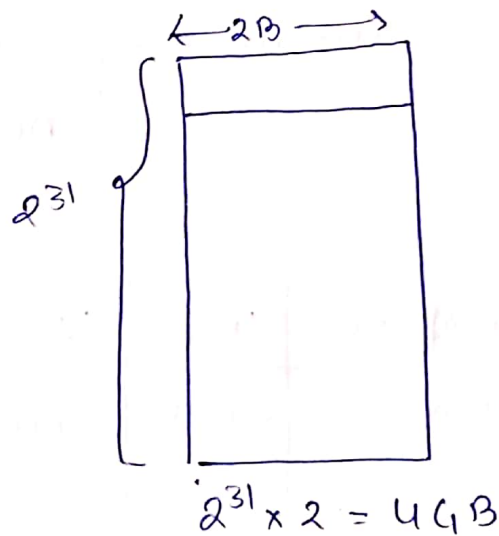
\therefore It will require 32 bit address size.

ie to have capacity of 4GB we require atleast 32 bit address size.

85
20 } 4GB
word addressable.
word size = 2B

$$\begin{aligned} \text{no. of words} &= \frac{4\text{ GB}}{2} \\ &= 2\text{ GB} \\ &= 2^{31} \text{ Bytes} \end{aligned}$$

address = 31-bits



Main memory.

It is used to store current running processes (Instructions) & their data.

Two types of Main memories.

\rightarrow ROM \Rightarrow Non-volatile.

\rightarrow RAM \Rightarrow Volatile.

~~or~~ During Booting process CPU does the POST & then it takes OS from the HDD-ROM and stores it to the RAM this

Bootling process instructions are also stored in ROM & ROM is also necessary in computer

Program responsible for Booting is Bootstrap program.

Bootstrap program is stored in ROM

OS is stored in HDD

But in modern days computer Bootstrap program is stored in HDD and Bootstrap loader is stored in ROM so that you can get privilege to change Bootstrap program.

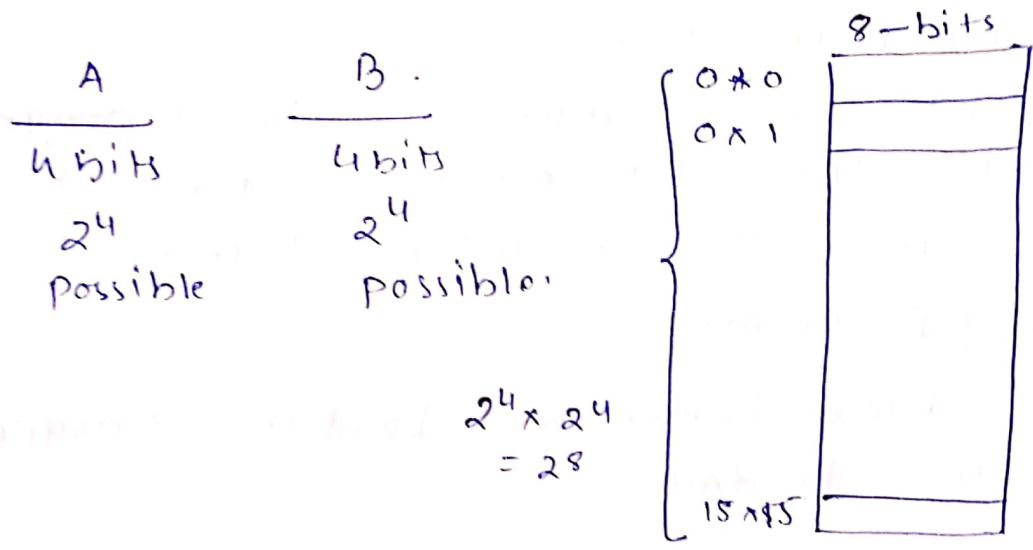
This Bootstrap loader will load the program from HDD to RAM.

Types of ROM :-

- (1) **EROM (Erasable ROM)** :- With the help laser or UV rays the content can be erased
- (2) **PROM (Programmable ROM)** user can perform one time writing in such ROM.
eg: CD-ROM.
- (3) **EPROM (Erasable & Programmable)** :- user can perform erase operation using UV rays which may require certain lab environment.
- (4) **EEPROM (Electrically erasable programmable ROM)**
Electric current can erase the content.
- (5) **EA PROM (Electrically alterable programmable ROM)**
Electric current can be used to alter the content.

GATE: The amount of ROM needed to store the table for multiplication of two 4-bit unsigned integer is?

For multiplication of 2 4bits will require 8 bits.



∴ Size of ROM required = $2^8 \times 8$
 $= 2^{11}$ bits
 $= 2k$ bits.

A	B	$A \times B$	$A + B$	Multiplication table size	Addition table size
4 bits	4 bits	8-bits	5 bits	$2^8 \times 8$ bits	$2^8 \times 5$ bits
n-bits	n bits	2n-bits	(n+1)-bits	$2^{2n} \times 2n$ bits	$2^{2n} \times (n+1)$ bits

Types of RAM

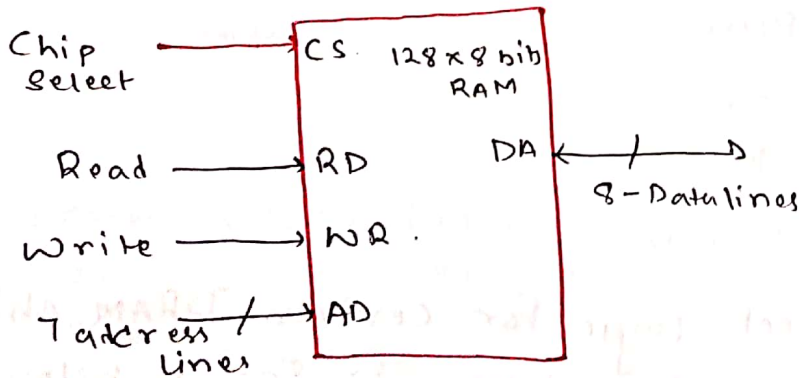
SRAM (Static RAM)

- Built using flip flops
- No Refresh or recharge
- Faster
- Expensive
- Used for cache memory

DRAM (Dynamic RAM)

- Built using capacitors. (In form of Electric charges & this charge tends to discharge.)
- Periodic refresh or recharge is required.
- Slower as refresh operation requires time.
- During refresh no read write is performed.
- Inexpensive.
- Used for Main memory.

RAM chip:

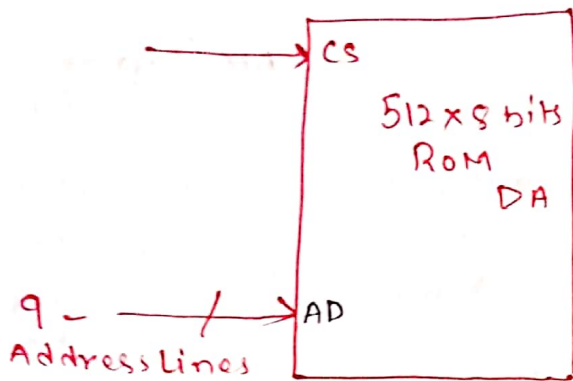


tree
Cross on line says many lines are involved
CS - 1 pin
RD - 1 pin
WR - 1 pin
Add - 7 pins
DA - 8 pins.

CS	RD	WR	Operation
0	X	X	no operation
1	0	0	no operation
1	0	1	write
1	1	X	Read

X → means any input doesn't matter

ROM chip:



CS	operation
0	No operation
1	Read.

8 Data Lines.

CS - 1 pin
AD - 9 pins
DA - 8 pins
18 pins.

no need of Read/write column as only Read is performed.

question - A ROM chip of 1024 x 8 bits has 4 select lines and operates from 5 volt power supply how many pins are needed for the IC package.

address = 10 pins.
data = 8 pins.

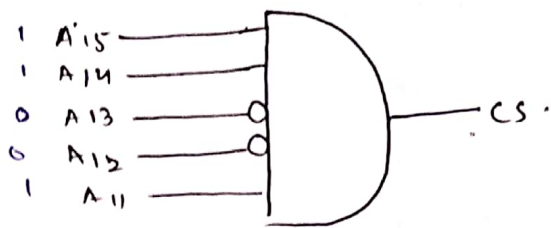
CS = 4 pins.

Power = 1 pin.
ground = 1 pin

24 pin.

$$\text{size} = 1024 \times 8 \text{ bit} \\ = 2^{10} \times 8$$

gate: The chipselect logic for certain DRAM chip in a memory system design is shown below. assume that the memory system has 16 address lines denoted by A₁₅ to A₀. what is the range of addresses (In Hexadecimal) of the memory system that can be enabled by the chipselect (CS) signal



- (A) C800 to CFFF
- (B) C800 to C8FF
- (C) DA00 to DFFF
- (D) CA00 to CAFF

to enable CS i.e. $CS = 1$ all the given lines should be ANDed to give out 1 and remaining can be any thing.

$$\therefore \begin{array}{cccc} 1 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 0 \end{array} \begin{array}{cccc} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{array} \text{ i.e. } \boxed{C800 \text{ to } C8FF}$$

$$\begin{array}{cccc} 1 & 1 & 0 & 0 \\ \hline C & F & F & F \end{array}$$

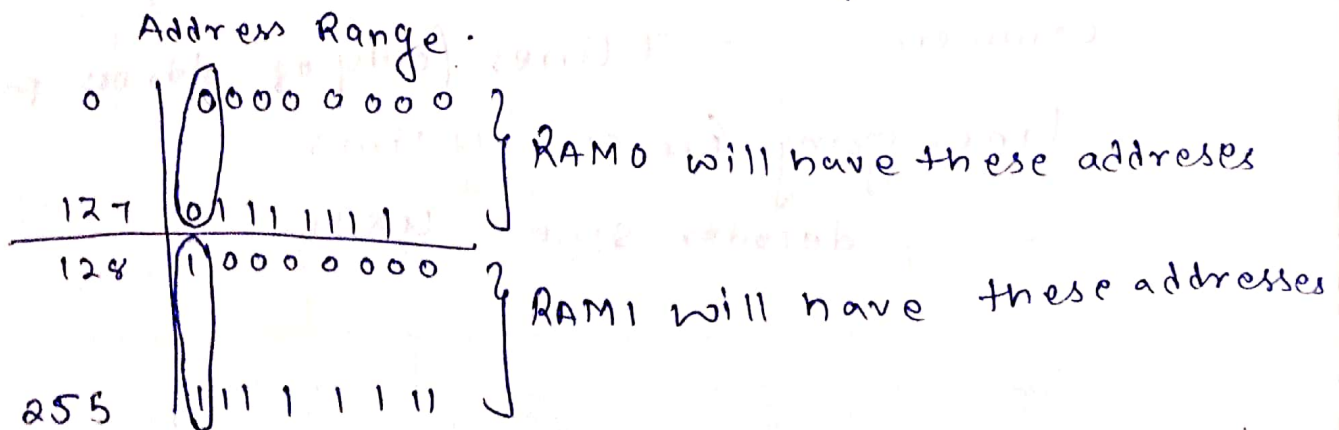
* Organising multiple RAM chips.

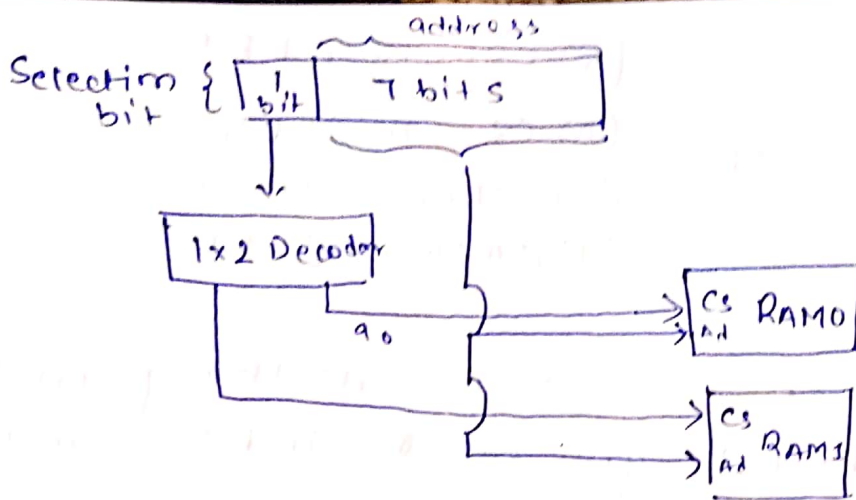
Total memory capacity = no. of chips * 1 chip capacity

Consider 2 RAMs of size 128×8 -bits each

$$\begin{aligned} \therefore \text{total capacity} &= 128 \times 8 \text{ bit} \rightarrow 7 \text{ bit address} \\ &+ 128 \times 8 \text{ bit} \rightarrow 7 \text{ bit address} \\ &\hline &256 \times 8 \text{ bit} \rightarrow 8 \text{ bit address} \end{aligned}$$

Previously only 7 bits were required to store addresses now 8 bit will be required.





question (i) How many 128×8 bits RAM chips are needed to provide a memory capacity of 2048 bytes?

(ii) 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 chips?

(iii) How many lines must be decoded for chip select specify the size of decoder.

$$\rightarrow \text{no. of chips} = \frac{\text{Total memory capacity}}{\text{1 chip capacity.}}$$

$$= \frac{2048 \text{ B}}{128 \text{ B}} = \frac{2^{11}}{2^7} = 2^4 = 16.$$

address lines = 11 lines.

common = 7 lines (only of address part)

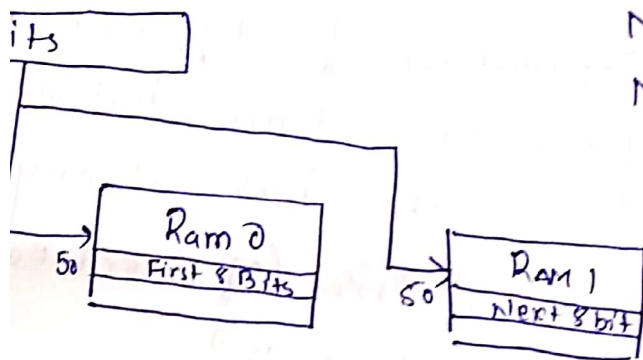
Lines going for CS = 4 lines.

decoder size = 4×16 .

How many 128×8 bit ram chips are needed to provide the memory capacity of 128×16 bits

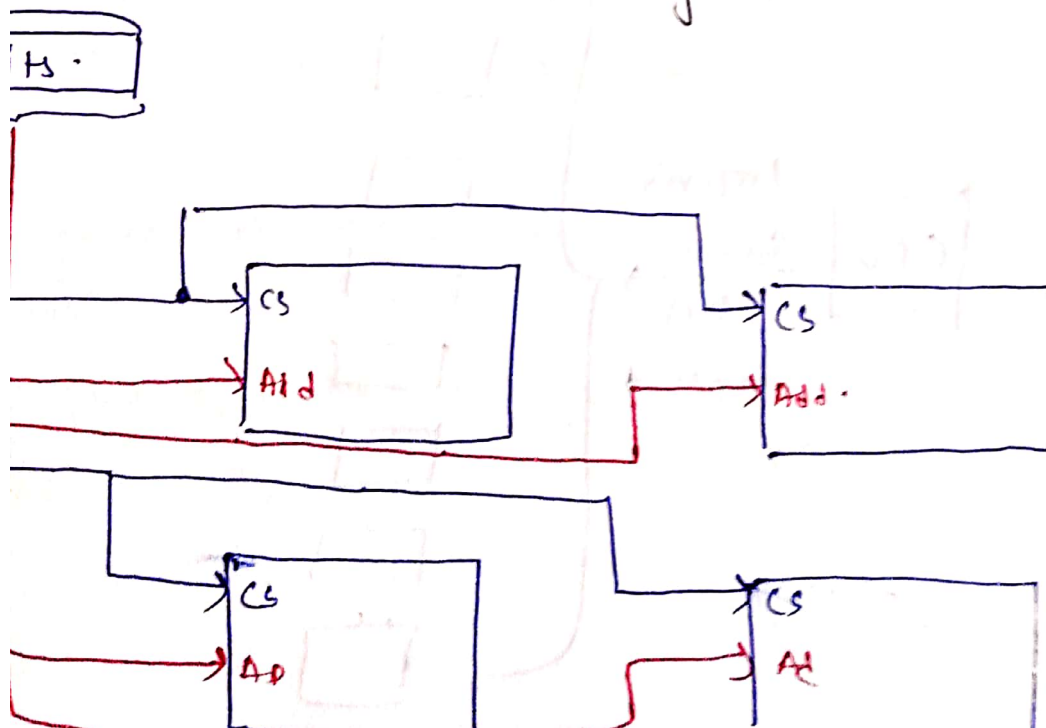
$$\begin{aligned} \text{No. of chips} &= \frac{\text{Total memory capacity}}{\text{1 chip capacity.}} \\ &= \frac{128 \times 16 \text{ - bits}}{128 \times 8 \text{ - bits.}} \\ &= 2. \end{aligned}$$

address no. of address is same but the stored per address is doubled hence we use horizontal arrangement where we use no. of bits of Address among 2 chips naturally



No. of addresses \Rightarrow same
No. of data bits per address \Rightarrow double

How many 128×8 bits ram chips are needed to provide a capacity of 256×16 bit.

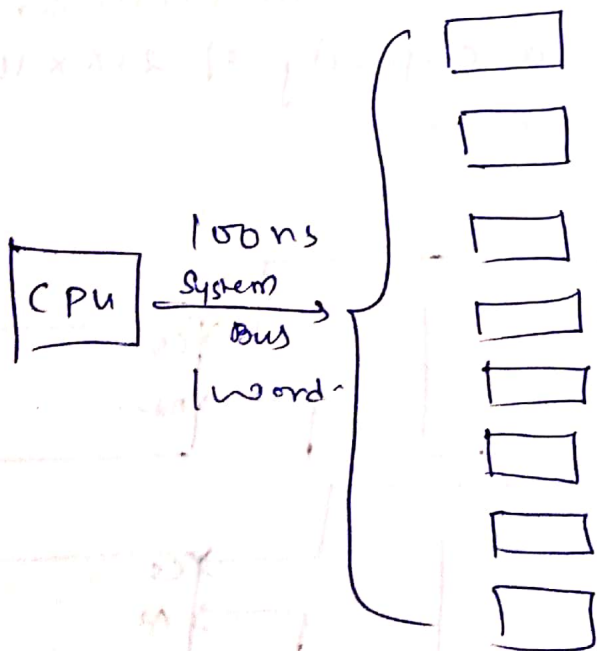


How many $32\text{ k} \times 1$ ram are needed to provide memory capacity of 256 k bytes.

Total capacity = $256\text{ k} \times 8$ (Default unit of capacity is bits)
 1 chip capacity = $\underline{32\text{ k} \times 1}$

\therefore Ram required = '64'

→ Consider a main memory system that consist of 8 memory modules attached to the system bus which is 1 word wide. When a write request is made, the bus is occupied for 100 ns and 500 ns for 500 ns there after, the addressed memory module executes 1 cycle accepting and storing the data. The (internal) operation of different memory modules may overlap in time, but only one request can be on the bus at any time. The maximum no. of stores (of 1 word each) that can be initiated in 1 ms is?



\therefore after every 500 ns the modules are freed.
 \therefore no ~~request~~ Request will be stopped.

$\therefore \frac{1\text{ ms}}{100\text{ ns}} = \frac{10^4}{100 \times 10^3} = \frac{10^4}{10^5} = 10^4 = 10000$

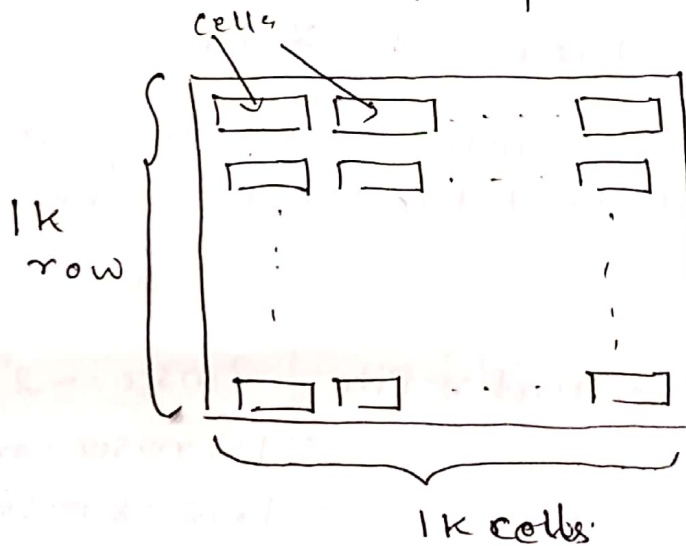
Arrangement of cells in DRAM chip:-

⇒ In 1 refresh operation 1 row of cells is refreshed.

$$1 \text{ chip refresh time} = \text{no. of rows in chip} \times 1 \text{ refresh time.}$$

memory system refresh time = 1 chip refresh time.
[all chips can be refreshed parallelly]

1MB DRAM chip (Byte addressable)

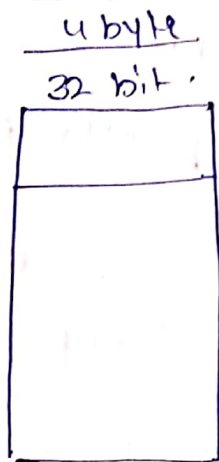


1 refresh operation time = 1 row refresh time
i.e. whole row is refreshed at a time.

→ a main memory unit with a capacity of 4MB is build using $1M \times 1\text{bit}$ Dram chips. each DRAM chip has 1K rows of cells with 1K cells in each row. the time taken for a single refresh operation is 100 ns. The time required to perform 1 refresh operation on all the cells in the memory unit is?

- (a) 100 ns (b) 100×2^{10} ns (c) 100×2^{20} ns.
(D) 3200×2^{20} ns.

Q 24
P 23



Requirement

1 GB \rightarrow 256M x 4 Byte.

Given.

256M x 4 Bits.

= 8 chips (Horizontal only)

no. of row = 2^{14}

1 operation = 50×10^9

1 chip refresh time = $2^{14} * 50$ nsec.

But after every 2 millisecond. 1 refresh time starts at that time nothing read & write will occur.

Time remaining for read & write = $2 \text{ msec.} - 2^{14} * 50 \text{ nsec.}$
= 1.2 msec. or
= 1.1808 msec.

$2^{14} * 50 \text{ nsec.}$

= $2^4 * 2^{10} * 50 \text{ nsec.}$

2^{10} can be taken as 1000 or 1024

= $16 * 1k * 50 \text{ nsec.}$

or $16 * 50 * 1024$

= 819200 nsec

= 800 Usec

= 0.8192 msec.

= 0.8 msec.

\therefore Time remaining

% time for read/write = $\frac{1.2}{2} * 100\%$

= 60%

or

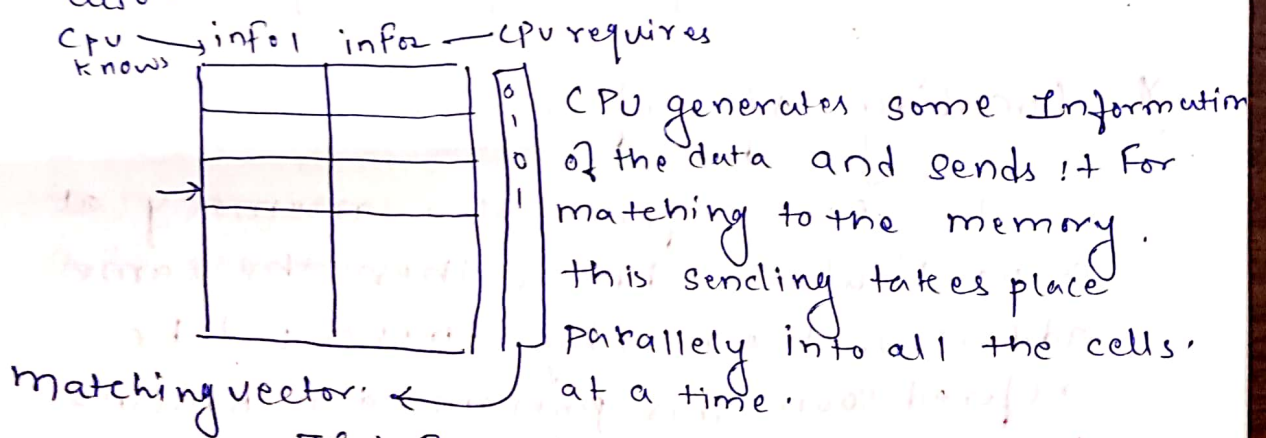
1.1808

* Associative Memory:-

No any address for the cells in this memory hence searching is performed based on the content stored in the cells. That is why this memory is known as content addressable memory (CAM)

Every cell ^{of memory} contains 2 info? 1 which CPU knows & 1 which CPU requires.

Every cell of memory contains a matching logic hardware to match the content in cells. Hence this memory becomes very faster but expensive also.

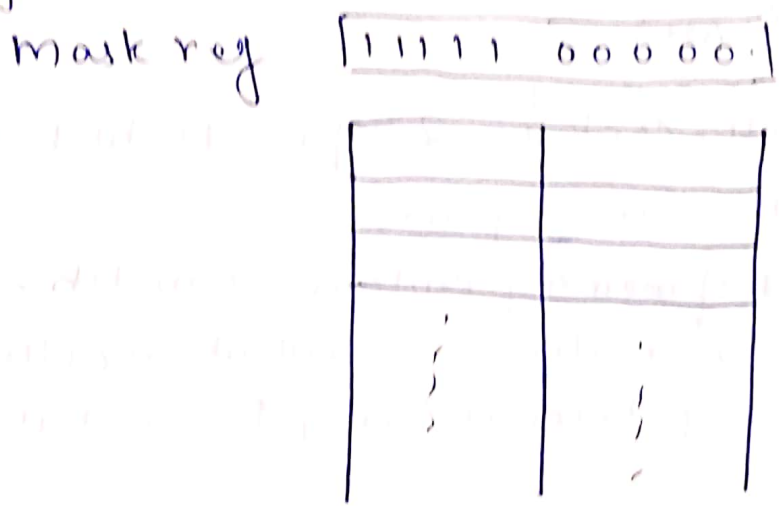


If information matches with any cell its corresponding cell information is copied into CPU.

This memory is used to implement Cache & TLB

matching is done by looking at the matching vector. if multiple matching occurs then all the corresponding values are sent to CPU sequentially

Comparison or information blocks Info 1 & Info 2 can be divided according to users choice. This variable or variation in information block is looked by "Mask register"



* Locality of reference.

IF CPU referce the memory at a particular address then the same address or the nearby addresses will be refered soon. This phenomena is known as

Locality of Reference.

Temporal
(time)

IF same address
Is refered again

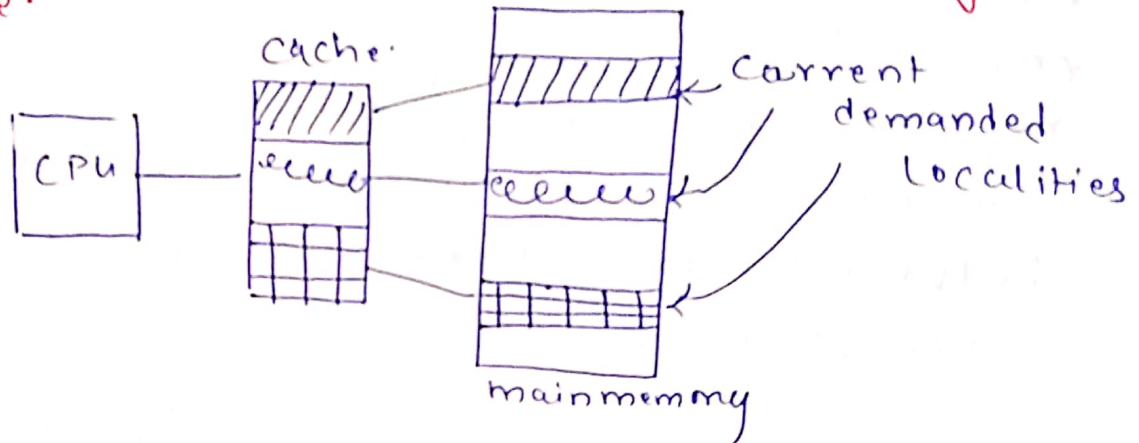
Spatial
(space)

IF nearby address
are refered.

* cache memory.

Based on Locality of reference concept: the current demanded localities (Block) are kept into a smaller and faster memory known as cache memory.

use of cache reduces over all memory access time.



Working of cache:-

If CPU's demanded content is present in cache it is known as cache hit.

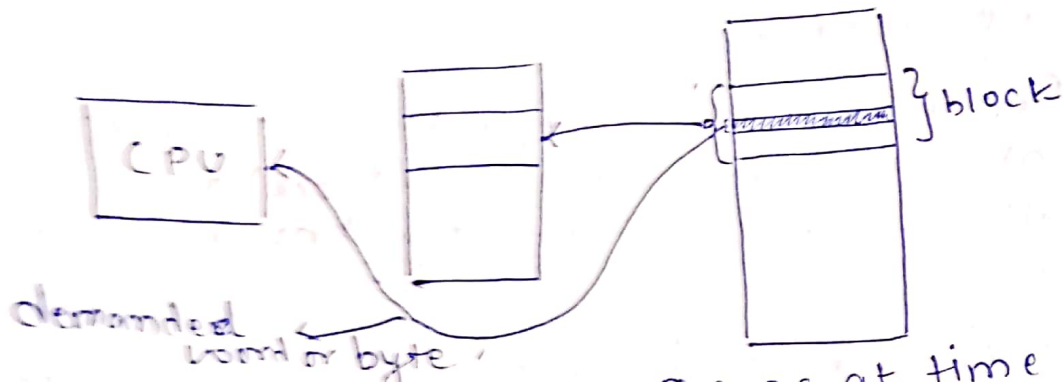
∅ If not present in cache its cache miss.

Performance of cache is given by cache Hit ratio

$$\Rightarrow \text{cache hit ratio (h)} = \frac{\text{no. of hits in cache}}{\text{total memory references}}$$

$$\Rightarrow \text{miss} = 1 - h$$

In case of cache miss 1 block (containing the demanded content) is copied from main memory to cache. because of this if the other content of block CPU demands then there will be hit in cache.



while copying a block ~~are~~ or at time of miss ~~the~~ the block is transferred to cache and the required info is transferred to CPU at the same time. So total time taken will only be the total time taken to copy the whole block to cache. as it is maximum

$$\therefore \text{average memory access time} = \text{Hit ratio} * \text{mem access time for hit}$$

$$\text{average memory access time} = \text{Hit ratio}(h) * (\text{memory access time for hit}) + (1-h) * (\text{mem access time for miss})$$

Types of cache access :-

1) Simultaneous access.

request for both cache & main memory generated simultaneously.

When there is hit then only cache access time is required and RAM access request is discarded but in case of Miss the RAM access is already present.

$$\therefore t_{avg} = H * t_{cm} + (1-H) * t_{mm}$$

Where t_{avg} = avg memory access time

$H =$ hit ratio; $t_{cm} =$ cache memory access time.
 $t_{mm} =$ main (RAM) memory access time.

(ii) Hierarchical access:-

only the faster memory is accessed first i.e. cache. If hit then great or if missed then the miss phenomenon will be initiated to access main memory

$$t_{avg} = H * t_{cm} + (1-H)(t_{cm} + t_{mm})$$

$$= H t_{cm} + 1 t_{cm} + t_{mm} - H t_{cm} - H t_{mm}.$$

$$t_{avg} = t_{cm} + (1-H) t_{mm}.$$

If in question Hierarchy or level given then use Hierarchical access else by default Simultaneous is used.

* If locality of reference should be used

(i) Simultaneously access:-

at time of miss the whole block is accessed

$$t_{avg} = H * t_{cm} + (1-H) t_{block}$$

$$t_{block} = \text{mm. access time for a block}$$

$$= \text{block size} * t_{mm}.$$

(ii) Hierarchical access:

$$t_{avg} = t_{cm} + (1-H) t_{block}.$$

* If memory access and ~~to~~ block transfer time should be needed then. Previously only reading was done now write or transfer

(i) Simultaneous access:-

$$t_{avg} = H * t_{cm} + (1-H)(t_{block} + t_{cm})$$

(ii) Hierarchical access:-

$$t_{avg} = t_{cm} + (1-H)(t_{block} + t_{cm})$$

~~***~~ in one access of cache memory entire block can be accessed

→ If t_{cm} & t_{mm} is not given \Rightarrow use general formula

→ If t_{cm} & t_{mm} given \Rightarrow use (i) If word hierarchy or level given then use hierarchical.
else use simultaneous

Q 18
P 22

$$0.8 * 5 + 0.2 * 50 = 14 \text{ nsec.}$$

Question: In the two level memory hierarchy if the top level has an access time of 8 ns & the bottom level has an access time of 60 nsec. What is the hit rate required on a top level to give an ~~access~~ average access time of 10 ns.

top level \Rightarrow mem which is closest to CPU.

$$t_{avg} = t_{cm} + (1-h)t_{mm}$$

$$10 = ~~60~~ 8 + (1-H)60$$

$$H = 0.96 \text{ i.e. } 96.6\% \text{ of Hit ratio}$$

→ a computer system contains a cache, uncached memory access takes 8 times longer than to access cache if the cache has hit ratio of 0.9 then ratio of cached memory access time to uncached memory access time is?

$$\begin{aligned} t_{cm} &= \frac{1}{8} t_{mm} \\ h &= 0.9 \end{aligned}$$

Uncached memory means the memory system without cache i.e. only main memory.

Cached means memory system with cache (cache + m.m)

$$t_{cm} = \frac{1}{8} t_{mm}$$

$$\text{or } t_{mm} = 8 t_{cm}$$

$$h = 0.9$$

$$\frac{\text{Cached access time}}{\text{Uncached access time}}$$

$$= \frac{(0.9) \times 1 + (0.1) \times 8}{1} = \frac{1.7}{1} = \underline{\underline{0.2125}}$$

* Cache writing :-

IF CPU perform write on a cache content. then corresponding main memory content should also be updated.

→ Write through.

→ Write back.

Write through: IF ~~cache~~ CPU updates the cache simultaneously main memory should also be updated. i.e.

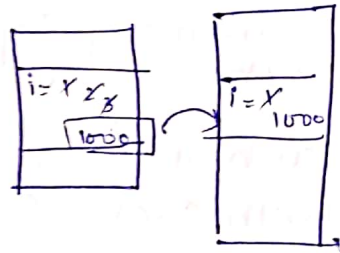
i.e. as soon as the cache memory content is updated corresponding main memory content is also updated immediately.

😊 Same value for cache content in cache & main memory always

😞 time consuming.

→ Write back:-

If CPU updates the cache contents then let it happen the main memory content is not updated immediately. When the updated block is to be replaced then only it is copied back to main memory.

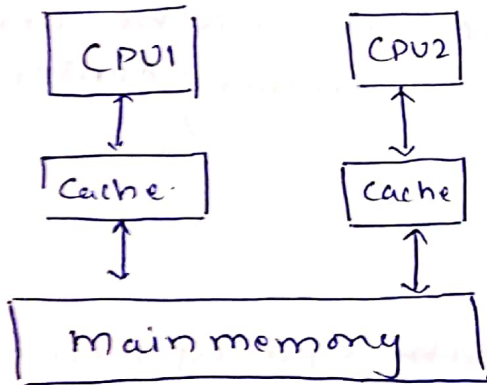


every operation gets completed on the cache then only the result is forwarded to RAM.

Cache memory will only be accessed by CPU.

Two different devices trying to access same content from two different memory will lead to problem of coherence.

* Cache coherence: In multiprocessor system.



In write through:-

$$t_{avg \text{ read}} = t_{avg} \text{ formulas}$$

$$t_{avg \text{ write}} = \max(t_{cm}, t_{mm})$$

where always t_{mm} will be higher
= t_{mm} .

$$t_{avg}(\text{read} + \text{write}) = \% \text{ of read} * t_{avg} \text{ read} + \% \text{ of write} * t_{avg} \text{ write}$$

Question

$$t_{cm} = 100 \text{ nsec.}$$

$$t_{mm} = 1000 \text{ nsec}$$

$$H_{\text{read}} = 80\%$$

Write through policy.

$$\% \text{ of read} = 70\% \quad \% \text{ of write} = 30\%$$

$$t_{avg} \text{ read} = ?$$

$$t_{avg} \text{ write} = ?$$

$$t_{avg} \text{ read} + \text{write} = ?$$

$$\begin{aligned} t_{avg} \text{ read} &= (H) t_{cm} + (1-H) t_{mm} \\ &= (0.80) 100 \text{ n.} + (0.20) 1000 \\ &= 280 \text{ ns.} \end{aligned}$$

$$t_{avg} \text{ write} = 1000 \text{ nsec.}$$

$$t_{avg} \text{ read} + \text{write} = 0.7 * 280 + 0.3 * 1000 = \underline{496 \text{ ns}}$$

gate 2014 The memory access time is 1 ns for a read operation with a hit in cache; 5 ns for a read operation with a miss in cache; 2 ns for a write operation with a hit in cache & 10 ns for a write operation with a miss in cache. execution of a sequence of instruction involves 100 instruction fetch operations 60 memory operand read operations & 40 memory operand write operations. the cache hit ratio is 0.9 the average memory access time (in nanoseconds) in executing the sequence of instructions is?

	hit	miss
read	1 nsec	5 nsec
write	2 ns	10 nsec.

$$H = 0.9$$

$$t_{avg \text{ read}} = 0.9 * 1 + 0.1 * 5 = 1.4 \text{ ns}$$

$$t_{avg \text{ write}} = 0.9 * 2 + 0.1 * 10 = 2.8 \text{ ns}$$

execution of seq of inst?

↳ 100 instruction fetch \Rightarrow read = 100

↳ 60 operand read \Rightarrow 60 read

↳ 40 memory write.

$$\therefore 100 + 60 = 160 \text{ read operations}$$

$$40 \text{ write}$$

\therefore 80% read & 20% write.


$$\therefore t_{avg} = \frac{80}{100} * 1.4 \text{ nsec} + 0.2 * 2.8 \text{ ns}$$

$$= \underline{\underline{1.68 \text{ nsec}}}$$

* Cache mapping

CPU always generates main memory address the size of main memory is larger which cannot be stored in cache memory fully

There should be a pattern with the help of which the main memory content from a specific address is copied to a specific location in cache memory this pattern is known as cache mapping

CPU  This pattern helps us to search in cache using the main memory address i.e. if that pattern is present in cache then its a hit else its missed.

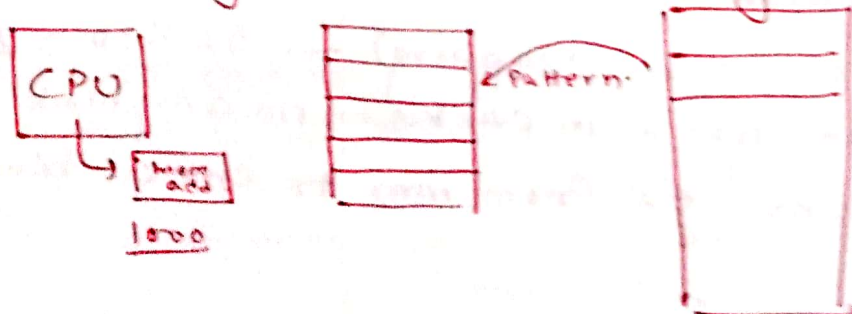
Transformation of main memory data into cache memory is known as cache mapping.

* there are three types of mapping.

(i) Direct mapping

(ii) Set associative mapping

(iii) Fully associative mapping.

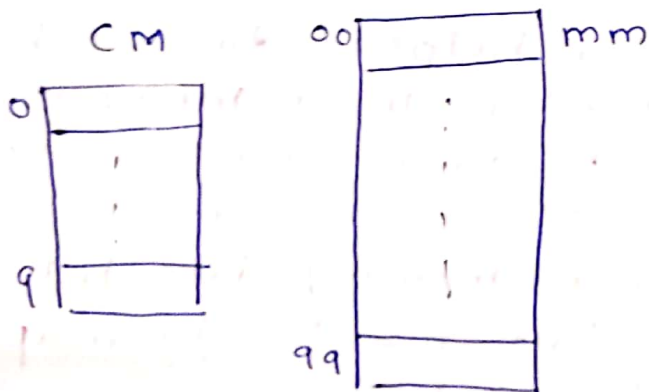


* Direct mapping :-

No. of blocks in CM = 10 (0-9)

no. of blocks in MM = 100 (00-99)

CM	0	1	2	3	4	5	6	7	8	9
MM	00	01	02	03	04	05	06	07	08	09
	10	11	12	13	14	15	16	17	18	19
	20	21	22	23	24	25	26	27	28	29
	30	31	32	33	34	35	36	37	38	39
	40	41	42	43	44	45	46	47	48	49
	50	51	52	53	54	55	56	57	58	59
	60	61	62	63	64	65	66	67	68	69
	70	71	72	73	74	75	76	77	78	79
	80	81	82	83	84	85	86	87	88	89
	90	91	92	93	94	95	96	97	98	99



$$\text{CM block no} = \text{mm block no.} \% \text{ no. of blocks in CM (mod)}$$

CPU demands mm blocks — 52.

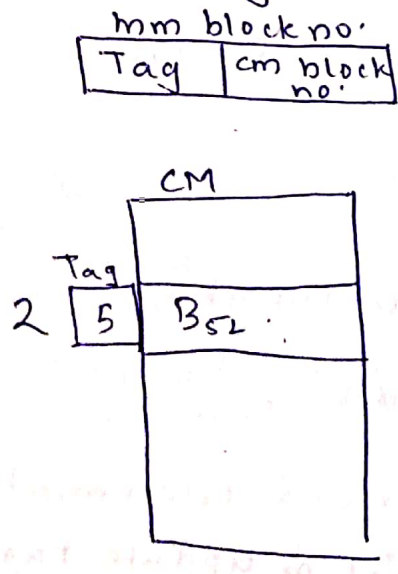
CM block no. (mapping) — $52 \% 10 = 2$.

check block no. 2 in cache \Rightarrow no any block \Rightarrow miss

'bring block no. 52 from mm to cm at block '2'

cpu demands mm block	cm block no. mapping	Operation
52	$52 \% 10 = 2$	Check block no. 2 in cache \Rightarrow no any block is present \Rightarrow miss bring block 52 from mm to cm at block no. 2.
96	$96 \% 10 = 6$	miss & bring block no. 96 from mm to cm at block no 6.
32	$32 \% 10 = 2$	check block no. 2 in cache \Rightarrow Block 52 present \Rightarrow miss bring block no. 32 from mm to cm at block block 2 by replacing block 52.

So 2, 12, 22, ... 92. all these blocks are competitors for being in cache memory at block no. 2. therefore every block we have 10 competitors each. But cache memory cannot differentiate the competitors block ID present in Block as it is in Binary. So with every Block ID some another field is attached to distinguish between every competitors. That distinguisher is known as Tag.



eg: mm block no \Rightarrow 52

5	2
---	---

Tag cm block no.

\downarrow \downarrow

goto block no. 2.

\downarrow

no block present

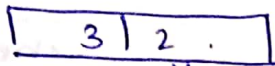
\downarrow

miss

\downarrow

bring block 52 in cm block 2

mm block no. 32



↓
goto CM block no. 2.



Tag present = 5

Tag in request = 3

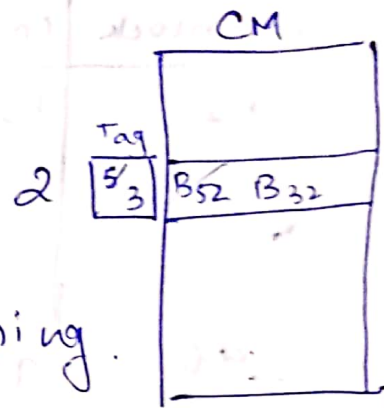
no matching



miss

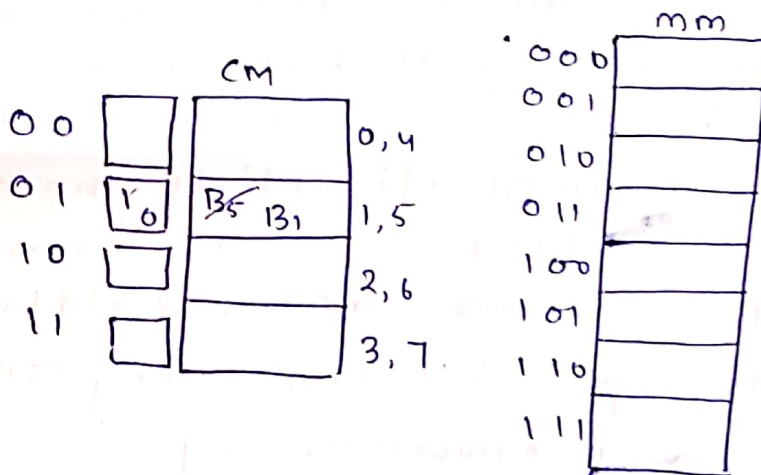


bring block 32 in cache & update tag to 3.

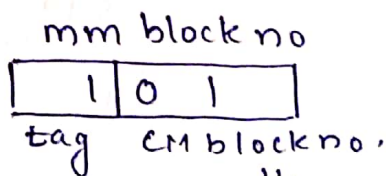


no. of blocks in CM = 4

no. of blocks in mm = 8



CPU Req => MM block no. required S = 101



↓
goto block no. 01 in cache



no block => miss



Bring block no. 5 from mm to CM in
block no. 01 & update tag by 1

CPU req \Rightarrow mm block no. = 1 \Rightarrow 001

mm block

0	0	1
---	---	---

tag

\Downarrow
goto block no. 01 in cache.

\Downarrow

tag present = 1

tag req = 0

} no match

\Downarrow
miss

\Downarrow

Bring block 1 from mm &

replace block no. 5 & update

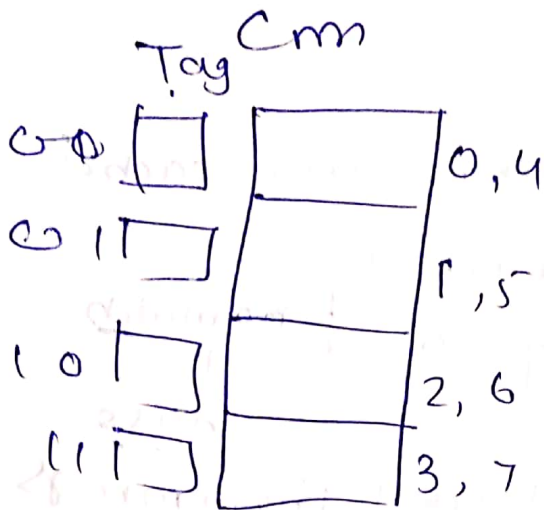
tag by 0.

The tags are distributed and is stored inside the cache controller.

Cache controller maintains one tag information for each block in cache.

\therefore Size of tag directory = no. of blocks * Tag bits
in CM

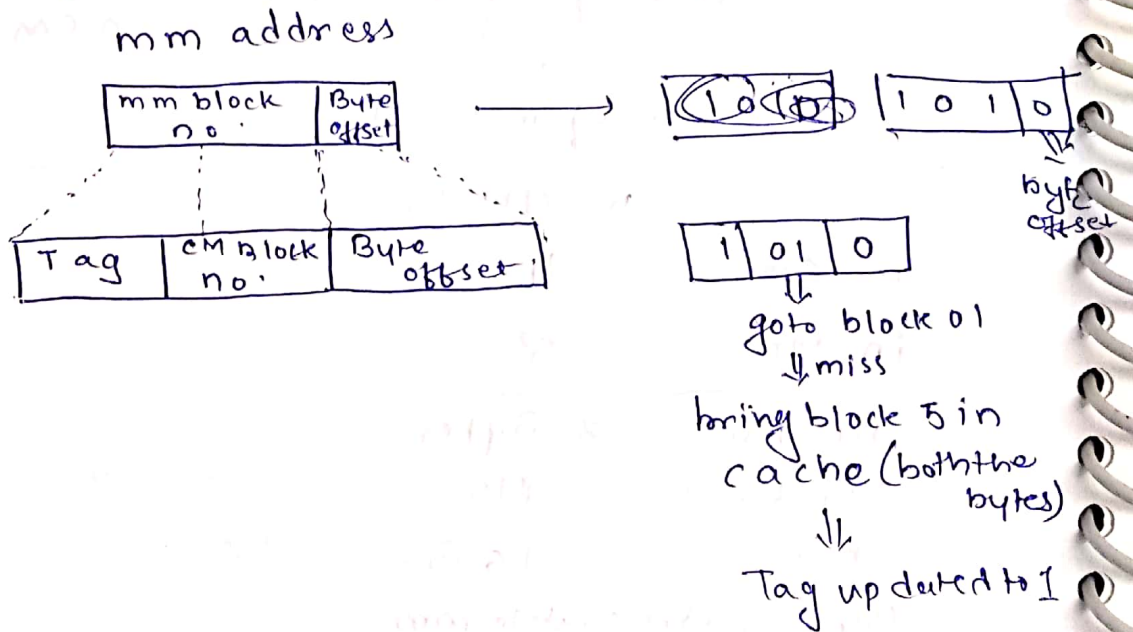
or
Size of meta-data.



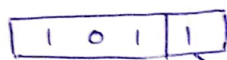
0000	Block 0
0001	
0010	Block 1
0011	
0100	Block 2
0101	
0110	Block 3
0111	
1000	Block 4
1001	
1010	Block 5
1011	
1100	Block 6
1101	
1110	
1111	Block 7

If CPU generates any address each block has 2 Byte which can be represented by 0, 1 in the LSB of Address. Byte no. \Rightarrow Byte offset.

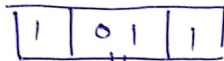
\therefore CPU generates address \Rightarrow 1010



CPU generates add 1011



block byte offset = 1



goto cm block 01 \Rightarrow tag present = 1

tag required = 1

Hit

Send byte no. 1 to CPU
(nothing to replace)

In this case replacement will only take place when entire block is mismatched.



bits in byte offset = $\log_2(\text{Block size})$

$$\text{Tag directory size} = \text{No. of blocks in cache} * \text{Tag size}$$

$$= 2^{12} * 16\text{-bits}$$

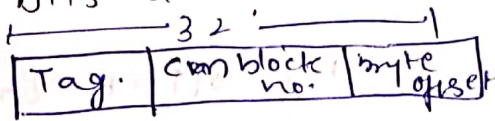
$$= 64 \text{ k bits}$$

$$= 8 \text{ k bytes}$$

Cachememory block no. is also known as index

question consider a direct mapped cache of size 2048 32 KB with block size 32 bytes.

CPU generates 32 bit addresses the no. of bits need for cache indexing & the no. of tag bits are?



$$\text{no. of bits in byte offset} = \log_2(2^5) = \underline{5} \text{ bits}$$

$$\text{no. of bits in CM block no.} = \log_2(2^{10}) = \underline{10} \text{ bits}$$

$$\text{no. of block} = \frac{\text{Cache size}}{\text{block size}} = 2^{10}$$

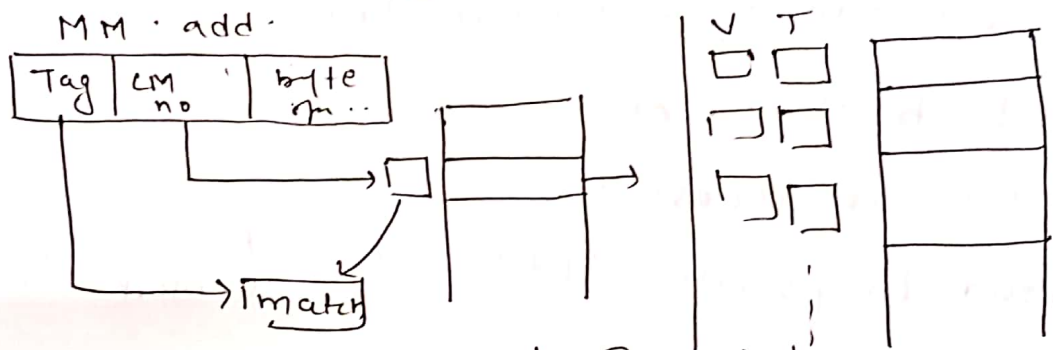
$$\therefore 32 - 15 = \underline{17} \text{ bits for tag}$$

→ cache initialization.

Initially when CPU is turned on the tag bit may have certain garbage value and also the cache memory.

hence after the CPU send some request the cache may generate 'hit' and CPU will get garbage value.

∴ we introduce another bit (Valid bit) combining 0 & 1

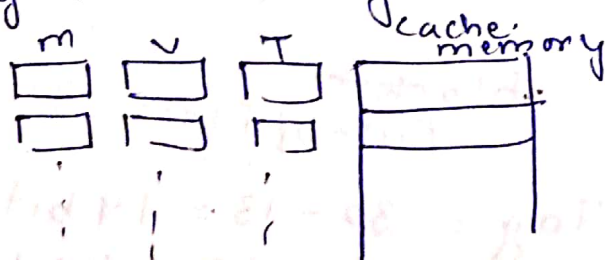


Cache is initialized from '0' on every power on which says it is invalid.

→ Performance improvement of write back cache.

There is block present in cache where CPU has only performed Read hence there is no need to replace the corresponding main memory location.

So only those locations are supposed to be modified which uses write operation hence to keep track of this we add modifying bit or dirty bit on the address.



If modifying bit is 1 then there is write operation and if modifying bit is 0 there is only read

$$\therefore \text{Tag directory size} = \text{No. of blocks in cache} * \left[\begin{array}{l} \text{Tag bits + valid bits} \\ + \text{modifying bit} + \\ \text{any extra bit} \end{array} \right]$$

only those among all the blocks accessed only those blocks are copied back to Main memory which has a dirty block or modification or has a write operation on them.

⇒ For write back cache.

Simultaneous access:

$$t_{\text{avg read}} = t_{\text{avg write}} = H * t_{\text{cm}} + (1-H) * (t_{\text{block}} + t_{\text{write back}})$$

$$t_{\text{write back}} = \alpha * t_{\text{block}}$$

α = Fraction of dirty blocks

Hierarchical :-

$$t_{\text{avg read}} = t_{\text{avg write}} = H * t_{\text{cm}} + (1-H) * (t_{\text{cm}} + t_{\text{block}} + t_{\text{write back}})$$

96
P20

$$\text{Cache size} = 8 \text{KB}$$

$$\text{block size} = 32 \text{B}$$

$$\text{MM add} = 32\text{-bits}$$

$$\text{No. of blocks in cache} = \frac{8 * 2^{10}}{2^5} = 2^8$$

8 bits to store it

$$\text{Block Byte offset} = 32 = 2^5 \Rightarrow \underline{5 \text{ bits}}$$

$$\text{Tag} = 32 - 13 = 19 \text{ bits}$$

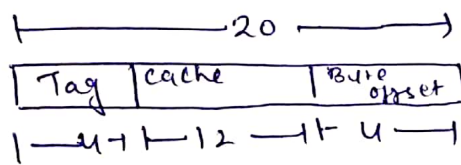
$$\Rightarrow \text{Tag director size} = 2^8 * (19 + 1 + 1) = 5376 \text{ bits}$$

Q 19
P 22

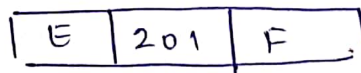
main memory size = 2^{20} Bytes of 1 Byte each
 \Rightarrow 20 Bits to store this address

No. of blocks in cache = 2^{12}

block size = 16 B = 2^4 \Rightarrow 4 bits for byte offset.



mm add \Rightarrow (E 201 F)₁₆ = convert in binary
 divide these bits according to block

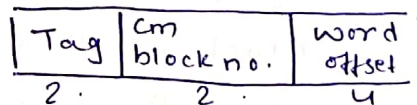


Q 15
21

Cache size = 64 words

mm size = 256 words \Rightarrow address in mm = 8-bit

block size = 16 words = 2^4 \Rightarrow word offset = 4-bits

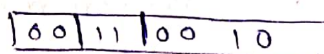


no. of blocks in cache = $\frac{64}{16} = 4 \rightarrow 2^2$

block	tag
0	10
1	10
2	00
3	01

cm block no. size = 2

word no. \Rightarrow 50 \Rightarrow ~~01000010~~
00110010

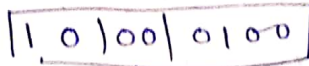


\Downarrow
block no. 3 in cache.

\Downarrow tag there = 01, tag required = 00
miss

\Downarrow
add 50 in 3rd block.

132 = 010000100



Block no 00

⇓

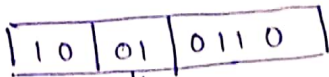
Tag there = 10 } hit

Tag request = 10

⇓

132 present

150 = 10010110



⇓
goto block 1

⇓

tag there = 10

Tag request = 10

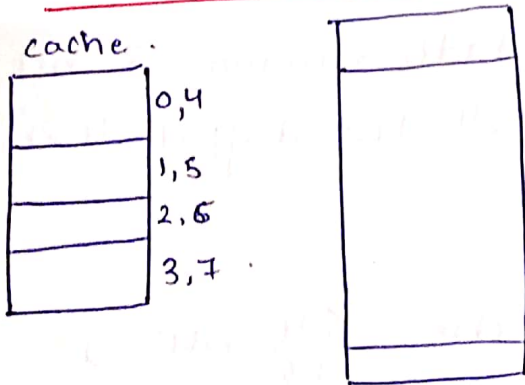
} hit

⇓

150 present

∴ 150 & 132 is present

* Set associative:



In previous arrangement if we get successive request of CPU toward a same block will increase the no. of misses in cache.

eg: 1,5, 1,5, 1,5, ...

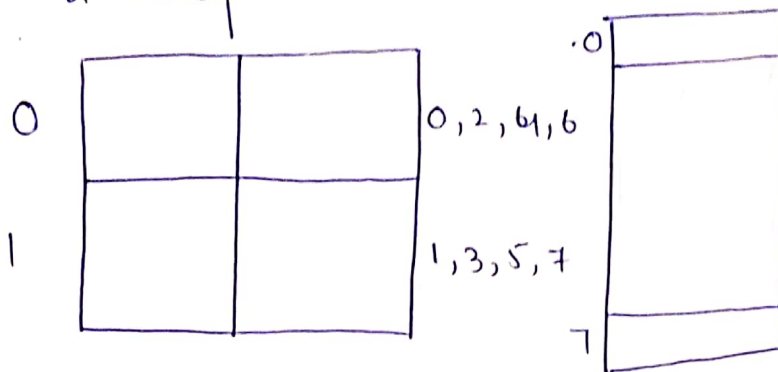
↓ miss | miss miss
replaces

main memory block no = 1,5,1,5,1,5

Cache memory block no = 1 always

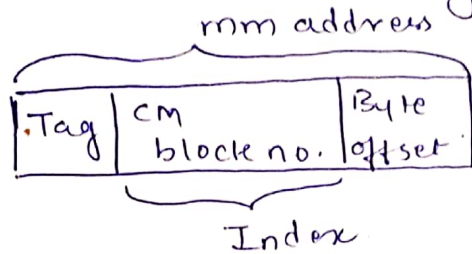
In set-associativity the cache blocks are divided into sets of blocks.

2-way set associative.

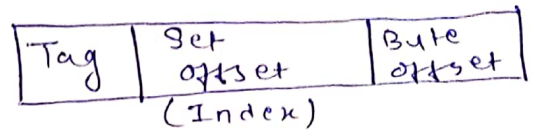


CPU Request:
 mm block no:-
 1, 5, 1, 5, 1, 5, ...
 ↓ miss ↓ miss
 hits

In direct mapping.



In set associative.



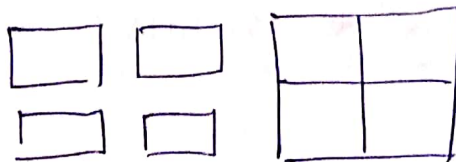
$$\text{no. of sets in cache} = \frac{\text{no. of blocks in cache}}{\text{associativity}}$$

$$\text{bits in set offset} = \log_2(\text{no. of sets in cache})$$

associativity \rightarrow k-associative (generally 2, 4, 8, ...)

\rightarrow Tags in set asso. is also stored in same way every block in a set has its own tag.

$$\therefore \text{Tag directory} = \text{no. of blocks in cache} * [\text{Tag} + \text{extra bits}]$$



But in set-associativity the no. of competitors per block increases. hence will the no. of bits in tag. i.e. size of tag block in set associative is more as compared to direct mapping.

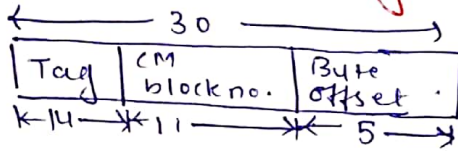
Question.

mm. add = 30-bits

Cache size = 64KB

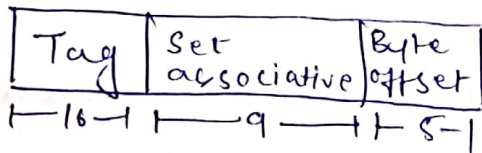
block size = 32 Bytes.

In direct mapping:



$$\begin{aligned}\text{no. of blocks in cache} &= \frac{64\text{KB}}{32\text{B}} \\ &= 2\text{K} \\ &= 2^{11}\end{aligned}$$

In 4-way set associativity.



$$\text{no. of set in cache} = \frac{\text{no. of blocks}}{\text{associative}} = \frac{2^{11}}{4} = 2^9$$

$$\begin{aligned}\text{Set\#} \\ \text{offset} &= 9\text{-bits.}\end{aligned}$$

$$\begin{aligned}\text{Size of tag directory} \\ &= \text{no. of blocks in cache} * \text{Tag bits} \\ &= 2^{11} * 16\text{-bits.}\end{aligned}$$

q7

~~mm-add = 256 kB~~

Cache size = 256 kB

Block size = 32 B

mm-address size = 32 bits

tag directory contains

address tag + 2 valid bits + 1 modified bit + 1 replacement bit

$$\begin{aligned} \rightarrow \text{no. of blocks in cache} &= \frac{256 \text{ kB}}{32 \text{ B}} \\ &= \frac{2^8 \times 2^{10}}{2^5} \end{aligned}$$

Q.13

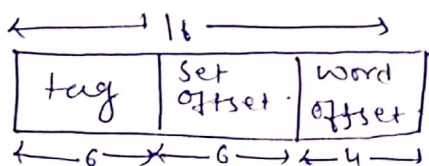
block size = 16 words.

no. of blocks in cache = 256 blocks.

mm blocks = 4096 blocks.

mm size = 4096 * 16 words
= 2^{16} word.

=> mm address = 16-bits.
word offset = $\log_2(2^4) = 4$.



no. of set in cache
= $\frac{\text{no. of blocks}}{\text{associativity}}$

$$= \frac{256}{4} = 2^6$$

Set offset = 6 bits.

∴ Tag bits = 16 - 10 = 6 bits

Q.9

Cache size = 16 kB.

mm-size = $4 \times 2^{30} = 2^{32}$.

address size = 32 bit.

block size = 8 word = $8 \times 4 = 32$ Byte.

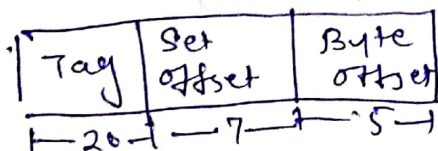
1 word = 32 bit = 4 byte.

∴ blocks in cache = $\frac{2^4 \times 2^{10}}{2^5} = 2^9$.

∴ no. of set in cache = $\frac{\text{no. of blocks}}{\text{associativity}}$

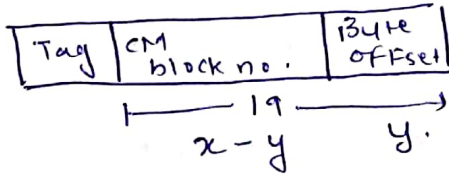
$$= \frac{2^9}{2^2} = 2^7$$

Set offset = 7 bits



q 22 cache size = 2^{19} B.

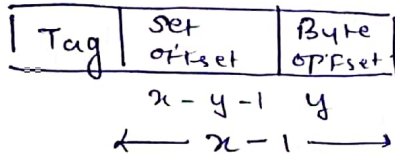
In Direct mapping



cache size $x = \log_2(\text{cache size})$

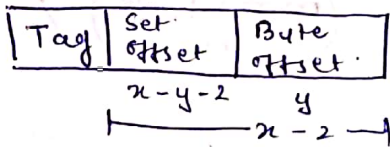
CM block no + Byte offset = $\log_2(\text{cache size})$
always in direct mapping

2-way:



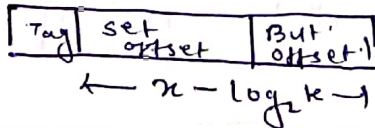
$$\text{no. of sets } s = \frac{2^{x-y}}{2} = 2^{x-y-1}$$

4-way:



$$\text{no. of sets} = \frac{2^{x-y}}{2^2} = 2^{x-y-2}$$

k-way



$$\text{Tag} = \text{mm-add} - [(\log_2(\text{cache size}) - \log_2 k)]$$

here

$$\text{cache size} = 512 \text{ kB} = 2^{19} \text{ B}$$

$$k = 8 \quad \text{mm-add} = 40 \text{ bits}$$

$$\therefore \text{Tag} = 40 - [\log_2(512 \text{ kB}) - \log_2 8]$$

$$= \underline{24 \text{ bits}}$$

Q 25

~~mm address = P bit~~

$$\text{Size of mm} = 2^P \times \cancel{2^{10} \text{ bit}}$$

$$\therefore \text{mm address} = \cancel{P} P \text{ bits.}$$

$$1 \text{ Word} = 2^W \text{ bytes. } k = k.$$

$$\text{Cache size} = \cancel{2^M \times 2^W \text{ bytes}} \cdot 2^N$$

$$\text{Block size} = 2^{W+M}$$

$$\text{Tag} = \text{mm add} - [\log_2(\text{cache size}) - \log_2 k]$$

$$= P - [\log_2(2^N) - \log_2 k]$$

$$= P - N + \log_2 k.$$

In direct mapping

$$\text{CM block no. will be mapping} = \left(\frac{\text{mm block no.}}{\text{no.}} \right) \% \text{ no. of block in cache}$$

In set associativity

$$\text{CM set no.} = \left(\frac{\text{mm block no.}}{\text{no.}} \right) \% \text{ no. of set.}$$

Q 4 Here blocks per set = 2 . i.e 2-way associativity

$$\text{no. of CM blocks} = 2c.$$

$$\therefore \text{no. of set} = \frac{2c}{2} = c$$

$$\text{no. of blocks in mm} = 2cm$$

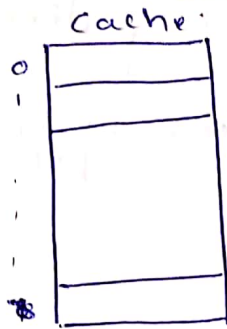
$$\text{CM set no.} = \left(\frac{\text{mm block no.}}{\text{no.}} \right) \% \text{ no. of sets in cache.}$$

$$= \lfloor k \% c \rfloor$$

* Fully associative mapping

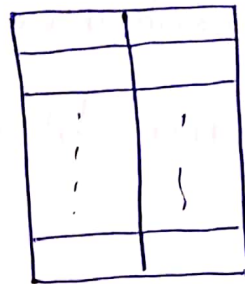
Direct:

Tag	CM block	byte offset
-----	----------	-------------



2way

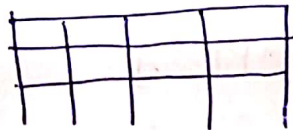
Tag	Set block offset	Byte offset
-----	------------------	-------------



If all the sets or blocks of main memory mapped to a single set in cache then it will be fully associativity.

4way

Tag	Set offset	Byte offset
-----	------------	-------------



∴ For k way k sets will be there.

$$\text{Set offset} = \frac{\text{CM block}}{\text{associativity}}$$

But at the end when k be equal to no. of blocks in cache then it is fully associative.

and set offset will need 0 bit.

i.e there is no set offset part in cache.

Tag	Byte offset
-----	-------------

$$\text{Set offset} = \frac{\text{CM block}}{\text{associativity}} = 1 = \underline{\underline{0 \text{ bit}}}$$

→ A fully associative cache contain 2^{12} blocks
 If main memory contains 2^{16} blocks then
 the size of memory (in kb) required at
 cache controller to store the tag directory

Tag = mm block number.

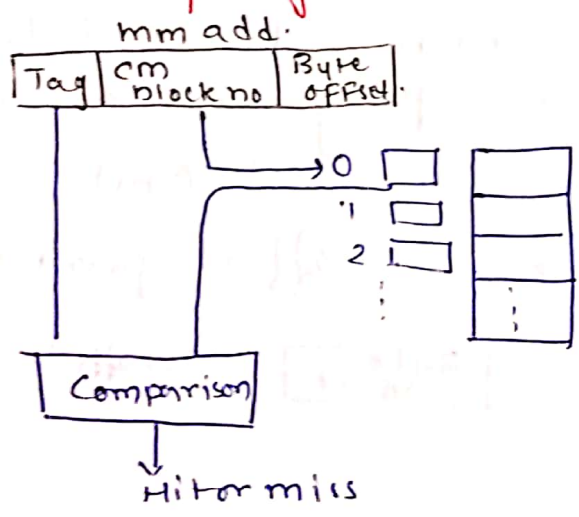
MM block no	Byte offset
-------------	-------------

no. of blocks in mm = 2^{16}
 mm block number = 16-bits.

Tag directory size = $2^{12} * 16$ - bits.
 = 2^{16} - bits
 = 64K bits.

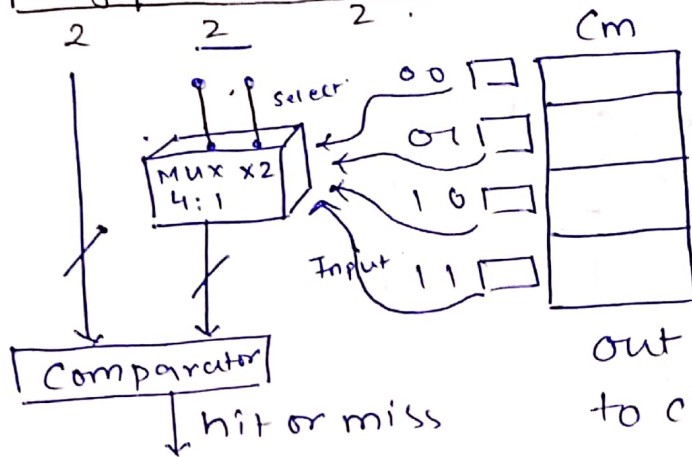
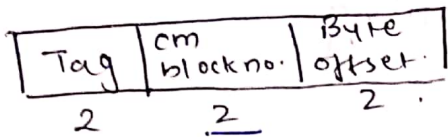
* Hardware implementation of Mapping.

→ Direct mapping:-



Assume; 2 bits

Tag	cm block	Byte offset
2	2	2



There will be 2 similar mux required as mux takes only one bit at a time.

out of 4 tags we need to choose 1

∴ In general.

A mux can forward only 1 bit of tag.

No. of comparators = 1

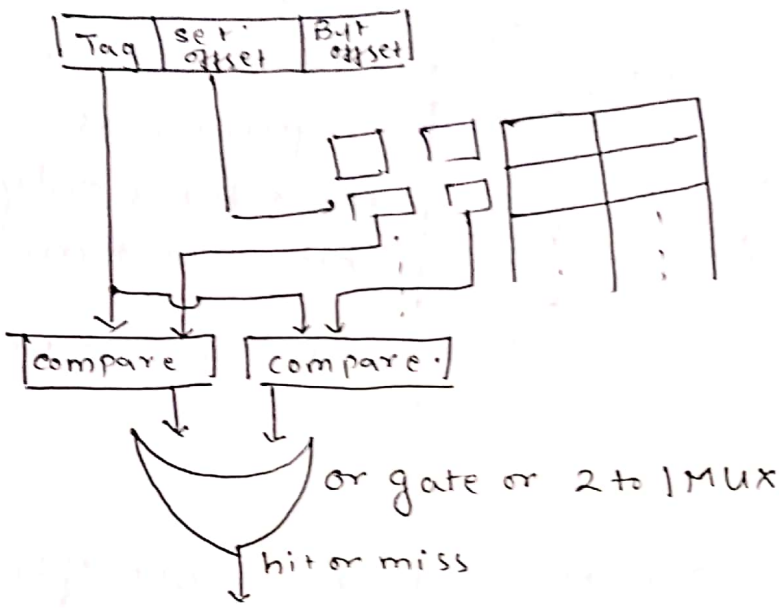
No. of size of comparator = Size of tag in address
(or tag bits in address)

No. of Mux = size of tag in address.
= size of comparator.

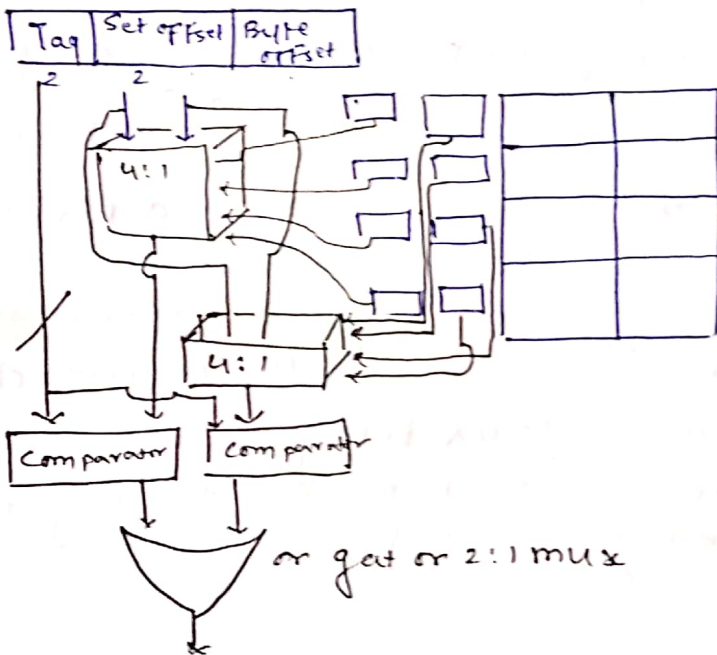
Size of Mux = (No. of blocks in cache : 1)

Hit latency time = MUX delay + Comparator delay
(Time required to declare hit or miss) (for selecting tag) (for tag comparison)

→ In set associative mapping.



eg:



For k -way set associative mapping

No. of comparators = k

Size of comparators = tag-bits

No. of MUX = $k * \text{tag-bits}$

(For tag selection)

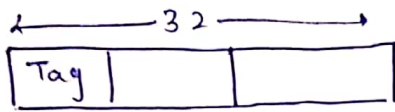
Size of MUX = No. of sets in cache :

Size of MUX = $k + 1$
(For hit or miss)

hit latency time = Mux delay + Comparator delay + Mux-or gate delay (hit or miss)

(For selecting tag)

Q1)



$$32 - 16 = 16 \text{ bits} = (\log_2 256 \text{KB}) - \log_2 u$$

$$\therefore \text{tag bits} = 16.$$

\therefore 4 comparators of size 16 bits each will be needed.

$$k = 4$$

$$\text{no. of comparators} = k = 4$$

$$\text{size of comparator} = \text{tag bits} = 16.$$

Q11

Cache size = 32KB.

2 way

block size = 32B

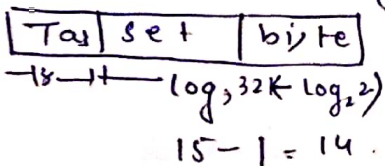
madd = 32-bits.

h_1

2 to 1 Mux \Rightarrow 0.6 ns.

k-bit comparator = $k/10$ nsec.

We need tag size first.



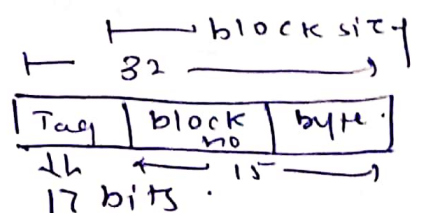
$$15 - 1 = 14.$$

$$\text{no. of blocks in cache} = \frac{32 \text{KB}}{32} = 2^{10}$$

$$\text{no. of sets} = \frac{2^{10}}{2} = 2^9$$

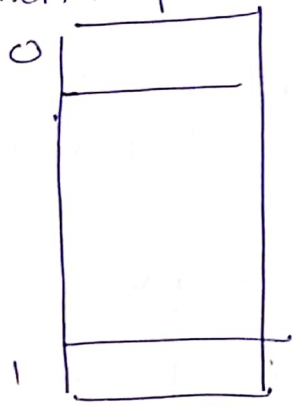
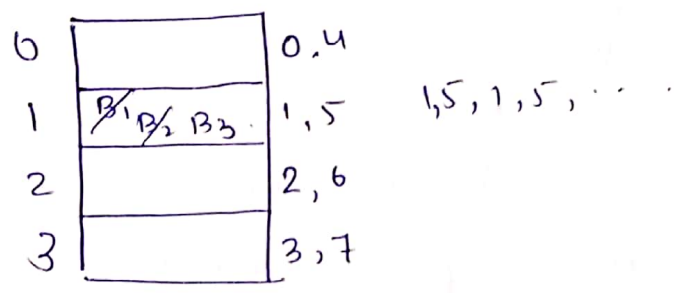
$$h_1 = 0 + \frac{18}{10} + 0.6 = 2.4 \text{ ns.}$$

$$h_2 = 0 + \frac{17}{10} + \dots = 1.7 \text{ ns.}$$

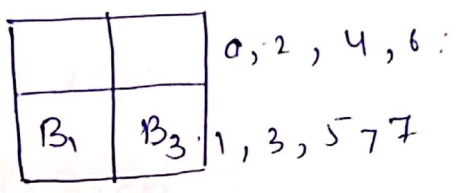


* cache block replacement *

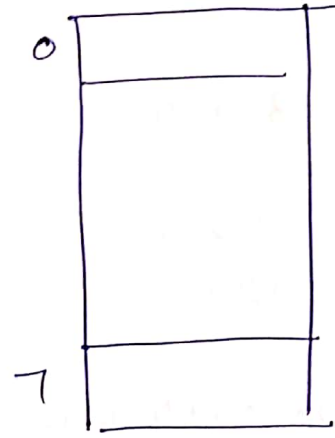
In direct mapping if there is ~~an~~ occurs any replacement there is only one choice to make. Hence no replacement policy needed



in two-way associative



1, 3, 5, 7
 ↓
 To replace any of B_1 or B_3



Replacement policy is used.

- Policies:
- (i) FIFO
 - (ii) Optimal
 - (iii) LRU

consider a direct mapped cache with 8 cache blocks (0 - 7). IF the memory block requests are in the following order:

3, 5, 2, 8, 0, 13, 9, 16, 20, 17, 25, 18, 30, 24, 2, 63, 5, 82, 17, 24.

Which of the following memory blocks will not be at the end of sequence in cache

- a) 3 b) 18 c) 20 d) 30

0	8 0 16 24
1	9 17 25 17
2	2 18 82
3	3
4	20
5	5
6	30
7	63

Q14)

0	48	32	9	216 92
1	1	133	129	73
2				
3	255 155	3	159	63

0, 4, 8, 216, 8, 48, 32.

∴ 216 is not there.

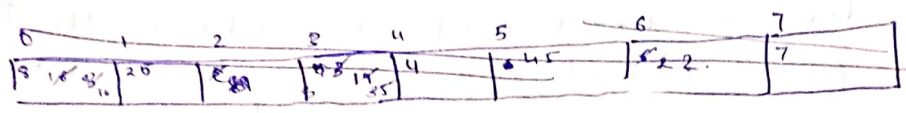
should be moded by k in k associative.

Q17)

0	8 8 12
1	

$\frac{8}{1}, \frac{12}{1}, \frac{8}{1}, \frac{12}{1}$
match

Q 16



0	1	2	3	4	5	6	7
4	3	25	8	14	6	16	35
45	22	✓	✓	3	7	✓	

Q 20

Byte addressable main memory

mm size = $2^{16}B \Rightarrow$ mm add = 16 bit

direct mapped cache

no. of blocks in cache = 32

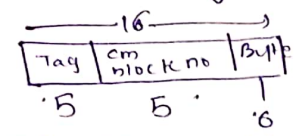
block size = 64B

50×50 2-d array \Rightarrow array size = 2500B

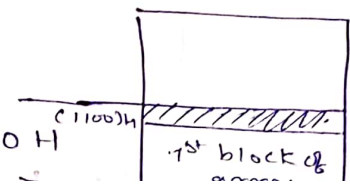
Base address of array = (1100H)

cache

0	24
1	30
2	31
3	22



starting add \Rightarrow 1100H



∴ in 2nd access 8 misses will occur (1 → 8) and upto 32 there will be hits again from (33 → 40) misses will occur

∴ total misses in 2nd access

$$= 8 + 8 = 16,$$

∴ total $40 + 16 = 56$ misses.

and at every time line 4 to line 11 will be replaced

GATE 2006. A CPU has 32KB direct mapped cache with 8 128 byte block size. Suppose A is 2 dimensional array of size 512×512 with elements that occupy 8 bytes each. Consider the following two code segments.

```
P1: For (i=0; i < 512; i++)
    {
        for (j=0; j < 512; j++)
            {
                x += A[i][j];
            }
    }
```

```
P2: For (i=0; i < 512; i++)
    {
        for (j=0; j < 512; j++)
            {
                x += A[j][i];
            }
    }
```

P1 & P2 are executed independently with same initial state. Mainly 'A' is not in cache & i, j, x are in registers. Let the no. of cache misses experienced by P1 be m_1 & that for P2 be m_2 .

(i), the value of m_1 is?

(ii), the value of m_1/m_2 is?

→ Direct mapping.

Cache size = 32 kB \Rightarrow blocksize = 128 B.

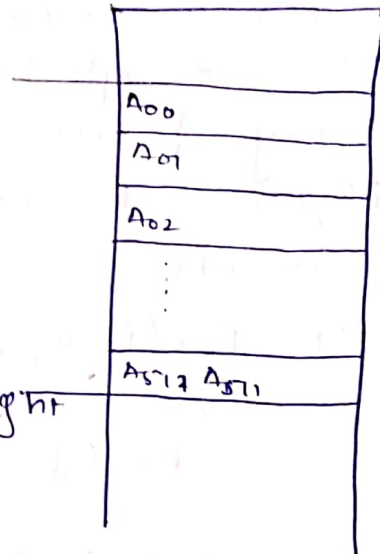
A : 512 \times 512 (each of 8 B)

no. of blocks in ~~array~~ cache = $\frac{32 \text{ kB}}{128 \text{ B}} = 2^4 = 16$

Array size = $(2^9 \times 2^9) \times 2^3 \text{ B} = 2^{21} \text{ B}$.

P₁

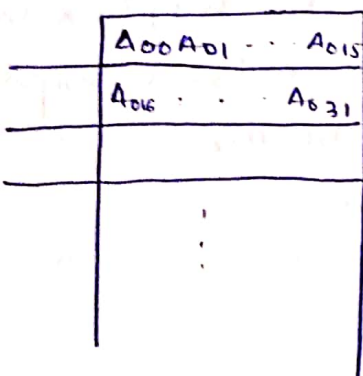
A₀₀ A₀₁ A₀₂ ... A₀₅₁₁
 A₅₁₁₀ ... A₅₁₁₅₁₁



In one block, 128 B can be brought at every miss

\therefore no. of array elements in 1 block = $\frac{128 \text{ B}}{8 \text{ B}} = 2^4 = 16$

\therefore at A₀₀ 1 miss and 15 elements will be no miss



When CPU access A₀₀ a block of 16 element A₀₀ - A₀₁₅ is copied in cache i.e. 1 miss for 16 elements.

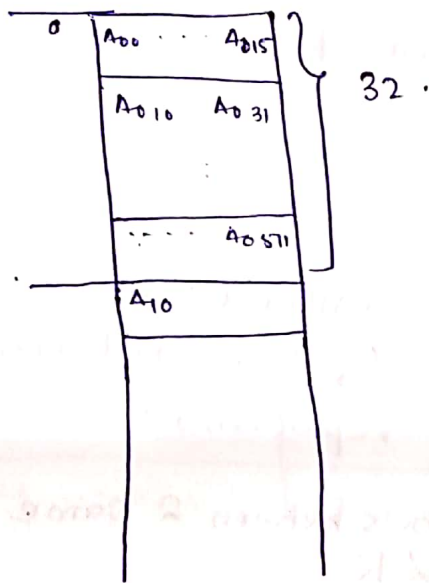
\therefore For Array $2^9 \times 2^9 \Rightarrow 2^{18}$ elements
 no. of misses = $\frac{2^{18}}{2^4} = 2^{14} = 16384$

no. of elements in 1 row = 512.

size of elements = $512 \times 2^3 = 2^{12} B$.

no. of blocks required for 1 row = $\frac{2^{12} B}{2^7} = 2^5 = 32$ blocks

∴ single row will take 32 blocks in cache.
after 32 blocks another row will start.



∴ no. of rows can be stored in cache = $\frac{256}{32} = 2^3 = 8$

∴ $A_{80} - A_{815}$ will replace $A_{00} - A_{015}$ on access.

∴ For each element of array.

⇒ 1 miss

$M_2 = 2^{18}$ misses.

∴ $M_1 / M_2 = \frac{2^{14}}{2^{18}} = 1/16$.

2514 gate → an access sequence of cache block addresses is of length N and contains n unique block addresses the no. of unique block addresses between two consecutive access to same block is bounded above by k . What is the miss ratio if the access sequence is passed through a cache of associativity $A \geq k$ exercising LRU replacement policy,
(A) n/N (B) $1/N$ (C) $1/A$ (D) k/n

Sequence of access = $\bullet N$

i.e. 'N' times CPU is accessing cache.
and \bullet out of them 'n' blocks are only
unique remaining are repeated.

1 3 4 2 3 4 2 3 4 2 1.
unique block!

no. of unique blocks betⁿ two consecutive
same block is $< k$.

\therefore unique block between two consecutive same
block should be maximum $k-1$

associativity $A \geq k$

if same $k=4 \Rightarrow A=4$

no. of misses that can be encountered
= $n +$ extra misses of repeated access
because of replacement.

as $A \geq k$ and unique symbols between 2 same
symbols $\bullet k$

then in A associative blocks.

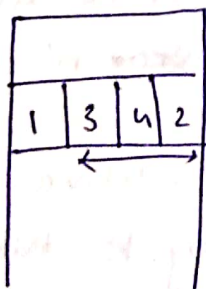
there will always be hits among
the replacement of repeated access.

\therefore extra misses = 0

Hence total miss encountered
will be n .

\therefore miss ratio = $\frac{n}{N}$

only unique no. will be having misses &
everytime the repeated access is there
it will be hit.



→ Cache miss penalty:-

The amount of time CPU spends to service the cache miss is known as cache miss penalty.

which means bringing the missed block from main memory to cache.

assume,

No. of cycles needed to send address to main memory = 1

No. of cycles needed to access main memory = 10

~~4~~ No. of cycles needed to ^{transfer} send 1 cell data to CM = 1

Block size | mm cell size | miss penalty

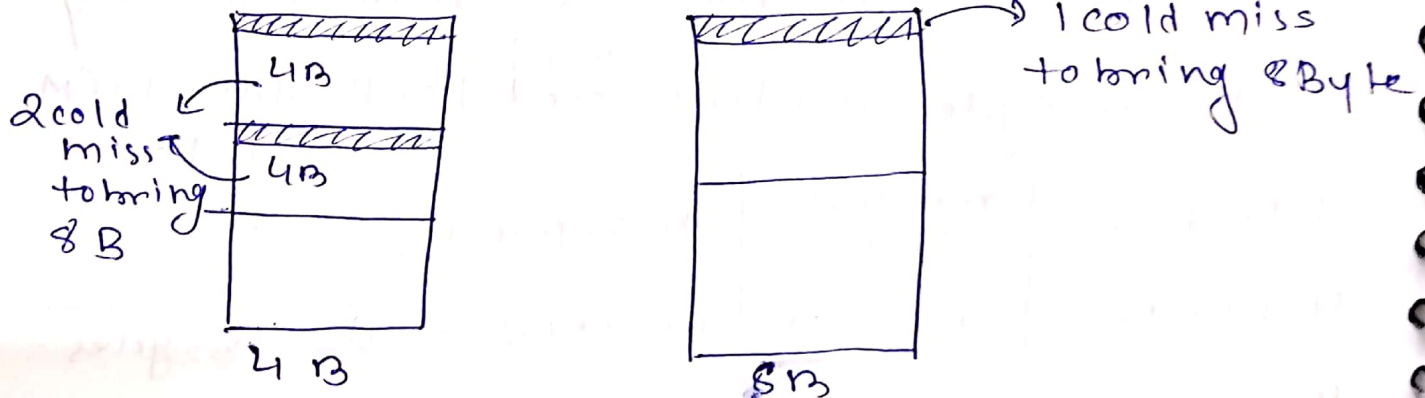
Types of cache miss

→ Cold or compulsory misses

If a cache miss occur due to first time access of any block.

to Reduce compulsory misses, increase the block size.

the more no. of block you can bring at '1' miss the lesser will be the compulsory miss.



→ Capacity miss.

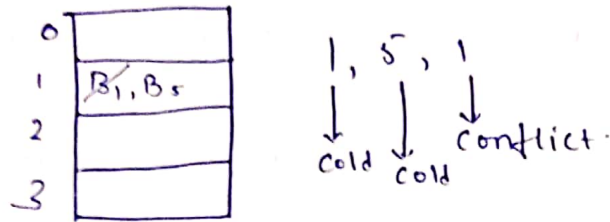
→ If the cache miss occurs because of the capacity of the cache which means when the cache memory is full

→ It should not be a cold or compulsory miss

To reduce capacity misses increase the cache size.

→ Conflict miss

If the cache miss occurs because of the mapping constraint (Tag mismatch).



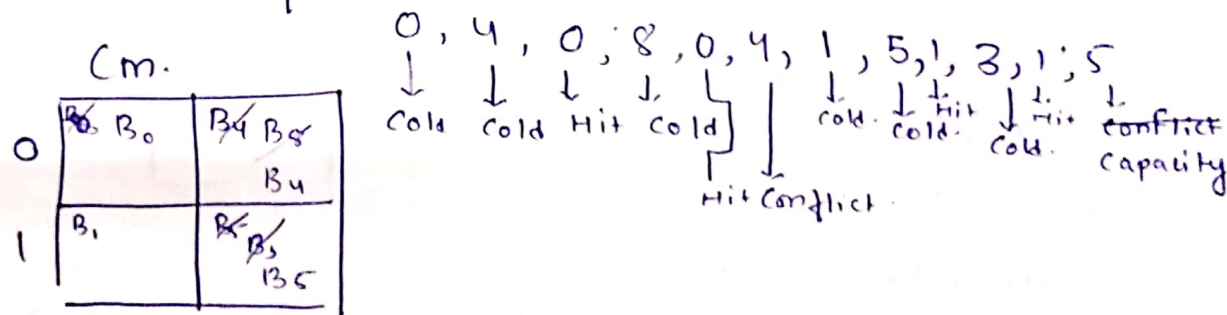
Ex: no. of blocks in cache = 4.

2way set associative

no. of set in cache = $4/2 = 2$.

with LRU replacement.

CPU requests for mm block



Q 23 no of sets = $\frac{256}{2} = 128$

	1 st time	2 nd time		1 st time	2 nd time
0	Cold	Conflict	1	cold	Conflict
128	Cold	Hit	129	cold	Hit
256	Cold	Conflict	257	cold	Conflict
129	Hit	Hit	129	Hit	Hit
0	Conflict	Conflict	1	Conflict	Conflict
128	Hit	Hit	129	Hit	Hit
256	Conflict	Conflict	257	Conflict	Conflict
129	Hit	Hit	129	Hit	Hit

1st time - 4 conflict miss
 2nd - 10th time → 8 conflict misses.
 ∴ Total = 4 + (9 × 8) = 4 + 72 = 76

→ Goal of using cache memory:

1. minimize access time.
2. Maximize hit rate.
3. Minimize miss penalty.

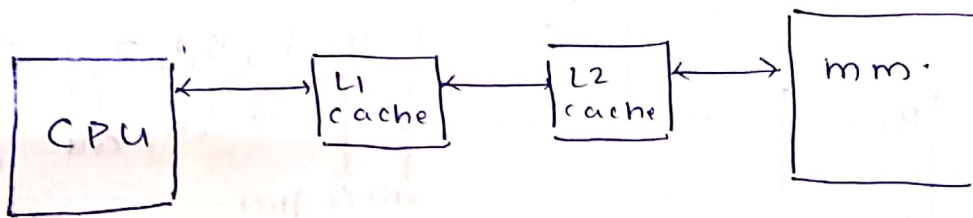
But to have minimum access time we use smaller cache.

And to have maximize hit rate, we need larger cache.

Hence the solution for this conflicts is to use multilevel cache.

to minimize miss penalty use non-block cache, multiport, multibank memories.

→ Multi-level cache:



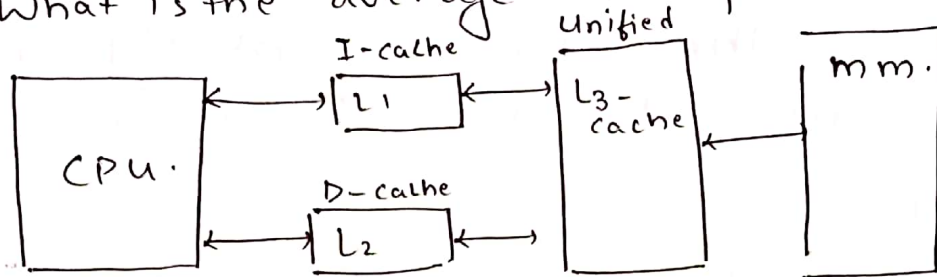
Simultaneous access:

$$t_{avg} = H_1 t_1 + (1 - H_1) [H_2 t_2 + (1 - H_2) t_{mm}]$$

Hierarchical:

$$\begin{aligned} t_{avg} &= H_1 t_1 + (1 - H_1) [H_2 (t_2 + t_1) + (1 - H_2) (t_1 + t_2 + t_{mm})] \\ &= t_1 + (1 - H_1) [t_2 + (1 - H_2) t_{mm}] \end{aligned}$$

question: Consider a multilevel memory Hierarchy as shown below. The hit ratio of L_1, L_2, L_3 and main memory are 0.8, 0.85, 0.9, & 1. Respectively the access times of respective memories are 10 ns, 10 ns, 50 ns & 500 ns. among total memory references 60% are for data access. What is the average memory access time.



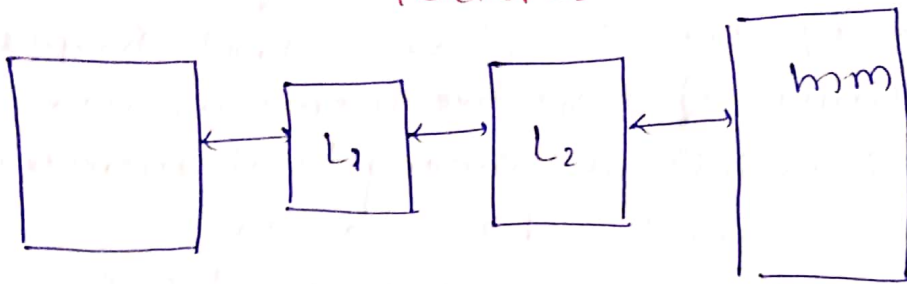
$$\begin{array}{l|l}
 t_1 = 10 \text{ ns} & h_1 = 0.8 \\
 t_2 = 10 \text{ ns} & h_2 = 0.85 \\
 t_3 = 50 \text{ ns} & h_3 = 0.9 \\
 t_{mm} = 500 \text{ ns} & h_{mm} = 1
 \end{array}$$

$$\begin{aligned}
 t_{\text{avg inst}} &= t_1 + (1-h_1) * [t_3 + (1-h_3) * t_{mm}] \\
 &= 10 + 0.2 * [50 + 0.1 * 500] \\
 &= 30 \text{ ns}
 \end{aligned}$$

$$\begin{aligned}
 t_{\text{avg data}} &= t_2 + (1-h_2) * [t_3 + (1-h_3) * t_{mm}] \\
 &= 10 + 0.15 * [50 + 0.1 * 500] \\
 &= 25 \text{ ns}
 \end{aligned}$$

$$\begin{aligned}
 \therefore t_{\text{avg}} &= (0.6 * 25) + 0.4 * 30 \text{ ns} \\
 &= 15 + 12 \\
 &= 27 \text{ nsec}
 \end{aligned}$$

→ Inclusion vs. Exclusion Policies.



~~The entire~~ The entire content of L_1 should be present in L_2 also this is inclusion policy.

- That means if there's a ~~hit~~ miss in L_1 & hit in L_2 then bring the block from L_2 to L_1
- Similarly IF there's a miss in L_1 & L_2 both then the block will be taken from mm to L_2 & then from L_2 to L_1

Exclusion policy :-

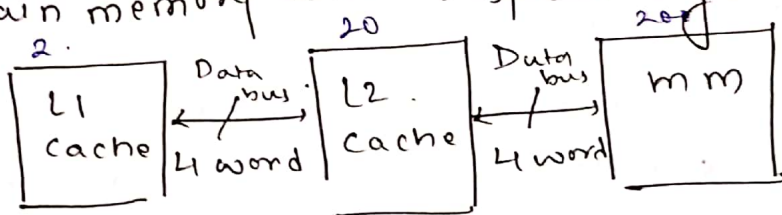
The entire content of L_1 should not be present in L_2 necessarily this is exclusion policy.

- IF miss in L_1 and hit in L_2 bring blocks from L_1 to L_2 .
- IF miss in L_1 & also in L_2 then bring the block from main memory to L_1 & ~~send the~~ copy replaced block of L_1 into L_2 .

The replaced block is called victim block and ~~vict~~ L_2 cache is known as victim cache. as L_2 gets populated by L_1 's victim blocks.

There might be difference in block size of both L_1 & L_2 then the amount of data copied from L_2 to L_1 will be according to L_1 block size.

A computer system has an L_1 cache & an L_2 cache and a main memory unit connected as shown. The block size in L_1 is 4 words the block size in L_2 is 16 words the memory access times are 2ns; 20ns; 200ns. For L_1 cache, L_2 cache & main memory unit respectively.



- 9.ii) When there is a miss in L_1 cache & Hit in L_2 cache a block is transferred from L_2 to L_1 cache. What is the time taken for this transfer.
- ii) When there is a miss in both L_1 & L_2 cache first a block is transferred from main memory to L_2 and then a block is transferred from L_2 to L_1 cache. What is the total time taken for these transfers.

Data bus transfers 4 words at time.

1) Here only transfer time is asked.

4 word in L_1
16 in L_2 .

∴ 4 words block is required
will be directly transferred
to L_2 .

$$\begin{aligned} \text{reading from } L_2 &= 20 \text{ ns} \\ \text{writing to } L_1 &= \frac{2 \text{ ns}}{2} \end{aligned}$$

(ii) any no. of words can be accessed from mm.

$$\therefore \text{from block from mm to } L_2 \Rightarrow 4 * [200 + 20] = 880 \text{ nsec.}$$

$$\text{from } L_2 \text{ to } L_1 \Rightarrow 22 \text{ ns.}$$

$$\therefore \underline{880 + 22 \text{ ns} = 902 \text{ ns.}}$$

→ Secondary Memory:-

↳ Magnetic Disk:-

Here storage is known as recording, and every platters has both surface as recording surface.

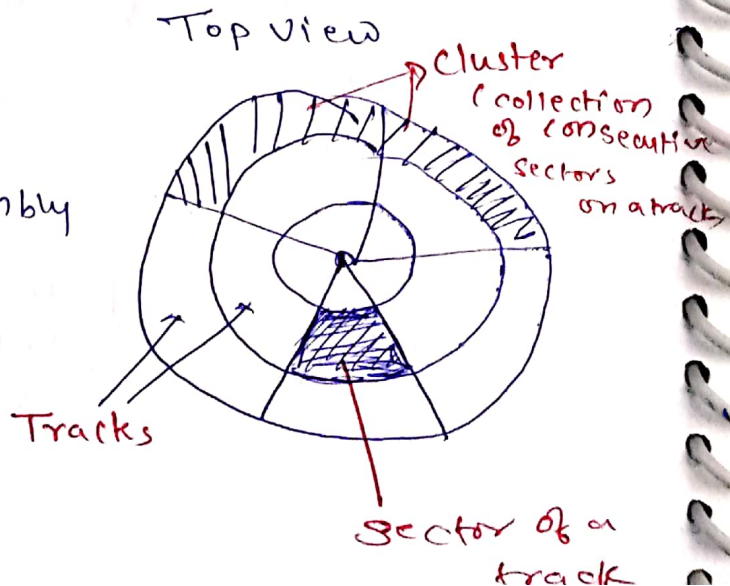
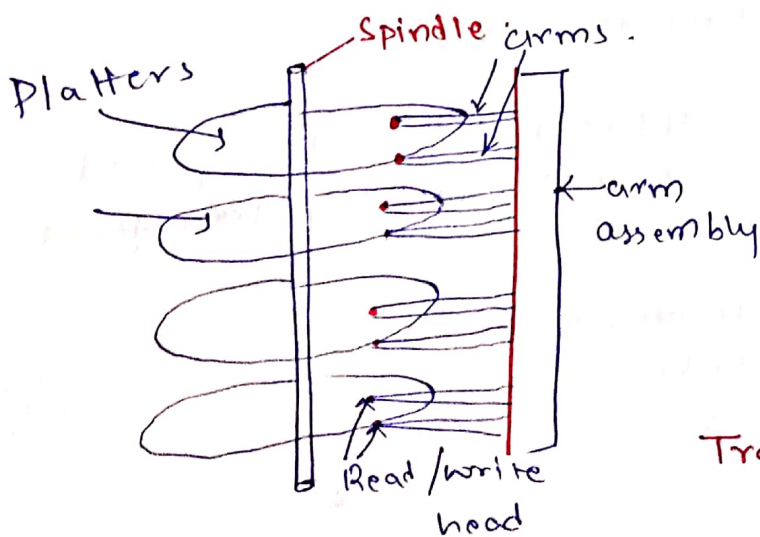
For every surface will require read/write heads which points at the surfaces.

They are attached through arm assembly

every disk will be divided into concentric tracks and these tracks will be divided into sectors

Here sector of a track is Sector within particular track.

Collection of consecutive sectors is cluster.



Storage.

used in modern days.

constant sector capacity or variable storage density

variable sector capacity. or constant storage density.

Disk capacity = $2 * \text{no. of plater in disk} * \text{no. of tracks per surface} * \text{no. of sector per track} * 1 \text{ sector capacity.}$
For constant sector capacity.

Here smallest unit of disk is sector. which can be read or written at once.

hence address will be generated for sectors i.e every sector has one address

→ Disk access time:

In one access time is to access 1 sector.

$\therefore \text{access time} = \text{Seek time} + \text{rotational latency} + 1 \text{ sector transfer time}$
(header should reach track on which sector is present) (to reach to the sector i.e the r/w header reaches to sector) + extra time.

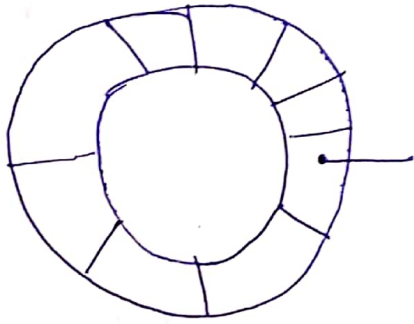
Seek time: Time required to move the arm over desired track.

rotational latency = Time required to rotate the desired sector under read/write head.

1 sector transfer time = Time required to read or write one sector

average rotational latency = 1 disk rotational time

Assume:- no. of sectors per track = 10
 1 rotational time = 100 msec.



If specific position of r/w head given calculate specific rotational latency otherwise use average rotational latency.

In one rotation 1 track can be transferred.

$$\therefore \text{1 sectors time} = \frac{\text{1 rotation time}}{\text{no. of sectors per track.}}$$

question consider a disk with 16 platter. each with 2 recording surfaces; 1k tracks per surface 2k sectors per track and 2048 B per sector. disk rotates with 3000 rpm, seek time of disk is 10 ms

- 1) find capacity of disk in GB
- 2) find the no. of bits required for disk addressing
- 3) find the disk access time.
- 4) what is disk transfer rate.

$$1) \quad 16 \times 2 \times 2^{10} \times 2^{10} \times 2^8 \times 2^{10} \\
 2^4 \times 2^4 \times 2^{10} \times 2^{10} \times 2^8 \times 2^{10} \\
 2^7 \times 2^{40} = 128 \text{ GB.}$$

$$(i) \text{ disk capacity} = 2 * 16 * 1k * 2k * 2048B$$

$$= 2^{37}B = 128GB$$

$$(ii) \text{ no. of sectors} = 2 * 16 * 1k * 2k$$

$$= 2^{26}$$

$$\text{address} = 26 \text{ bits}$$

$$\text{disk access time} = \text{Seek time} + \text{rotational latency} + \text{1 sector transfer time}$$

$$= 10 \text{ ms} + \frac{20}{2} + \frac{20 \text{ ms}}{2k}$$

$$\text{rotational speed} = 3000 \text{ rotation}$$

$$\therefore \text{for } 3000 \text{ rotation} \rightarrow 60 \text{ sec}$$

$$1 \text{ rotation} = \frac{60 \text{ sec}}{3000}$$

$$= 20 \text{ msec}$$

$$\therefore \text{disk access time} = 10 + 10 + 0.01 = 20.01 \text{ msec}$$

now 0.01 msec required to transfer $\rightarrow 2048 \text{ Bytes}$

then in 1 ms $\rightarrow 204800 \text{ Bytes}$

~~204800~~ $\rightarrow 200 \text{ kB/msec}$

~~200 kB/msec~~

then in 1 sec $\rightarrow \frac{2048}{0.01 \times 10^{-3} \text{ sec}} = 200 \text{ MB/sec}$

If n no. of sectors to be transferred

1) If all the sectors are stored sequentially

$$\text{Total disk access time} = \text{Seek time} + \text{rotational latency} + (n * \text{1 sector transfer time})$$

2) IF sectors are stored on random locations

$$\text{Total disk access time} = n * [\text{seek time} + \text{rotational latency} + \text{1 sector transfer time}]$$

Q5

$$\Rightarrow 6000 \text{ rotation} \rightarrow 60 \text{ sec.}$$

$$\therefore \text{rotation latency} = 10 \text{ msec}$$

$$= 100 * \left[10 \text{ msec} + \frac{10}{2} + 0 \right]$$

$$= 100 * 15 \text{ msec.}$$

$$= 1500 \text{ msec}$$

$$= 1.5 \text{ sec.}$$

Q7

$$1 \text{ rotation time} = \frac{60 \times 1000 \text{ ms}}{10000} = 6 \text{ ms}$$

$$\text{total file transfer time} = 2000 * \left[4 + \frac{6}{2} + \frac{6}{600} \right]$$

$$= 2000 * [4 + 3 + 0.01]$$

$$= 2000 * 7.01 \text{ ms.}$$

$$= 14020 \text{ msec}$$

Q8

$$1 \text{ rotation time} = \frac{60 \times 1000}{15000} = 4 \text{ ms.}$$

$$\text{Avg rotation latency} = 2 \text{ ms.}$$

$$\text{seek time} = 4 \text{ ms.} = 2 * 2 \text{ ms.}$$

$$\text{Disk To transfer } 50 \times 10^6 \text{ B} \rightarrow 1 \text{ sec.}$$

$$512 \text{ B} \rightarrow \frac{1}{50 \times 10^6} * 512.$$

$$1 \text{ sector transfer time } 0.01 \text{ msec.}$$

$$\text{controller transfer time} = 10 * 0.01 = 0.1 \text{ ms}$$

$$\therefore 1 \text{ disk access time} = \text{seek} + \text{rotation latency} + 1 \text{ sector transfer time} + \text{controller.}$$

1260th <1, 0, 0>

<1, 19, 62>

no. of sectors in 400 cylinders (0-399) = 400×1260
= 504000

no. of sector in 16 surface = $16 \times 63 = 1008$
+ 29.

$$\begin{array}{r} \Rightarrow 504000 \\ + 1008 \\ + 29 \\ \hline \boxed{505037} \end{array}$$

no. of sectors per cylinder = n_c .

no. of sector per track = n_t .

for address $\Rightarrow (c, h, s)$

$$\text{Sector no.} = (c * n_c) + (h * n_t) + s.$$

$$c = (\text{Sector no.}) / n_c.$$

$$h = (\text{Sector no.} \% n_c) / n_t.$$

$$s = (\text{Sector no.} \% n_c) \% n_t.$$

$$\therefore 1039 \quad n_c = 1260 \quad n_t = 63$$

$$1039 / 1260 \Rightarrow 0$$

$$(1039 \% 1260) / 63 \Rightarrow 16$$

$$(1039 \% 1260) \% 63 \Rightarrow 31$$

Q6

no. of surfaces = 16 (0-15)

no. of cylinders = 16384 (0-16383)

no. of sectors per track = 64.

Sector Capacity = 512B

File 42797 KB

Starting add = <1200, 9, 40>

Sector no. of 1st sector = (1200 * 1024) + 9 * 64 + 40.

no. of sector per cylinders = 16 * 64 = 1024

Sector no. of 1st sector = 1229416.

no. of sector required to store the file.

$$= \frac{42797 \text{ KB}}{512} = 85594 \text{ sectors.}$$

Sector no. of last sector of file =

$$\begin{array}{r}
 1229416 \\
 + 85594 \\
 \hline
 1315009
 \end{array}$$

C = 1315009 / 1024 = 1284.

H = (1315009 % 1024) / 64 = 3.

S = (1315009 % 1024) % 64 = 1.

add = <1284, 3, 1>

→ Parallel Processing

technique for simultaneous data processing

types.

↳ Vector processing

↳ Array processing

↳ Pipeline "

→ Flynn's classification of computers

1) SISD :- Single Instruction stream
Single Data stream
at a time single Instruction is fetched
at a time single Instruction is Executed.

2) SIMD: Single Instruction stream Multiple
data stream
at a time single Instruction is fetched but
multiple Instruction is Executed.
eg: Pipeline.

3) ~~MISD~~ MISD: Multiple Instruction stream
Single data stream.
at a time multiple Instruction is fetched but
single Instruction is executed.

4) MIMD: Multiple Instruction stream multiple
data stream

multiple Instructions are fetched &
multiple Instructions are executed.

eg: Super scalar system, which

IIP — Instruction level parallelism.

Superscalar systems uses multiple pipelines

⇒ Pipelining:-

→ The technique to divide an operation into sub-operations so that the suboperations can be performed in parallel over different inputs.

→ Pipelining can be implemented only when the same processing should be applied over multiple inputs.

→ In pipelining entire processing is divided into multiple suboperations

→ Each suboperation is performed in a separate unit known as segment or pipeline stage.

→ all the segments can perform their respective suboperations in parallel on different different inputs.

→ entire processing performed on an input is known as a task.

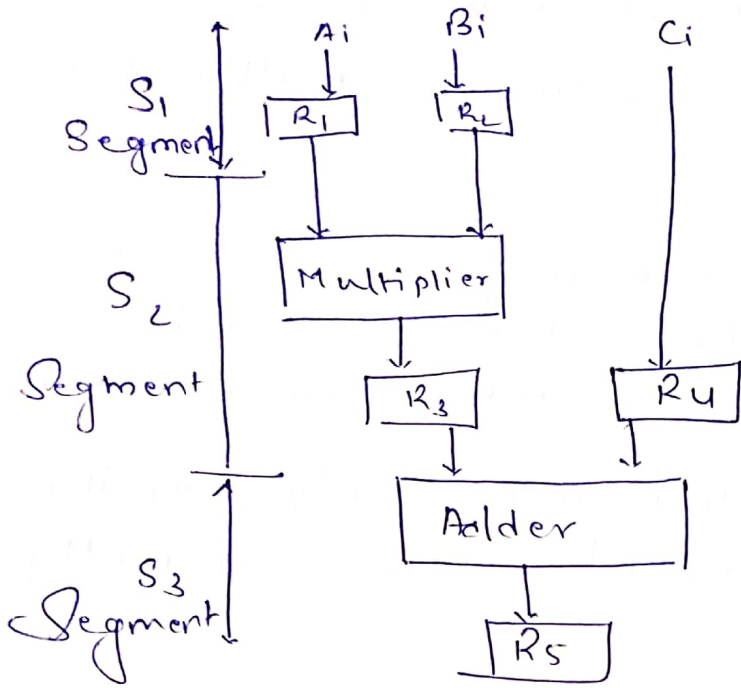
Consider an example.

$$A_i * B_i + C_i \quad \text{For } i = 1 \text{ to } 5$$

$$\text{Seg 1: } R_1 \leftarrow A_i ; R_2 \leftarrow B_i$$

$$\text{Seg 2: } R_3 \leftarrow R_1 * R_2 , R_4 \leftarrow C_i$$

$$\text{Seg 3: } R_5 \leftarrow R_3 + R_4$$



Segment clock cycle	Seg 1		Seg 2		Seg 3
	R ₁	R ₂	R ₃	R ₄	R ₅
1	A ₁	B ₁			
2	A ₂	B ₂	A ₁ * B ₁	C ₁	
3	A ₃	B ₃	A ₂ * B ₂	C ₂	A ₁ * B ₁ + C ₁
4	A ₄	B ₄	A ₃ * B ₃	C ₃	A ₂ * B ₂ + C ₂
5	B ₅	B ₅	A ₄ * B ₄	C ₄	A ₃ * B ₃ + C ₃
6	-	-	A ₅ * B ₅	C ₅	A ₄ * B ₄ + C ₄
7	-	-			A ₅ * B ₅ + C ₅

inputs	no. of cycles in non-pipeline	no. of cycles in pipeline
n = 5	3 * 5 = 15	7
n = 6	18	8
n = 7	21	9

Pipeline cycle time.

The amount of time in which all the segments can perform their respective suboperations.

General consideration about pipeline

Consider a k -segment pipeline with pipeline cycle time t_p . To perform 'n' no. of tasks.

Time required to perform first task equal to
 $= k * t_p$

Time required to perform remaining $(n-1)$ tasks =
 $= (n-1) t_p$.

∴ Time required to perform n tasks = $(k+n-1) * t_p$

↓
Total no. of cycles cycle time

→ Consider a non pipeline system which takes t_n time to perform 1 task

∴ Time required to perform n task in non pipeline system = $n * t_n$.

Performance of pipeline is given by speed up ratio: is $\frac{\text{Slower technique}}{\text{faster technique}}$

$$\therefore \text{Speed up} = \frac{\text{Non-pipeline time (s)}}{\text{Pipeline time}}$$

$$S = \frac{n * t_n}{(k+n-1) t_p}$$

as the no. of task increases $n \gg k$ such that $n \gg k$.

\therefore Time required to perform n tasks = " $n * t_p$ "
as $k-1$ can be ignored

$\therefore S_{max} = \frac{t_n}{t_p}$ S_{max} is also known as Speedup.

In ideal condition the initial $k-1$ task or $k-1$ cycles are ignored as this much time is negligible at higher rate.

If non pipeline & pipeline system take same time to perform 1 task

$$\therefore t_n = k * t_p$$

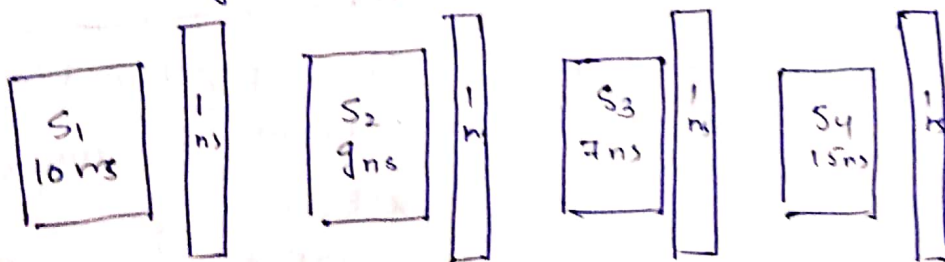
~~$$t_n = k * t_p$$~~

~~$$\therefore S_{max} = k$$~~

$$t_n = k * t_p$$

$$\boxed{S_{max} = k}$$

In pipeline system we use intermediate Buffer register to store output of each segment and to synchronise.



~~$$t_p = \max(\text{segment delay})$$~~

$$t_p = \max(\text{segment delay}) + \text{Reg. delay}$$

$$= \max(10, 9, 7, 15) + 1$$

$$= 15 + 1 = 16 \text{ ns}$$

Intermediate Reg. or buffer used in pipeline systems only.

To get the better performance, design the pipeline in such a way that all the segment should take almost equal time.

Q1
P30

$$k = 4 \text{ segments}$$

$$n = 1000$$

$$t_p = \max(\text{seg delays}) + \text{Reg. delay} \\ = 160 + 5 = 165 \text{ ns.}$$

$$\therefore \text{Cycle time} = (k+n-1) t_p = (4+1000-1) 165 \text{ ns.} \\ = 165495 \text{ ns.} \\ = 165.5 \text{ } \mu\text{s}$$

note:

Speed up

If value of n given

$$S = \frac{n * t_n}{(k+n-1) t_p}$$

If value of n not given

$$S = \frac{t_n}{t_p}$$

Q20

$$t_n = 12 \text{ cycles}$$

$$k = 6$$

$$t_p = 6$$

$$S = \frac{t_n}{t_p}$$

$$= \frac{12}{6} = 2.$$

D_1
 $k = 5$
 $3, 2, 4, 2, 3$
 $n = 100$
 $t_p = \max(3, 2, 4, 2, 3)$
 $= 4 \text{ ns}$
 $\text{time} = (k+n-1) t_p$
 $= (5+100-1) 4 = 416 \text{ ns}$
 $\therefore \text{time saved} = 416 - 214 = 206 \text{ ns}$

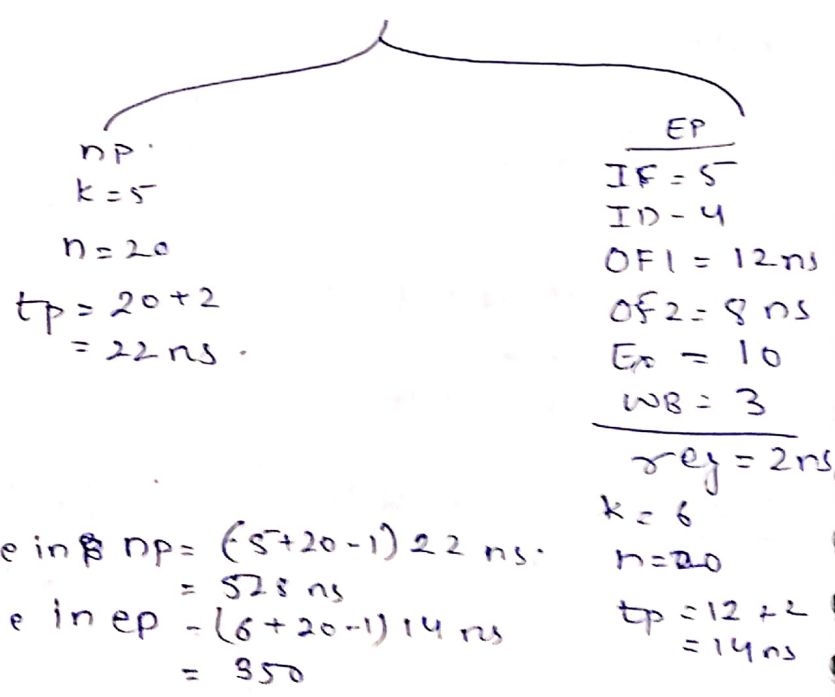
D_2
 $k = 8$
 $2, 2, 2, 2, 2, 2, 2, 2$
 $n = 100$
 $t_p = 2 \text{ ns}$
 $\text{time} = (8+100-1) 2 \text{ ns}$
 $= 214 \text{ ns}$

Q7 In steady stage.

$S_{\max} = \frac{t_n}{t_p}$
 $t_p = \max(5, 6, 11, 8) + 1$
 $= 12 \text{ ns}$
 $t_n = 5 + 6 + 11 + 8$
 $= 30 \text{ ns}$
 $\therefore S_{\max} = \frac{12}{30} = 2.5$

21
 $IF - 5$
 $ID - 4$
 $OF - 20$
 $EX - 10$
 $WB - 3$

 $\tau_{reg} = 2 \text{ ns}$
 $n = 20$
 ~~$t_p = 20$~~



$\therefore \text{time in } NP = (5+20-1) 22 \text{ ns}$
 $= 528 \text{ ns}$
 $\text{time in } EP = (6+20-1) 14 \text{ ns}$
 $= 350$
 $\text{Speed up} = \frac{528}{350} = 1.51$

Q17

$P_1 \quad k = 4$
 $t_p = \max(1, 2, 2, 1)$
 $= 2 \text{ ns}$

P_2
 $t_p = \max(1, 1.5, 1.5, 1.5)$
 $= 1.5$

$P_3 : t_p = 1 \text{ ns}$

$P_4 : t_p = 1.1 \text{ ns}$

$\therefore P_3$ has lowest clock time \Rightarrow highest frequency

Q18

non pipeline.

clock rate = 2.5 GHz

cycle time = 0.4 ns

no. of cycles per inst? = 4

$t_n = 4 \times 0.4 = 1.6 \text{ ns}$

Pipeline.

clock rate = 2 GHz

cycle time = 0.5 ns

$t_p = 0.5 \text{ ns}$

$\therefore \text{Speedup} = S_{\max} = \frac{t_n}{t_p} = \frac{1.6 \text{ ns}}{0.5 \text{ ns}} = 3.2$

In t_p whole one cycle is performed, hence nothing is required to store 'n' or 'k'

\rightarrow Instruction pipeline:

IF the pipeline processing is implemented on instruction cycle then that pipeline is known as instruction pipeline.

Assume,

5-stage inst? pipeline.

IF - Instruction fetch

DA - Decode & Add. calculation

OF - Operand fetch

EX - Execution

WB - Write Back

Cycles	1	2	3	4	5	6	7	8	9	10	11		
Instruction I_1	IF	DA	OF	EX	WB								
I_2		IF	DA	OF	EX	WB							
I_3			IF	DA	OF	EX	WB						
I_4				IF	DA	OF	EX	WB					
I_5					IF	-	-	-					
I_6													
I_7									IF	DA	OF	EX	WB
I_8										IF	DA	OF	-
I_9													

Stall cycles:

In 5th cycle, I_4 is decoded as branch instruction, then I_5 should be discarded from pipeline & CPU will not take any new instⁿ till EX phase of I_4 .

In EX phase of branch instruction (I_4) the condition is evaluated. Based on that in next cycle next instruction is fetched.

Total no. of instructions actually executed = 6.
For $n=6$ instructions:

no. of cycles required without any problem = $5 + 6 - 1 = 10$

extra cycles because of branch problem = 3 extra cycles

\therefore total cycles will be = $10 + 3 = 13$ cycles because of (I_5)

This extra cycles is known as stall cycles which is supposed to be discarded.

the branch result is available after
 the stage then no. of stall cycles because
 branch instruction is equal to "I-1"

I = 5
 I = 7
 D = 10
 I = 8
 D = 6.
 ⇒ 1 ns

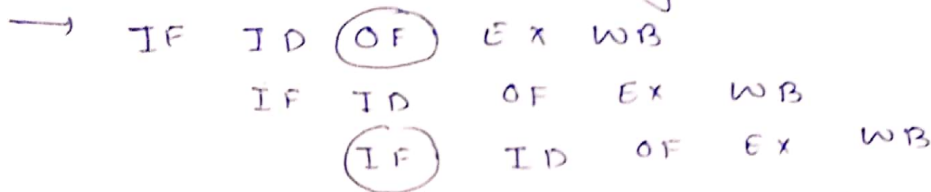
$$tp = 10 + 1 = 11 \text{ ns}$$

1	2	3	4	5	6	7	8	9	10	11	12	13	14
FJ	DI	FO	EI	WD									
	FJ	DI	FO	EI	WB								
		FJ	DI	FO	EI	WD							
			FJ	DI	FO	IE	WD						

→ Pipeline hazards :-

The reason which causes disturbance in smooth conduct of pipeline

1) Resource conflict (collision or structural dependency)
 If the same resource is required by two instructions simultaneously



Here, the 'IF' operation of Instruction I_3 wants to access memory at same time I_1 is using memory to access 'OF'

(ii) Data Dependency:

If the result of an instruction i is used as an input in next instruction

$i : R_1 \leftarrow R_2 + R_3$ ~~IF ID OF EX WB~~
 $i+1 : R_5 \leftarrow R_1 * R_4$

$i \rightarrow$ IF ID OF EX WB .

$i+1 \rightarrow$ IF ID - - OF EX WB

a general pipeline hardware cannot detect data dependency. In this case compiler generates the instructions which solves the data dependency problem.

This solution is known as delayed load

In above example if hardware cannot detect problem the compiler will send set of instructions in such a way that i & $i+1$ don't come one after another

i: $R_1 \leftarrow R_2 + R_3$	IF	ID	OF	EX	WB		
i+1: No operation Inst ⁿ	IF	ID	OF	EX	WB		
i+2: No operation Inst ⁿ	IF	ID	OF	EX	WB		
i+3: $R_5 \leftarrow R_1 * R_4$			IF	ID	OF	EX	WB

Delayed Load is also known as software solⁿ as compiler provides it.

IF pipeline hardware is intelligent enough to detect the data dependency. They might be amongy

- Hardware interlock.
- Operand Forwarding (bypassing)

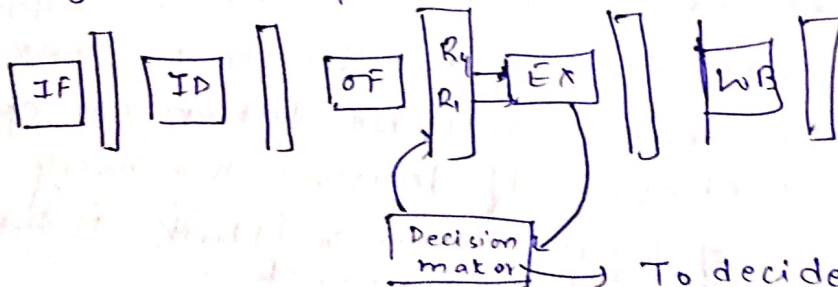
→ **Hardware interlock.**

Hardware will provide stall cycles to delay operand fetch as in last example.

→ **Operand Forwarding (bypassing)**

Avoiding the delay is extra delay. If one instruction is dependent on another then on execution time of I, the result is forwarded to 'WB' phase and also to the Execution phase of next immediate instruction as these two executions are not happening simultaneously.

$R_1 \leftarrow R_2 + R_3$	IF	ID	OF	EX	WB
$R_5 \leftarrow R_1 * R_4$	IF	ID	OF	EX	WB



To decide if dependency then only forward operand.

IF operand forwarding solⁿ is used then no stall cycles because of data dependency

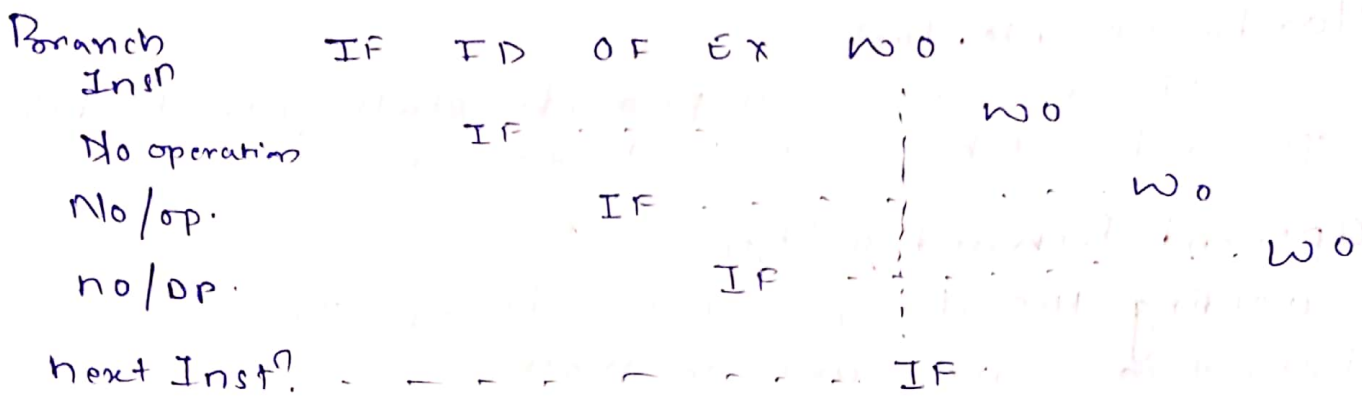
→ Branch difficulty. (Control dependency)
 Because of Branch instⁿ. if Hazard occurs

Normal h/w

Cannot detect branch problem
 Software solⁿ: => delayed branch
 (by compiler)

Intelligent h/w

can detect branch problem.
 H/w solⁿ?



H/w solⁿs => Branch prediction.

IF any Branch instruction is encountered the h/w will take prediction whether to take a branch or to execute the next immediate instruction. IF suppose 80% time the prediction is true then the branch will be taken an operation is continued if decision is wrong i.e.

20% of time then rollback is taken, and executed the branch

=> Prefetch target instruction.

IF any branch instruction is encountered then inside pipeline run 2 pipeline at a time 1 says branch to be taken other where branch is not taken after some cycles the result arrives which says which pipeline to adapt

and which to discard.

Branch target buffer:-



Branch Inst?	Target Inst?
I ₄	I ₉
I ₁₂	I ₂₅

This buffer is used to store the condition of instructions if the instⁿ condition arrives again then we can refer to this buffer.

Loop buffer:-

used for small loops.

Store the loop instruction inside the Buffer which is added to pipeline. and for every loop condition take those instructions from buffer

⇒ Data hazard classification

Consider two instructions i & j and i executes before j

(i) RAW (Read after write)

' j ' is trying to read source before ' i ' writes it. hence ' j ' incorrectly gets old value.

$$i: R_1 \leftarrow R_2 * R_3$$

$$j: R_6 \leftarrow R_1 + R_5$$

Data dependency

Sol: operand forwarding.

(2) WAW (write after write) (False dependency)

j tries to write a destination before i writes it

$$i: R_1 \leftarrow R_2 * R_3$$

$$ii: R_1 \leftarrow R_6 + R_5$$

Sol? Write Dependency.
Reg. Renaming.

$$i: R_1 \leftarrow R_2 * R_3$$

$$ii: R_8 \leftarrow R_6 + R_5$$

(3) WAR (write after read) (False dependency)

j tries to write a destination before i reads it hence i incorrectly gets new value.

$$i: R_1 \leftarrow R_2 * R_3$$

$$j: R_2 \leftarrow R_6 * R_5$$

Anti-dependency.

Sol? Reg. Renaming

$$(i) R_1 \leftarrow R_2 * R_3$$

$$(ii) R_4 \leftarrow R_6 * R_5$$

Q2. 1) Data dependency, 2) Branch conflict, 3) Resource conflict

S₁ → False.

S₂ → Yes

S₃ → No stalls created in Anti dependency

IF — 1 cycle Given Operand forwarding

ID — 1 cycle.

OF — 1 cycle.

ADD or SUB

MUL

DIV

PO ————— 1

3

6

WO — 1 cycle.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
IF	ID	OF	PO	PO	PO	WD						1		1

Q14 no. of inst? = $n = 4$

$$k = 4.$$

$$\text{Without hazard no. of cycles} = k + n - 1 \\ = 4 + 4 - 1 \\ = 7$$

$$\text{extra} = \frac{6.}{}$$

$$\text{Total} = \underline{\underline{13}}$$

Q22. no. of instruction = $n = 100$

$$k = 4.$$

IF - 1.

JD - 1.

OF - 1

PO

$$\begin{array}{ccc} \text{LD} & 35 & 25 \\ \hline 3 & 2 & 1 \end{array}$$

$$\text{Extra} = (40 \times 2) + (35 \times 1) + (25 \times 0) = 80 + 35 = 115$$

WB - 1.

$$\therefore \text{without hazard} = k + n - 1 = 100 + 4 - 1 \\ = 103.$$

$$\text{With extra's stall} = 115 \\ \text{cycles}$$

$$\therefore \text{total} = 103 + 115 = 219.$$

IF k segment pipeline.
 n instⁿ. to execute

$$\text{no. of cycles to execute } n \text{ inst}^n = k + n - 1 \\ \text{CPI} = \frac{k + n - 1}{n} \\ \text{(In normal case)}$$

If in case of hazard, total stall cycles = x .

$$CPI = \frac{(k+n-1) + x}{n}$$

In ideal conditions:

no. of cycles to execute n instⁿ = n .

$$CPI = 1 = n/n.$$

(In normal condition)

In case of hazard

$$CPI = \frac{n+x}{n}$$

Average instruction execution time = $CPI * t_p$

Q 16

$k = 6$. all stage takes ' x ' time, no register delay

$$t_p = x; \quad t_n = 6x.$$

In ideal condition $CPI = 1$

$$CPI_{avg} = 0.75 * 1 + 0.25 * (1+2) = 1.5$$

$$\text{avg inst}^n \text{ execution time} = 1.5x.$$

$$\text{Speedup} = \frac{\text{Time taken in Non pipeline}}{\text{Time taken in pipelined}}$$

$$= \frac{6x}{1.5x} = 4$$

Q 4

$$14 \text{ Hz} = \frac{1}{10^{-10}} = 1 \text{ nsec} \cdot 1 \text{ clock time.}$$

$$k = 5.$$

no. of stalls cycle because of branch = $3 - 1 = 2$.

$$CPI_{avg} = 0.8 * 1 + 0.2 * (1+2) = 1.4$$

$$1 \text{ instr}^n \text{ execution time} = 1.4 * 1 \text{ ns}$$

10^9 instructions execution time.

$$= 10^9 \times 10^{-9} \times 1.4$$
$$= 1.4 \text{ sec.}$$

Q9

$k = 5$

$t_p = 2 \text{ ns.}$

$CPI = 1$

no. of stall cycles = $5 - 1 = 4$.

$CPI_{avg} = 0.8 \times 1 + 0.2(1 + 4)$
 $= 1.8$

1 inst execution time = $1.8 \times t_p$
 $= 3.6 \text{ ns.}$

next instⁿ is Fetched only after completion of branch instⁿ \Rightarrow After 5th stage of branch instⁿ?

Q10

$CPI_{avg} = 0.8 \times 1 + 0.2 \times [0.2 \times (1 + 4) + 0.8 \times [5 \times (1 + 4) + 0.5 \times 1]]$

Annotations:
- non branch (under 0.8 * 1)
- branch (under 0.2 * [0.2 * (1 + 4)])
- unconditional (under 0.2 * (1 + 4))
- Conditional (under 0.8 * [5 * (1 + 4) + 0.5 * 1])
- branch taken (under 5 * (1 + 4))
- branch not taken (under 0.5 * 1)

$0.8 + 0.2 [0.2(5) + 0.8(0.5 * 5 + 0.5)]$

$0.8 + 0.2 [1 + 0.8(2.5 + 0.5)]$

$0.8 + 0.2 (1 + 0.8(3))$

$0.8 + 0.2(3.4) = 1.48 \text{ ~~sec~~ cycles.}$

\therefore 1 instⁿ execution time = $1.48 \times 2 \text{ ns.}$
 $= 2.96 \text{ ns.}$

12

k = 4

	1	2	3	4	5	6	7	8	9	10	11	12	13	14
I ₁	S ₁	S ₂	S ₂	S ₃	S ₄	S ₄								
I ₂		S ₁	S ₁	S ₂	S ₃	S ₃	S ₄							
I ₃			S ₁	S ₂	S ₂	S ₃	S ₃	S ₄						
I ₄				S ₁	S ₁	S ₂	S ₂	S ₃	S ₃	S ₄				
I ₁						S ₁	S ₂	S ₂	S ₃	S ₄	S ₄			

∴ I₁ on second time S₄ will be after 13 ns.

or

I ₁	1	3	4	6					
I ₂		3	$\max(3,3) + 1 = 4$	$\max(4,2) + 1 = 6$	$\max(6,3) + 1 = 7$				
I ₃			4	$4 + 1 = 5$	$\max(5,4) + 1 = 6$	$\max(6,4) + 1 = 7$			
I ₄				$4 + 2 = 6$	$\max(6,5) + 1 = 7$	$\max(7,5) + 1 = 8$			
I ₁					$6 + 1 = 7$	$\max(7,7) + 1 = 8$	$\max(8,7) + 1 = 9$	$\max(10,11) + 1 = 11$	$\max(11,11) + 2 = 13$

Q11

if 100% efficiency $\Rightarrow S = S_{max}$.

88% efficiency \Rightarrow Speedup = 88% S_{max}

100 MHz.

cycle time = $\frac{1}{100 \text{ MHz}} = 10 \text{ nsec.}$

$t_p = 10 \text{ ns.}$

k - segments

$t_h = k * 10 \text{ ns.}$

$S_{max} = k$

$\therefore S = 88\% \text{ of } S_{max}$

$6.6 = 88\% \text{ of } k$

$k = \frac{6.6}{0.88} = 7.5 = \boxed{8}$

In ideal conditions 1 results in 1 cycle,

hence in t_p time

$$\text{no. of results per unit time} = 1/t_p$$

= throughput of pipeline

throughput = no. of operations performed in unit time.

Q5 $k = 4$
800, 500, 400, 300
 $t_p = 800 \text{ ps}$

$$\text{throughput} = \frac{1}{800}$$

$k = 5$
600, 350, 500, 400, 300
 $t_p = 600 \text{ ps}$

$$= \frac{1}{600}$$

% of throughput increase

$$= \frac{\frac{1}{600} - \frac{1}{800}}{\frac{1}{800}} * 100\%$$

$$= (4/3 - 1) * 100\% = 33.3\%$$

⇒ Non-Linear Pipeline :-

Linear

S ₁	x			
S ₂		x		
S ₃			x	
S ₄				x

non-linear

S ₁	x		x	
S ₂		x		x
S ₃			x	
S ₄				

Initiation: Providing a new input in pipeline is known as initiation

Collision: resource conflict.

Latency: difference of no. of cycles b/w 2 initiation

Forbidden latency: latency which causes collision

Permissible latency: latency which does not cause collision.

1	2	3	4	5	6	7	8	
S ₁	S ₂	S ₃	S ₃	S ₂				} - latency = 1 (Permissible)
	S ₁	S ₂	S ₁	S ₃	S ₂			

S ₁	S ₂	S ₁	S ₃	S ₂] - latency = 3 (Forbidden)
			S ₁				↳ collision	

Problem in non-linear pipeline:-

To find a ^{minimum} latency cycle which does not ~~come~~ cause any collision.

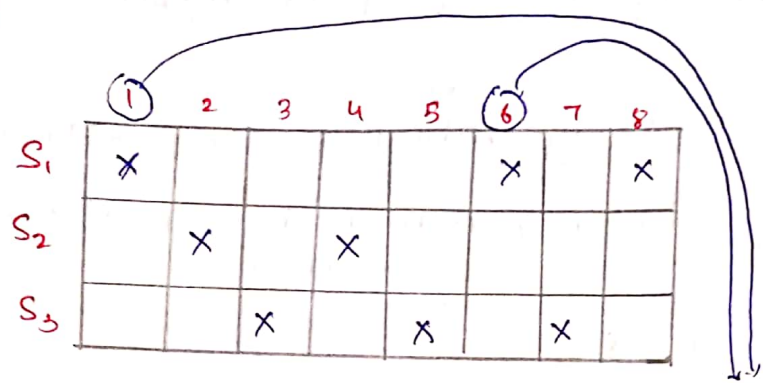
S₁ S₂ S₁ S₃ S₂ :

S₁ S₂ S₁ S₃ S₂ :

If we see above situation next initiation is occurring after whole instruction is executed. this situation is sequential set not pipeline hence to improve it we need minimum latency cycle which provide better outcome in lesser time without any collision.

Eg: Reservation table.

eg.



Forbidden latency = $(6-1), (8-1), (8-6), (4-2)$
 $(5-3), (7-3), (7-5)$
 $= 5, 7, 2, 2, 2, 4, 2$
 $= 7, 5, 4, 2$

Permissible latency = 8, 6, 3, 1

* Collision vector \rightarrow 0 - Permissible
 1 - Forbidden

$c_8 \ c_7 \ c_6 \ c_5 \ c_4 \ c_3 \ c_2 \ c_1$
 0 1 0 1 1 0 1 0 \rightarrow collision vector.

State transition Diagram

Take original collision vector as a state.
 For each 0 at position i in collision vector, right shift it '1' position and take OR with original collision vector; which can give you new state.

* For state 01011010

(i) For 0 at position 1

$$\begin{array}{r} 00101101 \\ (or) 01011010 \\ \hline 01111111 \end{array}$$

(ii) For 0 at position 3:

$$\begin{array}{r} 00001011 \\ 01011010 \\ \hline 01011011 \end{array}$$

(iii) For 0 at position 6:

$$\begin{array}{r} 00000001 \\ 01011010 \\ \hline 01011011 \end{array}$$

(iv) For 0 at position 8 :-

$$\begin{array}{r} 00000000 \\ 01011010 \\ \hline 01011010 \end{array}$$

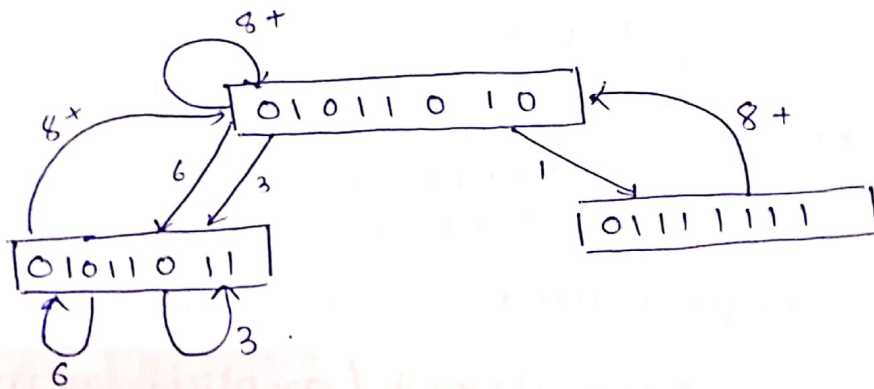
For state 01011011

(i) For 0 at position 3

$$\begin{array}{r} 00001011 \\ 01011010 \\ \hline 01011011 \end{array}$$

(ii) For 0 at position 6

$$\begin{array}{r} 00000001 \\ 01011010 \\ \hline 01011011 \end{array}$$



State diagram

→ **Simple cycle** :- A latency cycle in which every state appears exactly once.

$\{8\}, \{1, 8\}, \{3, 8\}, \{6, 8\}, \{3\}, \{6\},$

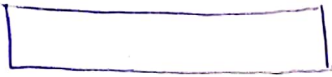
→ **Greedy cycle** :- A simple cycle in which only minimum latency from each state is considered.

$\{1, 8\}, \{3\}$

average latency $\Rightarrow \frac{1+8}{2} = 4.5, 3$

MAL (minimum average latency) = 3.

* Floating point representation *




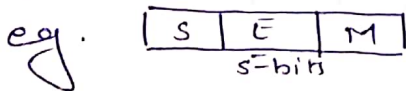
Motivation:- It can represent a larger range of numbers as compared to the fixed point representation using same no. of bits.

In floating point representation every no. is represented as follows.



S → sign ⇒ $\begin{cases} 0 \text{ - +ive} \\ 1 \text{ - -ive} \end{cases}$
 E → Exponent
 M → Mantissa.

- Exponent is represented in biased form.
- Mantissa is a normalized (explicit or implicit) fraction number
-  Bias form → storing only unsigned value of exponent or storing exponent in unsigned value.



S-bit ⇒ Range ⇒ (-16 to +15) ⇒ 0 to 31

original exponent (e)	stored Exponent (E)
-16	0
-15	1
-14	2
⋮	⋮
0	16
⋮	⋮
15	31

(excess -16)

$$E = e + 16 \rightarrow \text{bias}$$

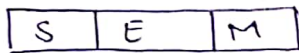
if exponent is represented by k-bits
 bias = 2^{k-1}

Normalization $\left\{ \begin{array}{l} \text{Explicit (convert no. into } 0.1\dots) \\ \text{Implicit (convert no. into } 1.\dots) \end{array} \right.$

100.01 $\left\{ \begin{array}{l} \text{Explicit normalization} \Rightarrow 0.\overbrace{10001}^M * 2^3 \\ \text{Implicit normalization} \Rightarrow \overbrace{1.0001}^M * 2^2 \end{array} \right.$
 ↳ compulsory '1' after
 ↳ compulsory '1' before

In explicit $\rightarrow M = 10001, e = 3$

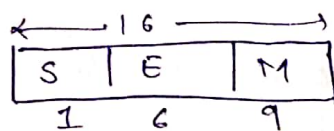
In implicit $\rightarrow M = 0001, e = 2$



Value = $(-1)^S * 0.M * 2^{E-\text{bias}}$
 (Explicit)

Value
 (Implicit) = $(-1)^S * 1.M * 2^{E-\text{bias}}$

\rightarrow Consider a 16-bit register which is used to represent floating point numbers. Exponent is represented in Excess-32 form. Mantissa is a normalized number. What is the 16-bit pattern for the value $-(13.25)_{10}$



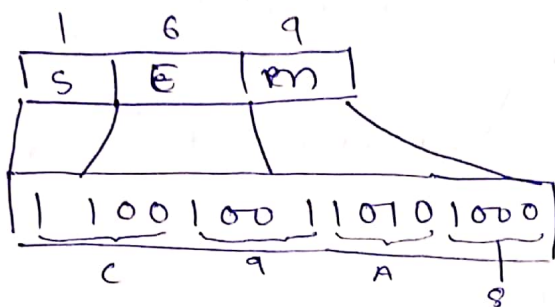
bias = $32 = 2^{k-1}$
 $k = 6$

$-(13.25)_{10} \Rightarrow 13.25 = 1101.01 \Rightarrow$ Explicit normalization
 ↳ (Default)
 $0.110101 * 2^4$

$$M = 110101 \Rightarrow 110101000$$

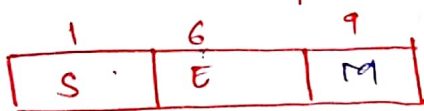
$$e = 4$$

$$E = 4 + 32 = 36 \Rightarrow 100100$$



C9A8

→ what is the maximum positive value which can be represented in this register.



$$S = 0$$

$$M = 111111111$$

$$E = 111111 = 63$$

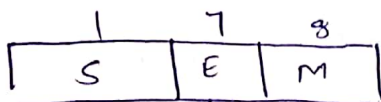
$$\therefore \text{value} = (-1)^0 * 0.111111111 * 2^{63-32}$$

$$= 111111111 * 2^{-9} * 2^{31}$$

$$= 111111111 * 2^{22}$$

$$= (2^9 - 1) * 2^{22}$$

Q4



↳ Excess 64

pure fraction — no normalization

$$0.239$$

~~2 * 2^-1~~ ~~2 * 2^-1~~

$$0.239 * 2 = 0.478 \rightarrow 0$$

$$0.478 * 2 = 0.956 \rightarrow 0$$

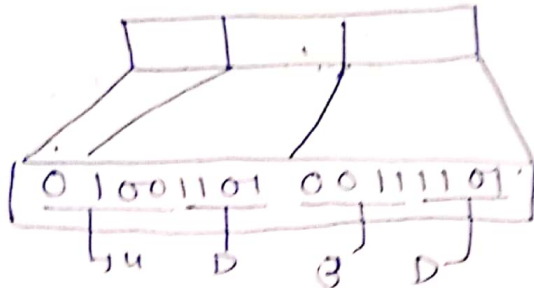
$$0.956 * 2 = 1.912 \rightarrow 1$$

$$0.239 \times 2^{13} \Rightarrow 0.00111101 \times 2^{13}$$

$$M = 00111101$$

$$e = 13$$

$$E = 13 + 64 = 77 = 1001101$$



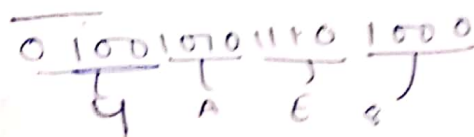
Q5 $0.00111101 \Rightarrow \text{Implicit} = 1.11101000 \times 2^{-3} \times 2^3$

$$1.11101000 \times 2^{10}$$

$$M = 11101000$$

$$e = 10$$

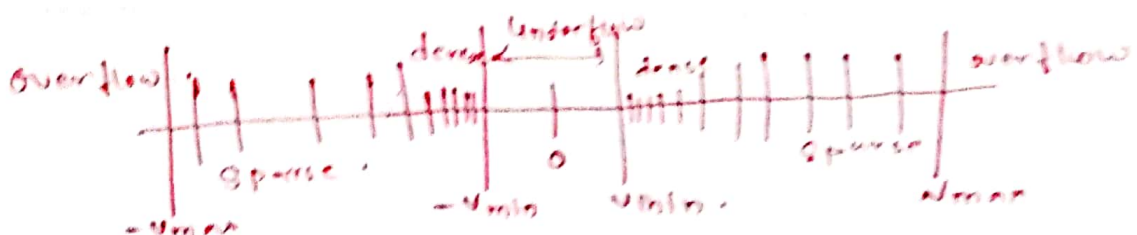
$$E = 10 + 64 = 74 = 1001010$$



Note: IF more no. of bits in mantissa means more precision.

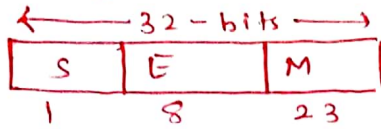
IF more no. of bits in exponent means larger range of numbers

This conventional floating point represent cannot store '0'



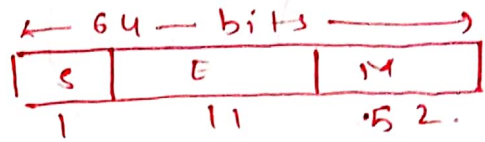
→ IEEE - 754 Floating point representation

Single precision



$$\text{bias} = 127 = 2^{8-1} - 1 = 2^7 - 1$$

Double precision



$$\text{bias} = 1023 = 2^{10} - 1 = 2^10 - 1$$

→ Some combinations of exponent are kept reserved to represent special numbers by keeping the bias value 127 or 1023.

→ Exponent combination all zeros or all ones represent special numbers

S	E	M	Notation.
0	000...0	0...0	+0
1	000...0	0...0	-0
0	11...1	0...0	$+\infty$
1	11...1	0...0	$-\infty$
0/1	11...1	$M \neq 0$	Not any number (NaN)
0/1	000...0	$M \neq 0$	Fraction or denormalized no.
0/1	$E \neq 00...0$ and $E \neq 11...1$	$M = xx...x$	Implicit normalized.

* Fraction or Denormalized no. :-

A very ^{very} small number which cannot be implicitly normalized.

min allowed exponent = -126.

upto this if we are getting normalized number well and good if not then store the mantissa as it is

and keep Exponent bits = all '0's

eg:

$$\rightarrow 0.000 \dots 0101$$

normalized form

$$1.01 * 2^{-130}$$

$$e = -130$$

~~$$E = -130 + 127 = -3$$~~

$$0.000101 * 2^{-126}$$

$$M = 000101$$

$$E = 0000 \dots$$

This could be Representation of Fraction no possible.

$$\rightarrow 0.000 \dots 0101$$

implicit normalized

$$1.01 * 2^{-127}$$

$$e = -127$$

$$E = -127 + 127 = 0 \Rightarrow \text{not allowed for implicit normal form.}$$

Value

$$(\text{implicit}) = (-1)^s * 1.M * 2^{E-\text{bias}}$$

$$\text{value (denormalized)} = (-1)^s * 0.M * 2^{-126 \text{ or } -1022}$$

Q7

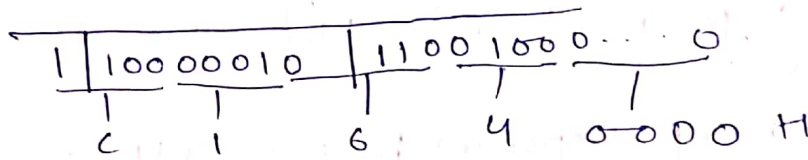
-14.25 = 1110.01 ⇒ Implicit normalisation

1.11001 × 2³

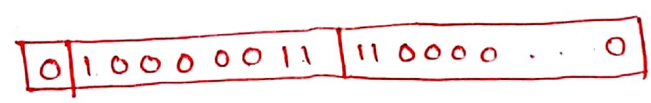
M = 11001000...0

e = 3

E = 3 + 127 = 130 = 10000010



→ Consider the following 32-bit register that denotes single precision floating point number in IEEE 754 format. What is its decimal value representation?

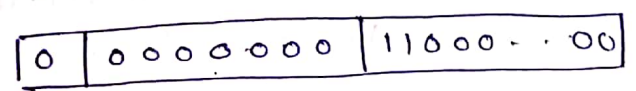


Exp → 131

Value = (-1)⁰ * 1.11 * 2¹³¹⁻¹²⁷

= 1.11 * 2⁴

= 11100 = +28



Exp → all '0' Man ≠ 0

∴ Denormalized no.

Value = (-1)⁰ * 0.11 * 2⁻¹²⁶

11.0 * 2⁻¹²⁷ * 2⁻¹²⁶

3 * 2⁻¹²⁸