

10

# COMPUTER ORGANISATION

What is the functionality of computer?

Computer:- It converts one form of data into another form under the control of program. While execution of program, one form of data is converted into another form. Therefore, functionality of computer is execution of the program.

Program:- It is the sequence of instructions along with data.

Instruction:- It is a binary sequence which is designed inside the processor to perform some task. i.e. binary sequence is bind with the operation.  $\left[ \text{Binary seq} \leftrightarrow \text{operat}^n \text{ rel at process design} \right]$

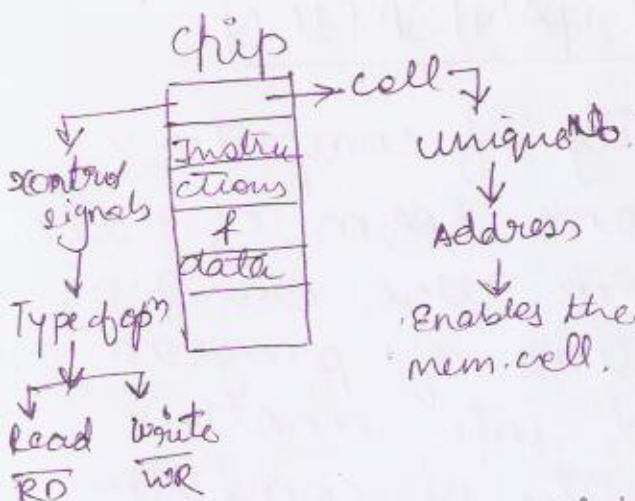
Data:- It is a binary sequence which is associated with the value based on the weighing factor.

To execute the program, the computer contains three (3) basic components named as: (1) CPU (2) Memory (3) I/O (External communication)  
(Processing) (Storage)

Program is initially stored in memory.

## MEMORY ORGANIZATION

- Memory chip is divided into equal parts called as cells.
- Each cell is able to be identified with unique no. called as address.
- Each cell is able to recognize 2 control signals, i.e. read and write.

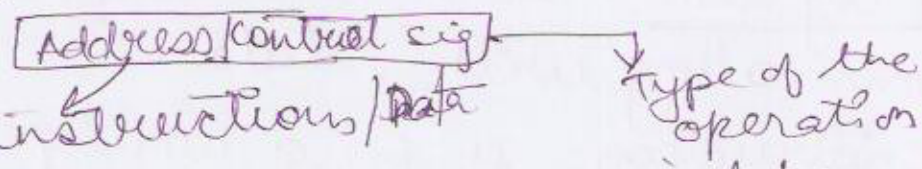


Additional CKT for address. Not in memory.

- To access any memory content, there is a need of address along with control signal.
- To execute the

program, ~~memory~~ CPU generates the memory request to access the memory to read the instructions.

The memory request contains



The no. of address bits required to access the memory is depending on the memory chip configuration.

MEMORY CHIP CONFIGURATION: The memory chips are configured as follows:-

64Kx8, 64Kx16, 64Kx32, ...

Let us consider 64.8 memory chip

(64K) x 8 → capacity of memory chip. i.e.

64K → # of memory cells present in the memory chip. i.e.  $2^{16}$  cells → # of address bits required to access 1 cell. i.e. 16 bits.

(NOTE :- If the memory chip contains  $2^n$  cells, n-bit address is required to access 1 cell among  $2^n$  cells) Address space/Address range/mem. map of the mem chip i.e.

000... (16 times) 0 to 111... (16 times) 1 →

0000H to FFFFH (Based on limited human recognition space is converted into hexadecimal address space. (0000)<sub>16</sub>/H/0x to (FFFF)<sub>16</sub>/H

64Kx8 → Cell Size.

Based on the cell size, the memory-chips are classified into 2 types → byte addressable memory and word addressable memory.

Byte addressable memory :- When the memory cell size is 8 bits, then the corresponding address space is called as byte address.

Word addressable memory :- When the memory cell size is based on the word length of the processor then the corresponding address space is called as word address.

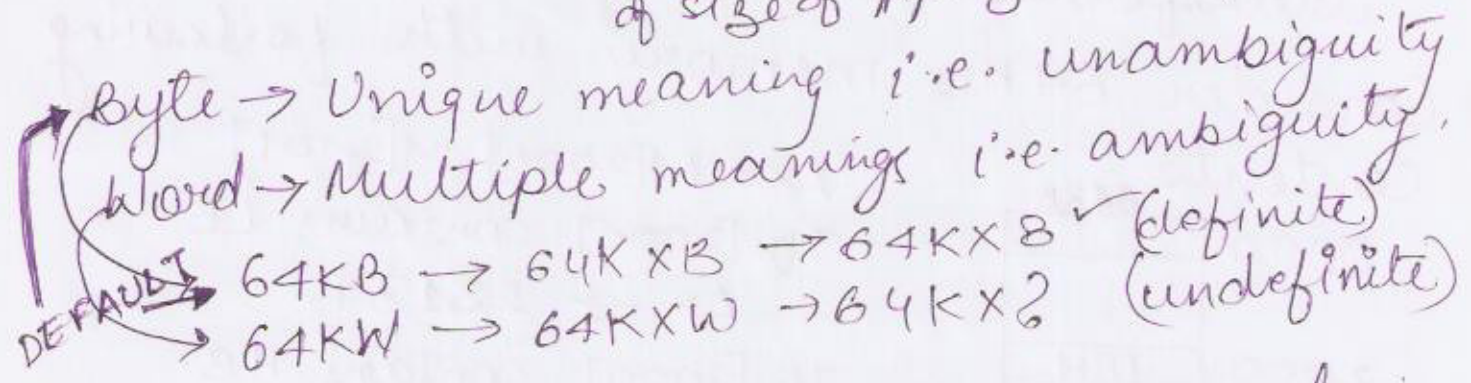
Processor                      Byte (b)

Word (w)  
word size = word length  
of the processor

8 bit  $\mu P$   
16 bit  $\mu P$   
32 bit  $\mu P$   
n-bit  $\mu P$

8 bit  
8 bit  
8 bit  
8 bit  
fixed irrespective  
of size of  $\mu P$  size

8 bit  
16 bit  
32 bit  
n bit  
↑  
variable



NOTE :- When the processor word size is greater than the memory cell size then there is a need of accessing the multiple memory cells to transfer the word from memory to processor. Eg. 16 bit processor :- MOV AX, [2000]

MOV  $\rightarrow$  Mnemonic  $\rightarrow$  User understandable format of the operation  $\rightarrow$  converted to opcode (M/c understandable format of the operation)  $\rightarrow$  Indicates type of operation, MOV indicates data transfer operation.

AX  $\rightarrow$  Destination <sup>reg</sup> (Registers are always referred by the names) ~~reg~~  $\rightarrow$  16 bits.

Register size is equal to processor size.

[2000]  $\rightarrow$  Source  $\rightarrow$  Memory (Mem. is referred by the addresses)  $\rightarrow$  Byte address.  $\rightarrow$  points to 8 bit data (cell size). But 2 cells are required  $\therefore$  {M[2000], M[2001]} data sent to AX.

Eg  $\rightarrow$  32 bit processor  
 MOV R0, [2000]  $\rightarrow$  operation  $\rightarrow$  Register (32bit)  
 $\rightarrow$  DEST  $\rightarrow$  SRC (Memory add) (8bit)  
 {M[2000], M[2001], M[2002], M[2003]} data transferred to R0 register.

\* Consider 64KB memory with following contents

64KB	
0000	
!	
2000	18H
2001	23H
2002	04H
2003	16H
!	
ffff	

AX = M[2000] M[2001]  
 If [2000] contains LB  
 AX  $\rightarrow$  2318H

If [2000] contains HB  
 AX  $\rightarrow$  1823H

R0 = M[2000] M[2001] M[2002] M[2003]  
 If [2000] contains LB  
 R0 = 16042318H  
 otherwise  
 R0 = 18230416H.

During the execution of above instructions, due to the lack of order of data storage, the instructions generates the multiple outputs.

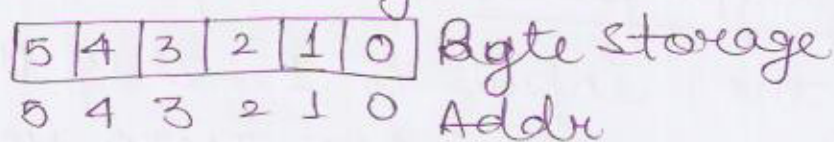
To avoid the above problem, there is a need of the memory address interpretation (Endian-ness) Mechanisms.

ENDIAN-NESS MECHANISMS:- These mechanisms show the order of data storage. They are 2 types:-

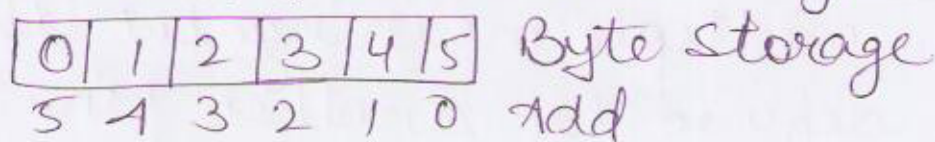
(i) little endian

(ii) Big endian

little endian mechanism shows the order of data storage as lower address contains the lower bytes & higher address contains the higher ".



Big endian mechanism shows the order of data storage as lower address contains the higher byte and higher address contains the lower byte.



NOTE:- Little endian is a default address interpretation <sup>tech</sup> used in processor design. ~~It is~~

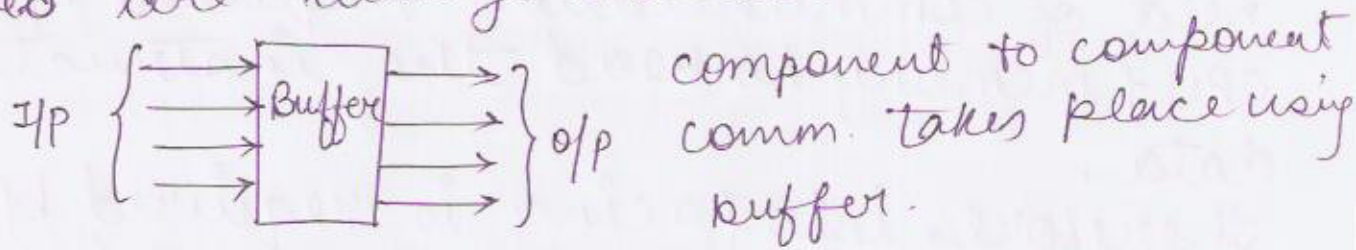
Therefore the above instructions produce the outputs as AX 

HB	LB
23	18
2001	2000

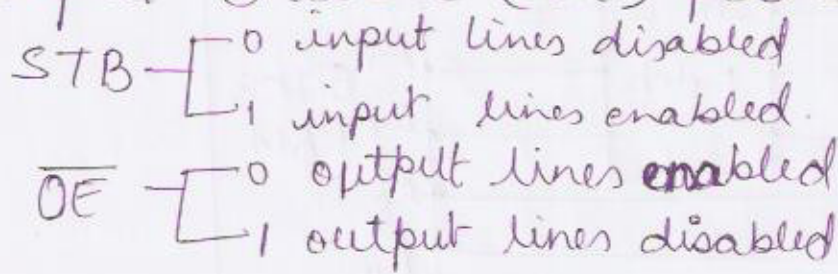
 $\Rightarrow$  2318H.



Buffer:- Buffer is a temporary storage component without locking. Input + output lines are always enabled.

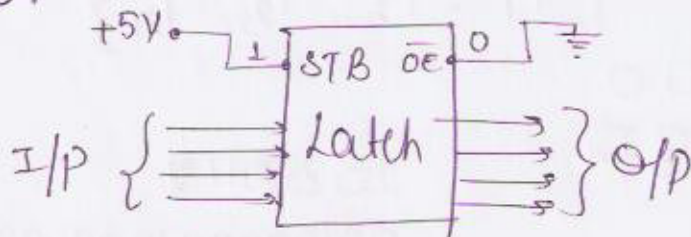


Latch:- Latch is a temporary storage component with locking i.e. ip + op lines ~~and~~ are there in the control of locking pins i.e. strobe (STB) and output enable ( $\overline{OE}$ ) ~~pins~~ pins respectively.



NOTE:- Latch may be configured as the buffer but buffer can never be configured as latch.

Latch is configured as the buffer by connecting STB pin to +5V and  $\overline{OE}$  pin to GND.



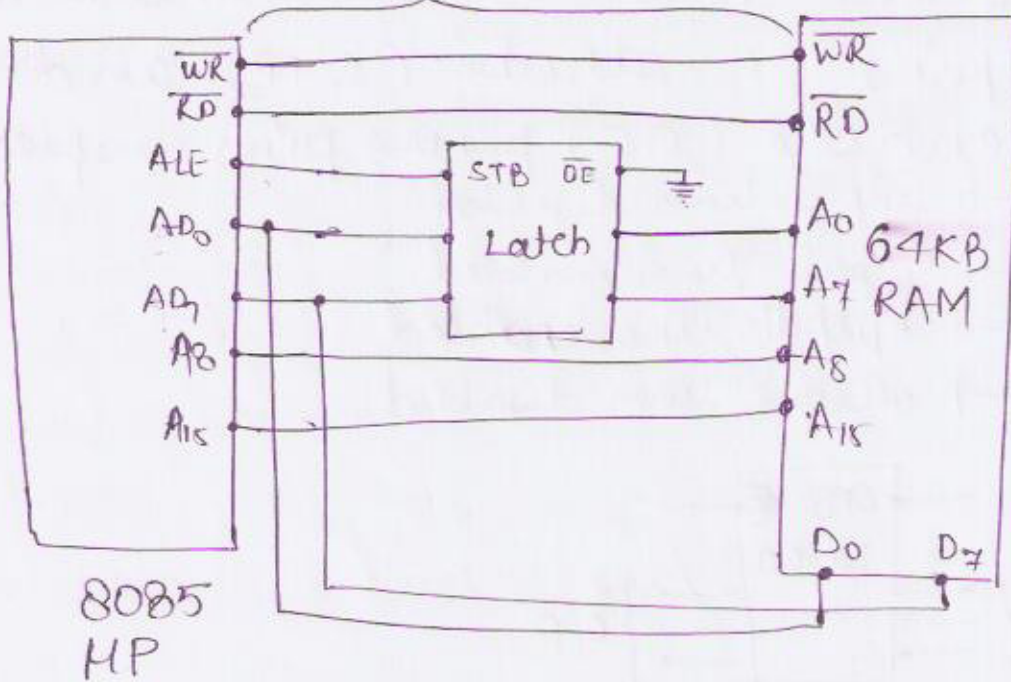
Latch as Buffer.

# MEMORY INTERFACING :

To execute the program, there is a need of communication required b/w CPU & memory to read the instructions & data.

Therefore, interfacing is required b/w the CPU & memory.

done using BUS in practical.



CPU generates the mem. request:  $2080H$   
 $\overline{RD}$

M/C cycles - The no. of clock cycles required to access memory (to read or write) is called machine cycles.

Machine cycles  $\Rightarrow \{T_1, T_2, T_3, T_4\}$

ALE  $\rightarrow$   $\downarrow$  0, 0, 0  
 $T_1, T_2, T_3, T_4$

$\overline{RD}$   $\rightarrow$  0  
 $T_2$

$\overline{WR}$   $\rightarrow$   $\downarrow$   
 $T_3$

$2080H \downarrow$

00100000 10000000  
 $\downarrow$                      $\downarrow$                      $\downarrow$   
 $A_{15}-A_8$                      $A_7-A_0$



T1  $\Rightarrow$   $A_{D_0}-A_{D_7} \rightarrow 80$ ,  $A_8-A_{15} \rightarrow 20 \leftarrow$  CPU, ALE  $\rightarrow 1$   
 A0-A7  $\rightarrow 80$ ,  $A_8-A_{15} \rightarrow 20 \leftarrow$  RAM, STB  $\rightarrow 1$

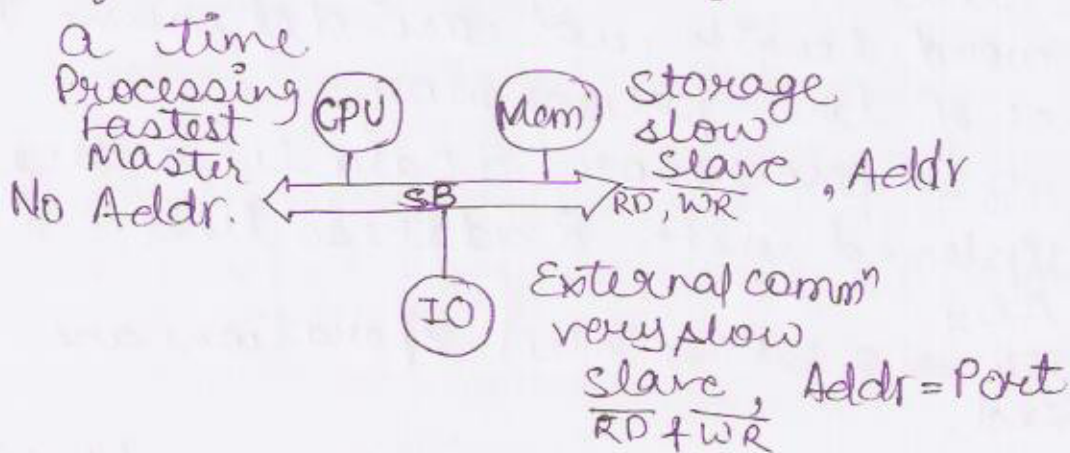
T2  $\Rightarrow$   $\overline{WR} - 1$   
 $\overline{RD} - 0$  } CPU, RAM ALE=0, STB=0.

T3, T4  $\rightarrow$  Binary sequence at  $D_0-D_7$  in RAM  
 to  $A_{D_0}-A_{D_7}$ . This sequence is not  
 forwarded to latch. Data is moving  
 between memory and CPU. 2 cycles  
 because memory is slower than CPU.

NOTE: In every m/c cycle, in first time  
 state, ALE pin is in the high/enable  
 state, & remaining time states, ALE is  
 disabled to carry the address & data  
 efficiently over the TM pins.

BUS: Bus is a communication channel  
 used between the major components  
 of the system i.e. CPU, memory and I/O.

- It is also called as the system bus.
- The key characteristic of the bus is  
 shared transmission media. The limitation  
 of the bus is only one transmission at



The communication channel contains 3 lines to provide communication.

- 1) Address lines
- 2) Data lines
- 3) Control lines

### Address lines

These lines are used to carry the address towards memory & I/O. They are unidirectional.

Based on the no. of address lines, we can determine the capacity of the memory system. Eg. 8085 supports 16 address lines i.e.  $AD_0 - AD_7$  &  $A_8 - A_{15}$ .

$\Rightarrow 2^{16}$  address cells  $\Rightarrow 64KB$  (64K cells).

8086 supports 20 address lines i.e.

$AD_0 - AD_{15}$   $A_{16} - A_{19}$ .  $\Rightarrow 2^{20}$  address cells  $\Rightarrow 1M$  cells

$\Rightarrow 1MB$  storage space.

Data Lines:- These lines are used to carry the binary sequence between CPU, memory & I/O. Therefore, data lines are bidirectional.

Based on the no. of data lines, we can determine the word length of the processor. Based on word length, we can determine the performance of the processor.

Eg: In 8085 processor, 8 data lines are time multiplexed with 8 address lines i.e.  $AD_0 - AD_7$ .

Word length  $\Rightarrow 8$  bit  $\Rightarrow 8$  bit operations are performed.

In 8086, 16 data lines are time multiplexed

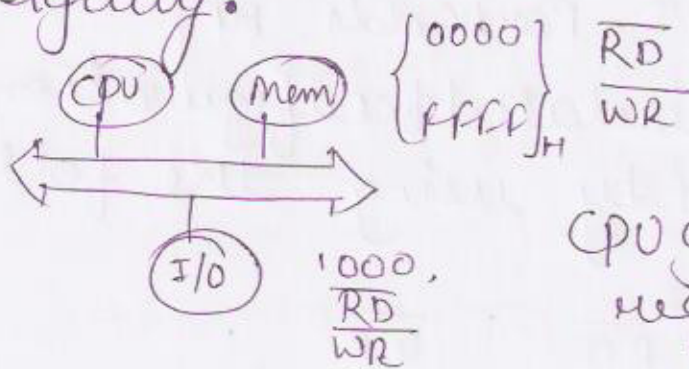
with 16 Address lines i.e.  $AD_0 - AD_{15}$ .

Total word length = 16 bit  $\Rightarrow$  16 bit operations are performed.

### Control lines:-

These lines are used to carry the control signals & timing signals. Control signals are used to indicate the type of the operation and timing signals are used to synchronise the memory & I/O operations with processor clock. Control lines are bidirectional.

NOTE:- When there is a common channel, common address space, and common control sequence signals maintained between the memory and I/O, then there is a possibility of ambiguity.



CPU generates request  $\left\{ \begin{array}{l} 1000 \rightarrow \text{AL's} \\ \overline{\text{RD}} \rightarrow \text{CL's} \end{array} \right\}$

(Conflict Ambiguity)

I/O controller  $\left\{ \begin{array}{l} \text{both valid} \\ \text{both valid} \end{array} \right\}$

To avoid the above problem, there is a need of bus configuration

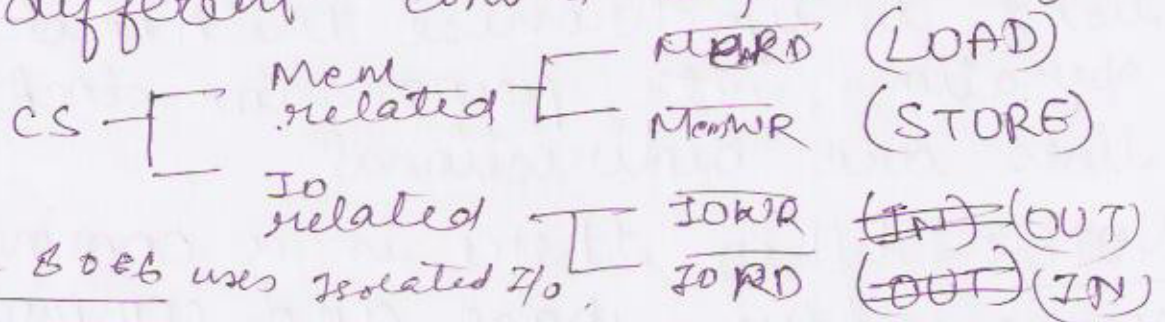
Bus Configuration:- There are 3 bus

configuration in systems design.

- (i) I/O Processor (IOP)
- (ii) Isolated processor
- (iii) Memory mapped I/O

① IOP (Input Output processor):- This configuration uses the separate buses for I/O and memory. Introducing an additional ~~process~~ bus is an expensive process. So this config is used in the high performance processor design.

② Isolated I/O:- This configuration uses the same bus and common address space but different control ~~sequence~~ signals.



8085 & 8086 uses isolated I/O.

- LOAD → LOAD R0, [2000]
- STORE → STORE [2000], R1
- IN → IN R0, Port Addr
- OUT → OUT PortAddr, R4

In 8086, isolated I/O configuration is implemented by using the following 4 pins.

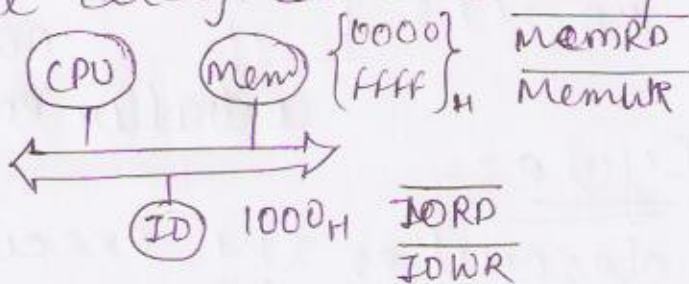
	M/IO (dual pin)	RD (unipin)	WR (unipin)	
$\overline{\text{MemRD}}$	1	0	1	Memory Read
$\overline{\text{MemWR}}$	1	1	0	Memory Write
$\overline{\text{IORR}}$	0	0	1	IO read
$\overline{\text{IOWR}}$	0	1	0	IO write

In 8085, isolated I/O configuration is implemented by using the following 4 pins:-

different from 8086

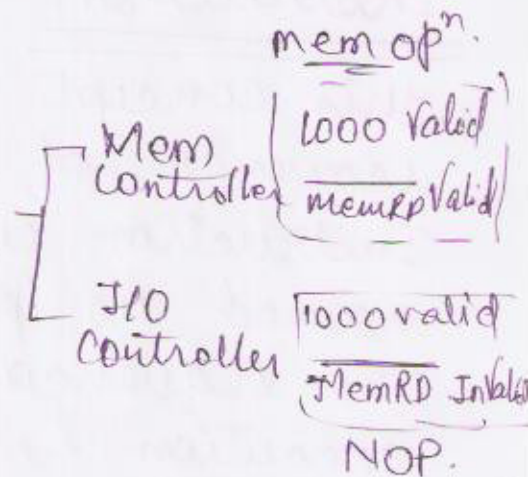
$\overline{IO/\overline{M}}$	$\overline{RD}$	$\overline{WR}$	
1	0	1	$\overline{IORD}$
1	1	0	$\overline{IOWR}$
0	0	1	$\overline{MemRD}$
0	1	0	$\overline{MemWR}$

So, the diagram changes to:



CPU generates mem req.

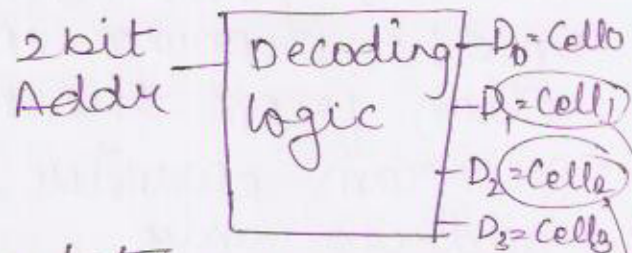
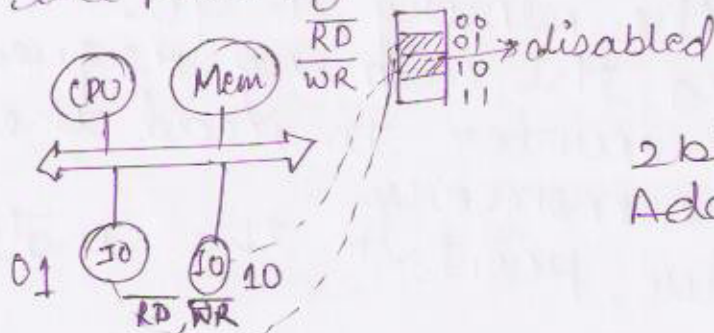
1000 → AL's  
 $\overline{MemRD}$  → CL's



### Memory Mapped I/O

This configuration uses the common bus and common control signal but unique address space.

Unique address space is maintained by taking some of the addresses from the pool of the address space, assigned them to I/O ports. that means, port addresses are the part of the memory address space.



Disadvt → Not complete memory space is not in use. (Not allowed to be used by user. reserved by processor)

The disadvantage in memory mapped I/O is that the complete memory is not in the use.

CPU generates:  $OO \rightarrow AI$ 's  
 Mem Req:  $\overline{RP} \rightarrow CI$ 's

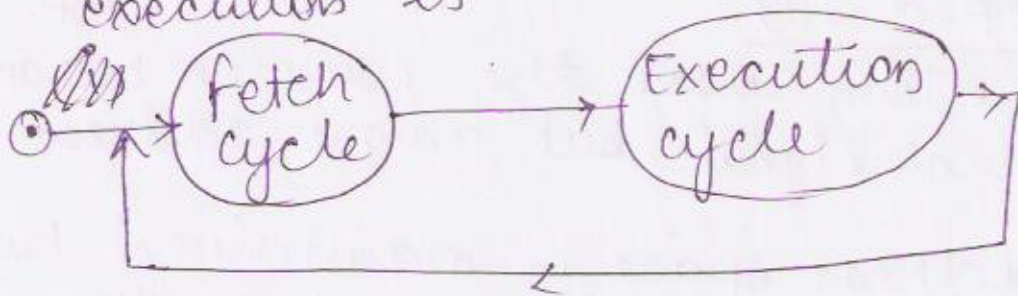
[	Mem Controller	$OO$ valid	]	Mem op <sup>n</sup>
		$\overline{RP}$ - valid		
[	I/O Controller	$OO$ - Invalid	]	Nop
		$\overline{RP}$ - valid		

### Instruction Cycle :-

This concept describes the execution sequence of an instruction.

Instruction cycle consists of 2 subcycles named as fetch & execution cycle.

The sequence diagram of instruction execution is



### Fetch cycle :-

- The objective of the fetch cycle is to read the instruction related binary sequence from the memory to CPU.
- The CPU generates the memory request based on the program counter to read the instruction from the memory.
- Program counter points to the instruction address only.
- The functionality of PC is to hold the starting instruction address and immediately points to the next instruction address.
- Here the starting address is provided by the



(b) Assume that memory is a word addressable with word size of 32 bits. The program has been loaded in memory with the starting address of 1000 (decimal) onwards. During the execution of  $I_8$  what could be the value of PC?

Addresses:- 1000-1001, 1002, 1003, 1004-5, 6, 7-8, 9.

NOTE:- When the processor is supported with fixed length instructions then step size is a fixed constant. When the processor is supported with variable length instructions, then step size is also variable.

Q. A computer has 24 bit instruction. The program has been loaded in the memory with the starting address of 300 (decimal) onwards. Which one of the following is the valid program counter values!?

- (a) 400
- (b) 500
- (c) 600 ✓
- (d) 700

Execution cycle:-

- The objective of the execution cycle, currently fetched instruction undergoes processing.
- To process the instruction there is a need of identifying the type of operation.

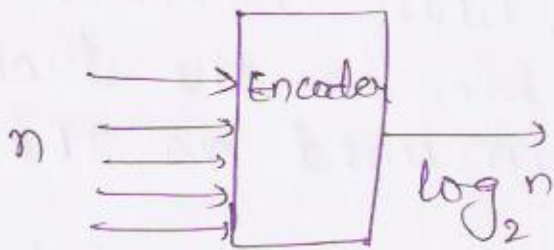


associated with the binary sequence.

- opcode indicates the type of operation
- opcode is present in the instruction but that is defined by using instruction format.
- Instruction format shows the layout of an instruction, that means, it describes the internal structure of the instruction.

### Encoding (Many to 1 Mapping)

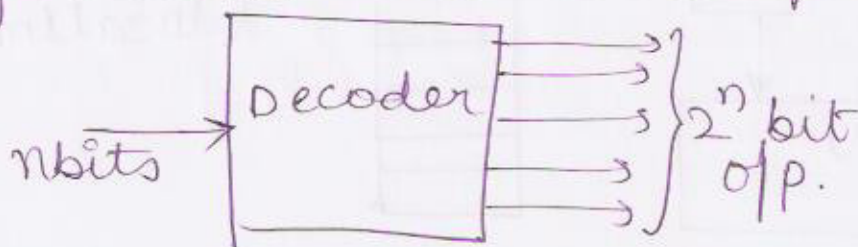
In the encoding process,  $n$  different combinations can be represented using  $\log_2 n$  bits to indicate only one operation at a time.



NOTE:- Instructions are always represented in the form of encoded binary format. This means when the processor supports 16 different operations, then 4 bit opcode is required to indicate 1 operation at a time.

### Decoding (1 to many mapping)

In the decoding process,  $n$  bit decoder generates the  $2^n$  output combinations.



NOTE During the execution of the instructions decoder is required to identify the operation associated with the encoded binary format.

Instruction format is classified into 5 types based on the type of CPU organization.

The CPU organization is classified into 3 types based on the availability of ALU operands.

### Classification of Instruction set Architecture

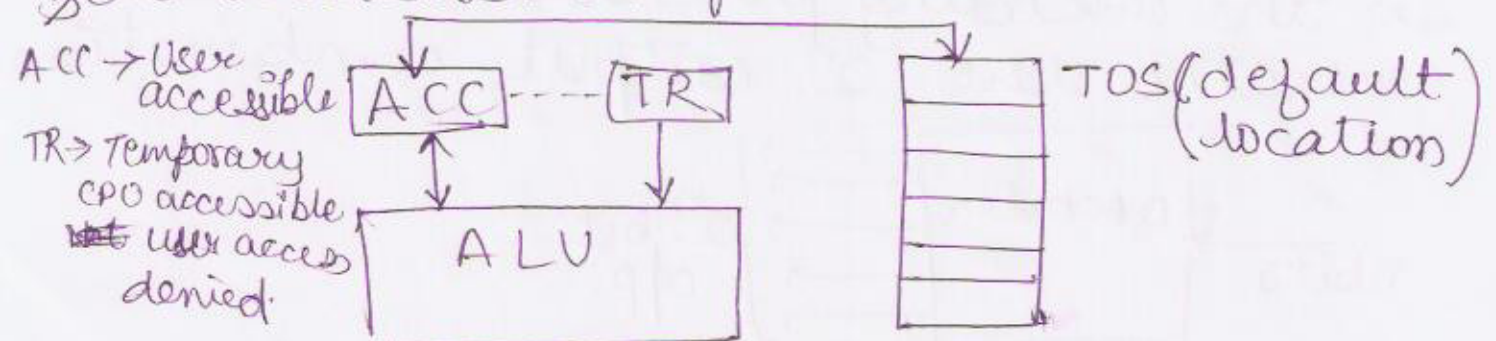
- (1) Stack machine
- (2) Accumulator machine
- (3) General register organization

#### Stack Machine

In this organization, both the ALU operands are available in the stack. The same stack is also used as the destination.

The stack contains only one entry point and only 1 exit point i.e. top of stack (tos). So tos is default location. Therefore, None need to pass the tos address explicitly in the instruction.

- The instruction format contains only 1 field i.e. opcode
- The above instruction format is known as zero instruction format.



OPCODE

Type of op<sup>n</sup>

ALU op<sup>n</sup>: src1 src2 Dest  
TOS TOS TOS.

POP → POP → operate → push.

Eg. ADD.

Sequence of operations → POP, pop, +, push

NOTE: - In the stack machine ALU operations are the zero address instructions but data transfer operations are not zero<sup>addr.</sup> instructions

Q. Consider the following program segment used to execute on a stack machine. Each arithmetic operation pops the 2<sup>nd</sup> operand, pops the 1<sup>st</sup> operand, operates on them & push back the result onto the stack.

push b  
push x  
ADD  
pop c  
push c

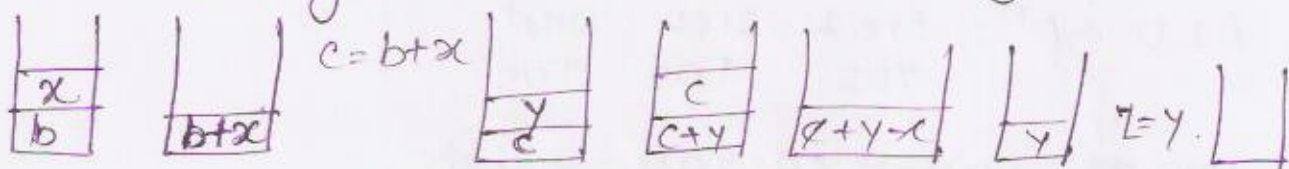
push y  
ADD  
push c  
SUB  
POP z

Which of the following statements is/are true?

1. If each push and pop instruction requires 5 bytes of storage and each arithmetic instruction requires 1 byte of storage, then the whole program requires 40 bytes of storage.
2. At the <sup>end of</sup> execution, Z contains the same value as y.
3. At the end of execution stack is empty. (Consider stack is initially empty)

- a) 1 only (c) 2 & 3 only ✓  
 b) 2 only (d) 1, 2, 3.

Total storage =  $7 \times 5 + 8 \times 1 = 38$  bytes.



**2D**  
 (2) Accumulator Machine

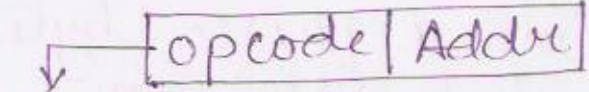
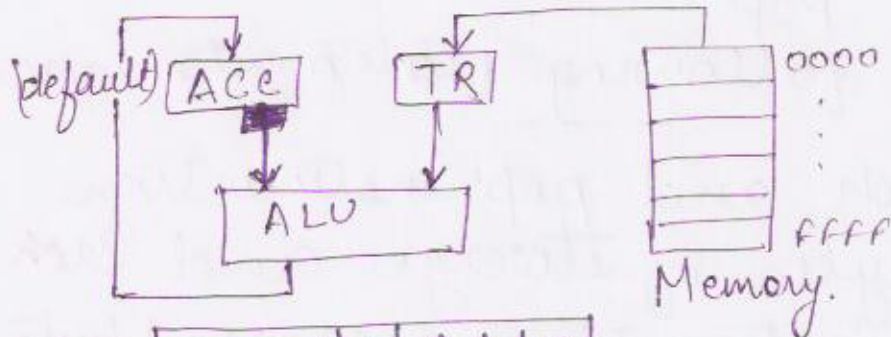
In this organization, the first ALU operand is default available in Accumulator. The same accumulator is used as the destination.

The CPU contains only 1 accumulator register so it becomes the default location.

The 2<sup>nd</sup> ALU operand, may be present in the memory or register. So, memory address or register name is required in the instruction to access the operand 2.

Therefore, the instruction format contains 2 fields i.e. opcode | Addr

The above format is known as 1 address inst<sup>r</sup> format.



Type of operation

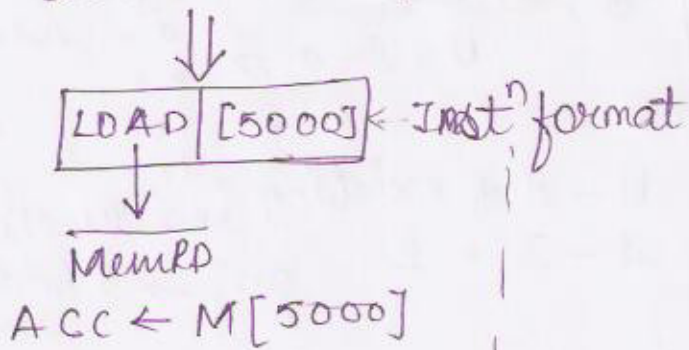
Data Transfer operat<sup>n</sup>

Source	Dest <sup>n</sup>
Acc	Mem (write) STORE
Mem	Acc (Read) LOAD

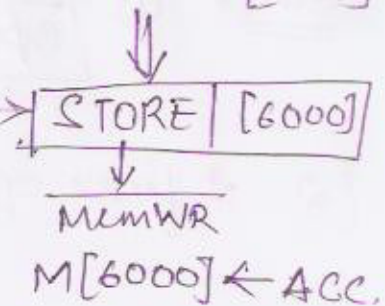
ALU OP<sup>n</sup>:

SRC 1	SRC 2	DSTN
ACC	Mem/Reg	ACC

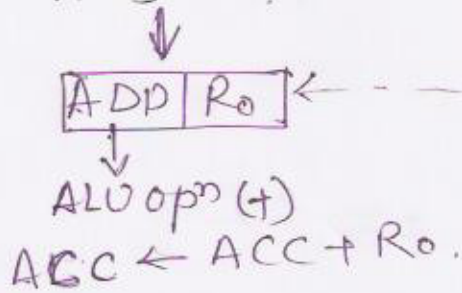
Eg 1: LOAD [5000]



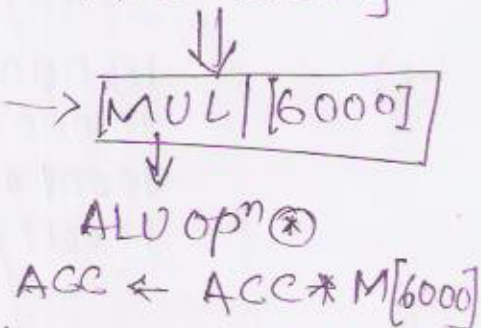
Eg 2: STORE [6000]



Eg 3: ADD R<sub>0</sub>



Eg 4: MUL [6000]



Q. Consider a hypothetical processor. It supports both 1 addr. & 0 addr instructions. It has 6 bit instructions & 4 bit addresses. If there exists 2 one address inst<sup>n</sup>s, how many 0 address instructions can be formulated?

NOTE: 1 Address m/c also supported by 0 address instructions if there is a free space after allocating the 1 address inst<sup>n</sup>. This technique is called as expand-opcode technique.

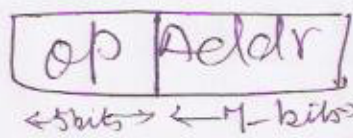
Step 1:- Identify the highest instruction format according to the problem statement. <sup>1 addr</sup>

S2:- Identify the no. of possible operations. (Instructions)

S3:- Identify the no. of free combinations (free opcodes) after allocating the high order instructions.

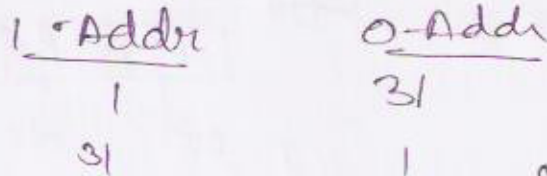
S4:- Calculate the no. of low order Inst<sup>n</sup> by multiplying the decoded value of the Addr field with the free opcode.





$$N = 2^5 = 32$$

Share the 32 op's. B/w the 1 Addr & zero add.



Range of 1 address = {1 to 31}  
 Range of 0 address =  $\{1 \times 2^4 + 2 \times 2^3 + \dots + 31 \times 2^1\}$

3) General Register Organization: - Based on the no. of registers supported by the processor, the architecture is divided into 2 types:-

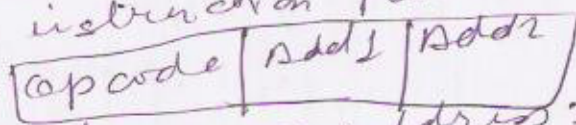
- (a) Register to memory reference architecture
- (b) Reg-Reg reference architecture.

Register to memory reference Architecture

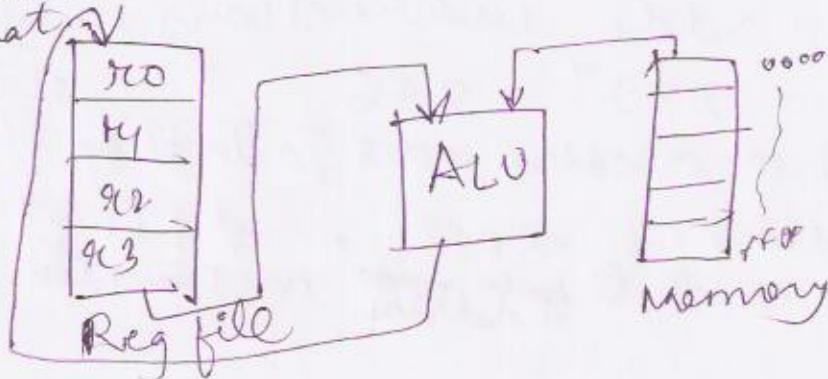
In this architecture, processor supports less no. of registers. That means, register file size is small.

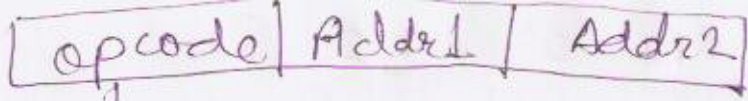
In this organization, the first ALU operand is always present in the registers, the same register is also used as the destination.

The second ALU operand is available either in memory or in registers. So this architecture supports an instruction format with 3 different fields -

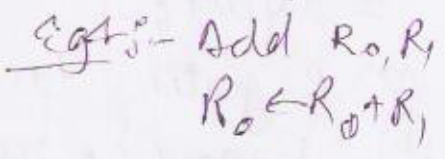
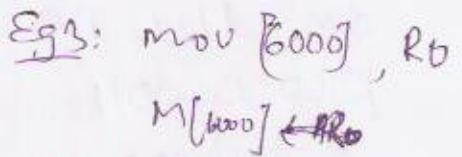
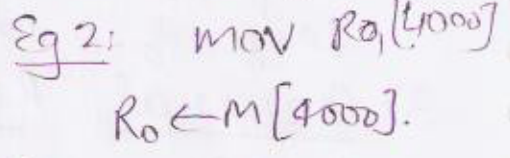
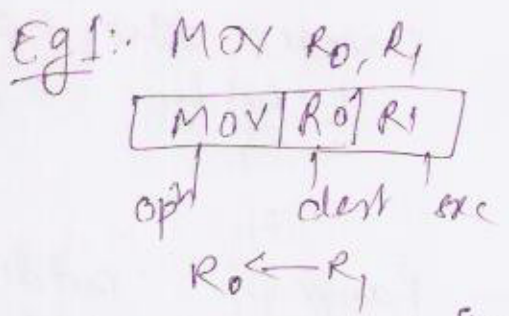
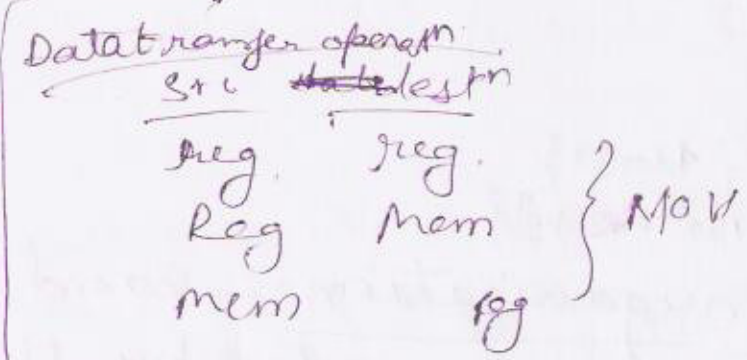


The above instruction is called as 2 address instruction format.

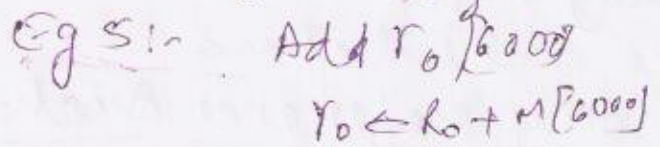
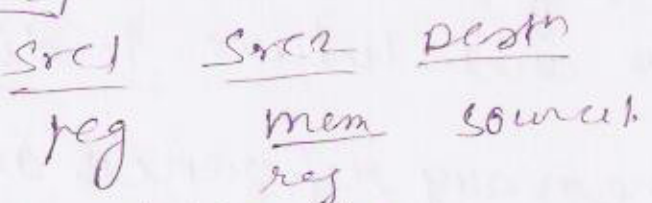




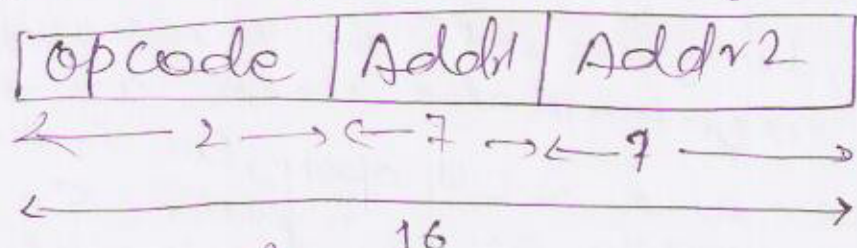
Type of operation



ALU op<sup>n</sup>:-



Q. Consider a processor that supports 2 addi, 1 addi & 0 addi. A 16 bit instruction is placed in 128M memory. If there exists 2-2 address instructions, 100-1 address instructions, how many instructions for 0 can be formulated?



$N = 2^2 = 4$

No of free opcodes after allocating the 2 addi instructions =  $4 - 2 = 2$

# of 1 addi instructions =  $\boxed{\text{op} \mid \text{Addr}_1 \mid \text{Addr}_2}$

$= 2 \times 2^7 = 256$

# of free opcodes after allocating the 1 addi instruction =  $256 - 10 = 156$

$\therefore$  NO. of ~~bits~~ instructions =  $156 \times 2^7$



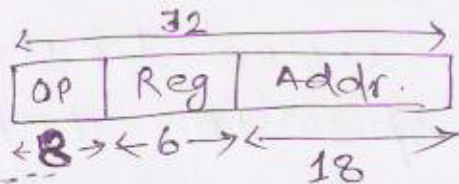
Q. A computer has an instruction format with 3 fields - i.e., opcode, register field which indicates 1 of the 60 processor register and memory address. A 32 bit instruction is placed in the 256KW memory.

(a) How many no. of operations are possible.

(b) If there exists  $n$  2-address instructions which uses both register and memory reference, then how many 1 address memory reference instructions are possible.  
(register merging with opcodes)

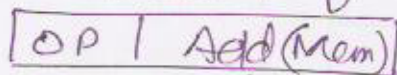
(c) If there exists  $n$  2-address instructions which uses both register & memory reference,  $m$  1 address register references then how many zero address instructions are possible.

Ans:-  
(a)



$N = 256 \rightarrow$  # free opcodes after allocating 2 addr. instructions =  $(256 - n)$

# of 1 address mem ref =  $(256 - n) \times 2^6$   $\rightarrow$  Reg merged.



# of 1 addr reg ref =  $(256 - n) \times 2^{18}$

# of free opcodes after allocating 1 addr mem ref  
Inst<sup>n</sup> =  $[(256 - n) \times 2^6] - m$

# of free opcodes after allocating 1 addr reg ref  
Inst<sup>n</sup> =  $[(256 - n) \times 2^{18}] - M$

# of free opcodes after 0-address =  $(((256 - n) \times 2^6) - m) \times 2^{18}$

0 addr  $\rightarrow ((256 - n) \times 2^{18}) - M$

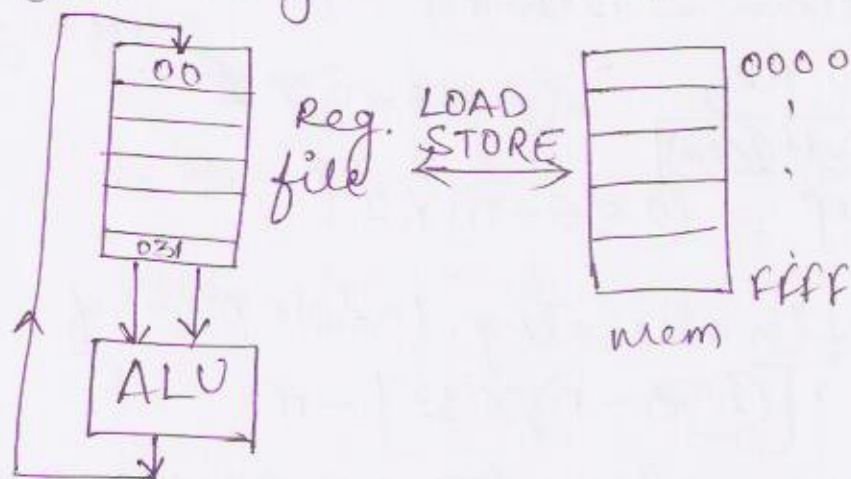
## Register to register Reference architecture

In this architecture, the processor supports more no. of registers i.e. register file size is large. No possibility of register overwriting.

In this architecture ALU operations are performed only on registers. i.e. Both the ALU operands are available in register. Due to the more no. of registers present in this processor, use a separate register to store the output of the ALU operations.

Therefore ALU operations require 3 address fields in the instruction. All the 3 address fields contains register names.

In this architecture, 2 address instruction archit format is used to implement the data transfer operation. LOAD and STORE instr are used to transfer the data b/w the registers and memory.



The 3 instr format is 

opcode	Addr1	Addr2	Addr3
--------	-------	-------	-------

opcode → Type of operation → ALU operation

Addr1 - Dst<sup>n</sup> - Reg

Addr2 - SRC1 - Reg

Addr3 - SRC2 - Reg



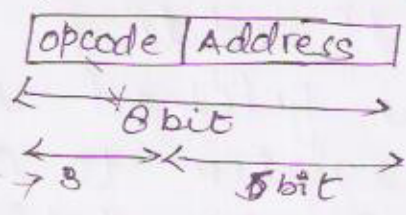
LOAD A	$Acc \leftarrow [A]$	$X \leftarrow (A+B) * (C+D)$	} Accumulator M/C.
ADD B	$Acc \leftarrow [A] + [B]$		
STORE P	$M[P] \leftarrow Acc$		
LOAD C	$Acc \leftarrow [C]$		
ADD D	$Acc \leftarrow [C] + [D]$		
MUL P	$Acc \leftarrow (A+B) * (C+D)$		
STORE X	$M[X] \leftarrow (A+B) * (C+D)$		

PUSH A	[A]	} Stack Machine
PUSH B	[B] [A]	
ADD	[A+B] [A]	
PUSH C	[C] [A+B] [A]	
PUSH D	[D] [C] [A+B] [A]	
ADD	[C+D] [A+B] [A]	
MUL	[A+B] * [C+D] [A+B] [A]	
POP X	$M[X] \leftarrow (A+B) * (C+D)$	

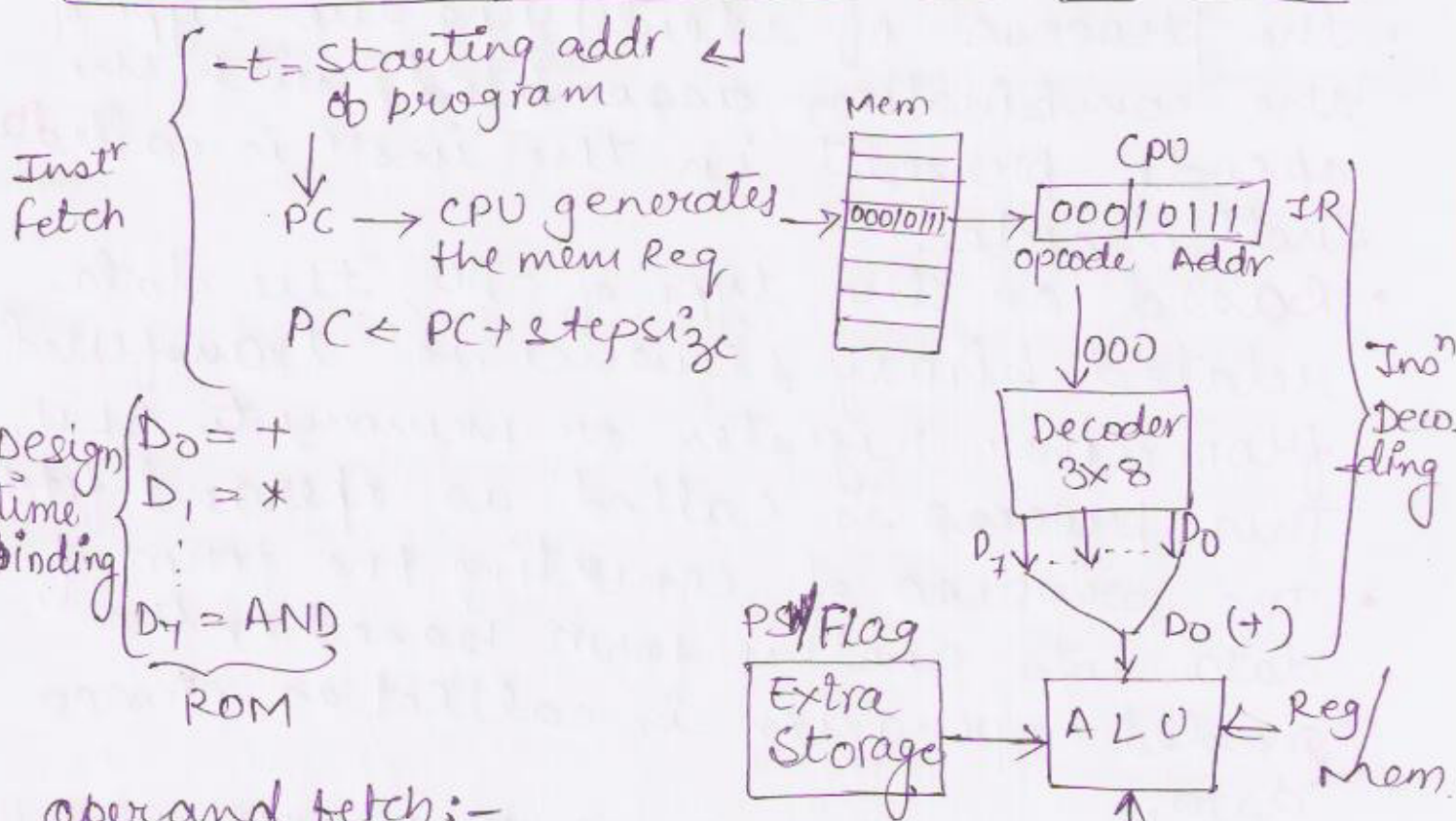
MOV R0, A	$R_0 \leftarrow [A]$	} Reg to Mem.
ADD R0, B	$R_0 \leftarrow R_0 + M[B]$	
MOV R1, C	$R_1 \leftarrow M[C]$	
ADD R1, D	$R_1 \leftarrow R_1 + M[D]$	
MUL R0, R1	$R_0 \leftarrow R_0 * R_1$	
MOV X, R0	$M[X] \leftarrow R_0$	

LOAD R0, A	$R_0 \leftarrow M[A]$	} Reg to Reg
LOAD R1, B	$R_1 \leftarrow M[B]$	
ADD R2, R0, R1	$R_2 \leftarrow R_0 + R_1$	
LOAD R0, C	$R_0 \leftarrow M[C]$	
LOAD R1, D	$R_1 \leftarrow M[D]$	
<del>ADD</del> <del>MUL</del> ADD R3, R0, R1	$R_3 \leftarrow R_0 + R_1$	
MUL R4, R2, R3	$R_4 \leftarrow R_2 * R_3$	
STORE X, R4	$M[X] \leftarrow R_4$	

NOTE:- To analyze the execution sequence of the program, let us consider an example CPU & its compatible instruction format i.e. accumulator m/c as an exg. CPU. It supports 8-bit instructions with 8 different operations. Its instr format is 1-address format.



EXECUTION SEQUENCE OF AN INSTR

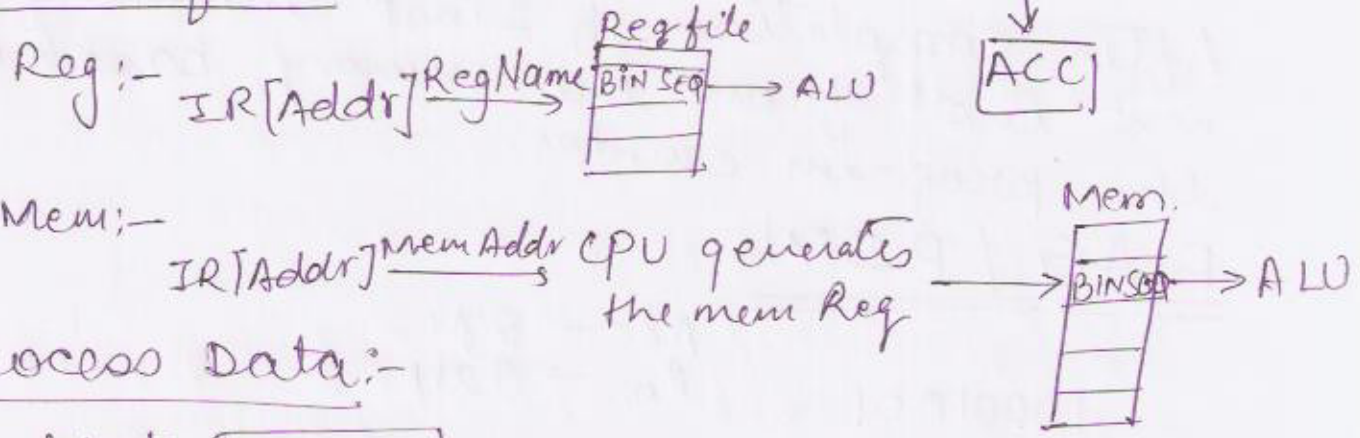


Design time binding

- $D_0 = +$
- $D_1 = *$
- $\vdots$
- $D_7 = \text{AND}$

ROM

operand fetch:-



Process Data:-



Based on the PC, the instr<sup>r</sup> related binary<sup>seq</sup> is transferred from memory to CPU. This process is called as instruction fetch. At the end of instr<sup>r</sup> fetch, PC is incremented to step size to point to the next instr<sup>r</sup> addr.

- All the fetched instr<sup>r</sup> is placed into the instr<sup>r</sup> register to decode because instr<sup>r</sup> format is predefined in this register (IR).
- The process of identifying the type of the combination associated with the opcode present in the instr<sup>r</sup> is called as instr<sup>r</sup> decode.
- Based on the type of CPU, the data related binary sequence is transferred from either register or memory to ALU. This process is called as operand fetch.
- The process of converting one form of data into another form based on the enabled operation is called as process data.

After completion of 1 instr<sup>r</sup> the CPU fetch the next instr<sup>r</sup> from the memory based on the program counter.

FLAG / PSW

ACC - 89H  
R0 → 82H

$$\begin{array}{r} 10001001 \\ 10000010 \\ \hline 10001011 \leftarrow \text{ACC.} \end{array}$$

↓ Extra, lost

To store the exception conditions, there exists extra storage provided in processor called as flag register / PSW (Program Status word).

Flag - flag is a flip-flop. Flip-flop is a bistate device. It is used to store the one bit information.

Flags are classified into 2 types:

- Conditional flag.
- Control flags

Conditional flag:-

These flags are set or reset based on the result nature of ALU. There are 6 conditional flags used in processor:-

- (i) Carry (ii) Parity (iii) Auxiliary carry (iv) zero
- (v) Sign flag (vi) overflow flag.

Carry:- It is a conditional flag set on "Is there an extra bit out of MSB".

$$\begin{cases} T = \text{set} = 1 = C \\ F = \text{reset} = 0 = \text{NC} \end{cases}$$

Carry flag is used to indicate the range exceeding conditions of the unsigned arithmetic operations.

In the unsigned representation, we can represent only the positive no. only.

n-bit unsigned range is 0 to  $2^n - 1$

Eg. 4-bit unsigned range is 0 to 15.

$$\begin{array}{r} 7 = 0111 \\ 3 = 0011 \\ \hline 10 = 1010 \\ \text{CY} = 0 \end{array}$$

↑  
Carry

$$\begin{array}{r} 9 = 1001 \\ 8 = 1000 \\ \hline 17 = \boxed{10001} \leftarrow \text{ACC} \\ \text{CY} = 1 \end{array}$$

Parity Flag:- "Is the ALU output (Acc content) contain even no. of 1's".

—  $T = \text{set} = 1 = \text{PE (even)}$   
—  $F = \text{reset} = 0 = \text{PO (odd)}$

This flag is used in the serial communication to prepare the error detection codes.

Auxiliary carry:- "Is there an extra bit from the lower nibble to higher nibble (3<sup>rd</sup> position to 4<sup>th</sup> position)"

Used in BCD arithmetic to indicate the range exceeding conditions.

Zero :- "Is the ALU output (Acc) zero"

—  $T = \text{set} = 1 = Z$   
—  $F = \text{reset} = 0 = NZ$

Zero flag is used to implement the control structures

Control Struct

- (i) Selection
- (ii) Iteration
- (iii) Subprogram

Sign:- "Is the MSB of ALU output (Acc) one"

—  $T = \text{set} = 1 = \text{NG (Negative logic)}$   
—  $F = \text{reset} = 0 = \text{P (Positive logic)}$

Overflow:- "There is a carry into the MSB & no carry out of the MSB or vice versa."

—  $T = \text{set} = 1 = \text{OV}$   
—  $F = \text{reset} = 0 = \text{NV}$

It is used in the signed arithmetic operations to indicate the range exceeding conditions.

- SIGNED NOS:-
- In the signed format we can represent both positive and negative numbers.
  - 2's complement no. representation is used to indicate signed nos in the processor.
  - when the MSB bit is 1, then the no. is negative. So take the 2's complement to conclude the value.



When the MSB bit is zero, no. is positive so no need to take the complement to conclude the value.

n-bit 2's complement range is  $(-2^{n-1})$  to  $+2^{n-1}$

Eg 4-bit 2's complement range =  $\{-8$  to  $+7\}$

2's comple. Binary

Equivalent Decimal

0000	+0
0001	+1
0010	+2
0011	+3
0100	+4
0101	+5
0110	+6
0111	+7
1000	-8
1001	-7
1010	-6
1011	-5
1100	-4
1101	-3
1110	-2
1111	-1

Eg.

$$\begin{array}{r} +7 \\ +3 \\ \hline +10 \\ \text{OV} = \text{set} \end{array}$$

$$\begin{array}{r} 0111 \\ 0011 \\ \hline 1010 \Rightarrow 01010 \Rightarrow (+10) \\ \text{OV} = 1 \end{array}$$

Eg

$$\begin{array}{r} +7 \\ -3 \\ \hline +4 \\ \text{OV} = 0 \end{array}$$

$$\begin{array}{r} 0111 \\ 1101 \\ \hline 0100 \Rightarrow (+4) \\ \text{OV} = 0 \end{array}$$

Eg

$$\begin{array}{r} -7 \\ +3 \\ \hline -4 \\ \text{OV} = 0 \end{array}$$

$$\begin{array}{r} 1001 \\ 0011 \\ \hline 1000 \Rightarrow (-4) \\ \text{OV} = 0 \end{array}$$

Eg

$$\begin{array}{r} -7 \\ -3 \\ \hline -10 \\ \text{OV} = 1 \end{array}$$

$$\begin{array}{r} 1001 \\ 1101 \\ \hline 10110 \Rightarrow (10110) \Rightarrow (-10) \\ \text{OV} = 1 \end{array}$$

$$OV(x, y, z) = \bar{x}\bar{y}z + xyz$$

$x$  - MSB of OP1,  $y$  - MSB of OP2,  $z$  - MSB of Result

condition of overflow set

Q. Consider the following 2 two's complement nos. and perform arithmetic op<sup>n</sup>. What could be status of the carry, sign, and overflow flags after the execution.

$$\begin{array}{r} 11001101 \\ 10110011 \\ \hline \boxed{1}10000000 \end{array}$$

$$\begin{aligned} \text{Carry} &= 1 \\ \text{Parity} &= 0 \\ \text{A.C} &= 1 \\ \text{Zero} &= 0 \\ \text{Sign} &= 1 \\ \text{overflow} &= 0 \end{aligned}$$

CONTROL FLAGS:- Based on the status of these flags, the program execution sequence is changed. Ins<sup>ns</sup> that changed by user. There are 3 control flags used:

- (i) Trap  $\begin{cases} \rightarrow 1 \rightarrow \text{Trace} \rightarrow \text{Single step execut}^n & (-t) \\ \rightarrow 0 \rightarrow \text{Go} \rightarrow \text{At a time prog}^n & (-g) \end{cases}$
- (ii) Interrupt  $\begin{cases} \rightarrow 1: \text{Enable interrupt} \\ \rightarrow 0: \text{Disable INTs} \end{cases}$
- (iii) Direction  $\begin{cases} \rightarrow 1: \text{Auto decrement (STD = Set Direct}^n \text{ flag)} \\ \rightarrow 0: \text{Auto increment (CLD = Clear " " )} \end{cases}$

ADDRESSING MODES:-

- It shows the way where the required object is present.
- Object may be an inst<sup>x</sup> / data.
- The o/p of the addressing mode is Effective address (EA).

- EA is the actual address where the required object is present  
i.e. object = content of EA = [EA]
- There are different conventions used to indicate different addressing modes i.e.
  - #/I = Indicates Immediate AM
  - Register Name = Register AM
  - [ ] = Direct AM
  - @/C ) = Indirect AM
  - Indexing register name = Indexed AM

#/I  $\Rightarrow$  Processor supports different funct<sup>ns</sup>. There are diff. AM used in computer. To use diff. AM, diff. user conventions. convenience codes are used during ~~exec~~ creating the program. These are above mentioned.

AM are basically classified into 2 types  
 (a) Sequential control flow AM (focussed on data)  
 (b) Transfer of control flow AM (focussed on Instr)

### SEQUENTIAL CONTROL FLOW AMs -

- These ~~AM~~ AMs are focussed on data.
- Data present in computer in 2 possible places i.e. processor register & memory.
  - $\therefore$  sequential control flow AM further divided into 2 types:
    - (i) Register based AM's: These are used only when the data is present in the register  $\therefore$  the EA is "register name".
    - (ii) Memory based AM's: These are used when the data present in the memory.  $\therefore$  EA is "Memory Address".

# Register based AM:-

## (1) Register AM:-

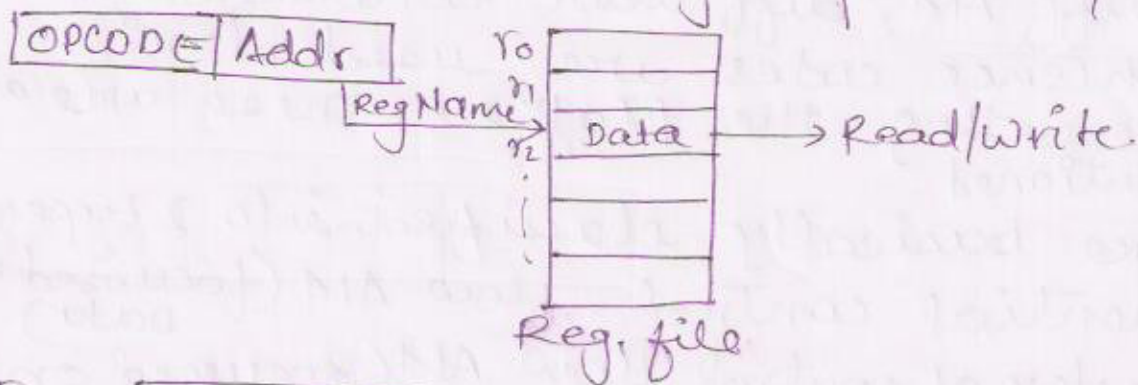
- This mode is used to access local variables.
- In this mode data is present in the register, the reg. name is available in the addr field of inst<sup>r</sup>.

∴ EA = Addr field of inst<sup>r</sup> value

↓  
Register Name.

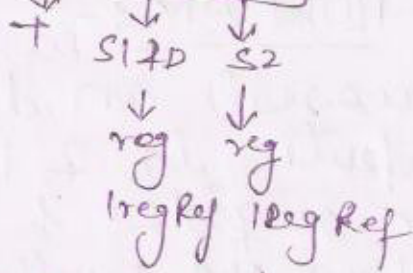
Data = [EA] = [Reg. Name]

↓ Reg. Ref → read/write.



Eg: 

ADD	R0	R4
-----	----	----



$R0 \leftarrow R0 + R4$   
 1 reg ref (write)    1 reg ref (read)    1 reg ref (read)

## (2) Memory based AM:-

There are diff. types of AM employed under this category.

(a) Implied / Implicit AM:- In this mode the data information is present in the opcode itself.

OPCODE
--------

Type of op<sup>n</sup>  
EA not required  

STC
-----

 Set carry  
C=1

operand info  

CLC
-----

 Clear carry  
C=0

ADD Stack M/C.

Here data is part of opcode.

POP  
POP  
+  
push.

NOTE; - All the zero address instr use the implied addressing mode.

(ii) Immediate:-

- This mode is used to access the constants.
- In this mode the data present in the address field of instr.

i.e. 

OPCODE	Address
--------	---------

↓  
Data  
Data part of instr.

Eg: 

ADD	R <sub>0</sub> , #23
-----	----------------------

↓      ↓      ↓  
+      S<sub>12D</sub>      S<sub>2</sub>  
         ↓      ↓  
         R<sub>reg</sub>    Imme.

$$R_0 \leftarrow R_0 + 23$$

NOTE:- The limitation of immediate AM is the range of possible const. are restricted based on the size of the address field. i.e. if the address field size of the instr is n-bit, then the possible range of unsigned constants are  $\{0 \text{ to } 2^n - 1\}$ , possible range of signed constants are  $\{-(2^{n-1}) \text{ to } +(2^{n-1} - 1)\}$ .

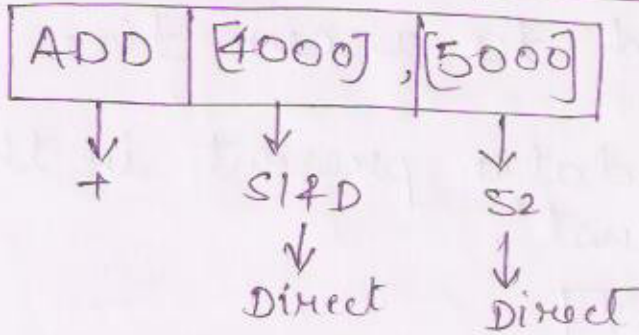
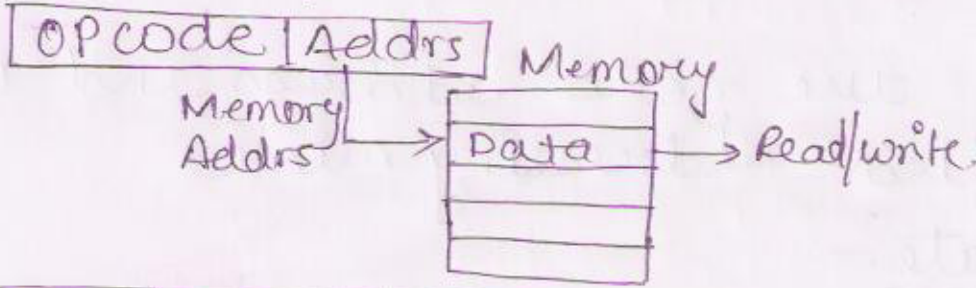
(iii) Direct AM (Absolute AM):-

- This AM is used to access static variables.
- In this mode data present in the memory, that memory address available in the address field of the instr.

$$\therefore EA = \text{Address field value} \Rightarrow \text{Mem. Addr.}$$

Data = [EA] = [Memory Addr]

- 1 Mem. ~~ref~~ reference is req. to read/write data to/from the memory by using this AM.



$$M[4000] \leftarrow M[4000] + M[5000]$$

1 Mem Ref      1 Mem Ref      1 Mem Ref  
 (3 Mem. Ref)

`ADD [4000], R0`  $\Rightarrow$  (2 Mem Ref)

`MOV [5000], [4000]`  $\Rightarrow$  (2 Mem Ref)

(iv) Indirect AM :-

- This is used to implement the pointer concept.
- In the indirect AM, data present in memory that memory Addr [EA] may be present in register/memory.

$\therefore$  Based on the availability of EA the indirect AM is further divided into two types :-

(iv)(a) Register indirect AM :- In this mode, EA is present in register, that reg. name is available in the addr field of the instr.

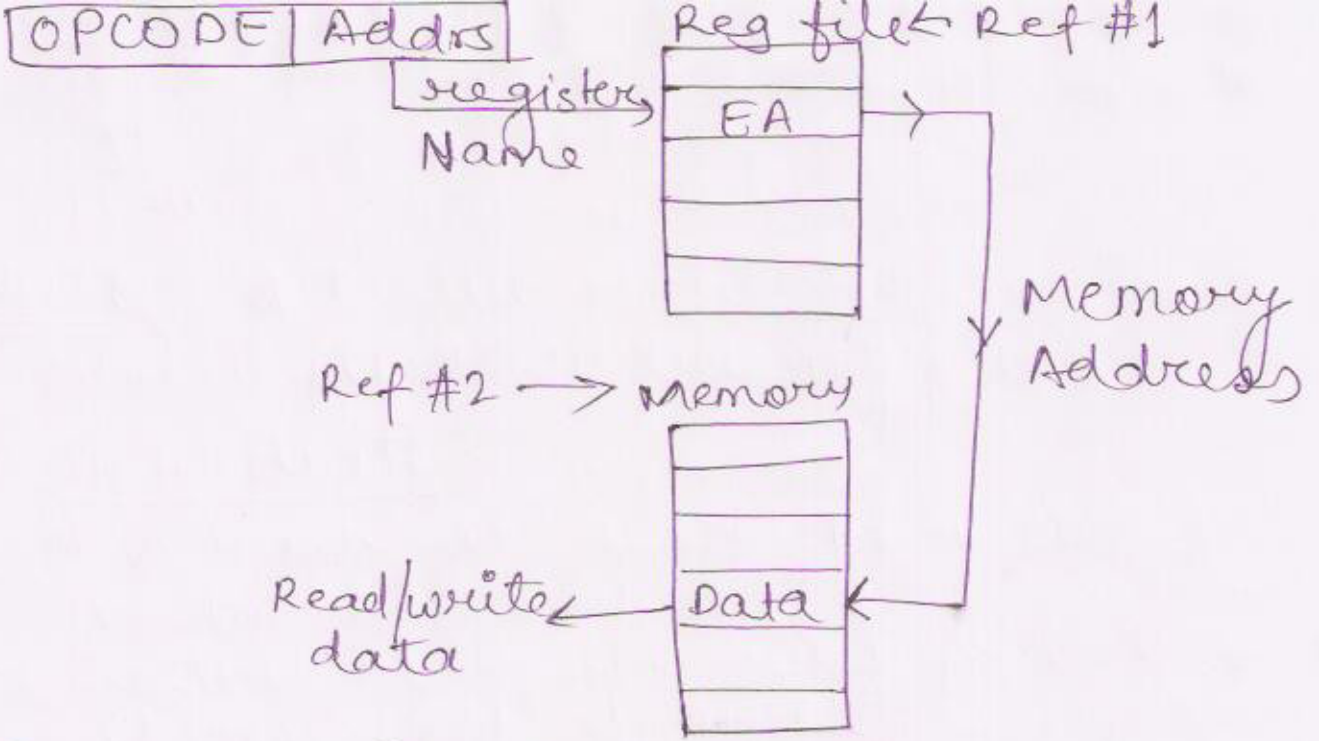
$\therefore$  EA = [Addr. field value]

↓  
Reg. Name

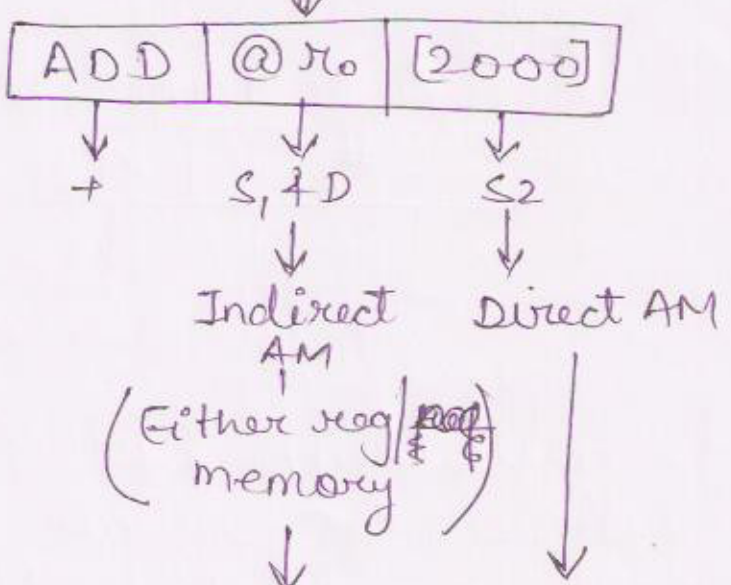
Data = [EA]  
= [Reg. Name]

2 Ref :-

- 1 Reg. ref  $\rightarrow$  EA
- 1 Mem Ref  $\rightarrow$  Read/write



Eg:- ADD @r0, [2000]



$$M[r_0] \leftarrow M[r_0] + M[2000]$$

1 reg. ref = EA      1 reg. ref = EA      1 mem. ref (read)  
 1 mem. ref = write      1 mem. ref = read

(3 references required)





### ③ Memory Indirect Addressing Mode:-

In this mode the effective address is present in the memory, that memory address is available in the address field of the instruction.

$$EA = [\text{Address field value}]$$

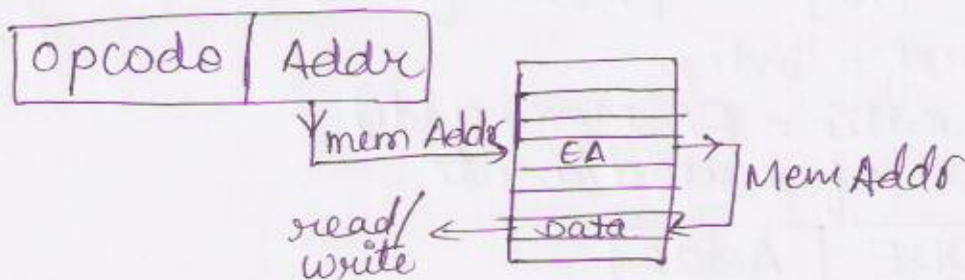
↓  
Mem. addr.

Register indirect is faster.

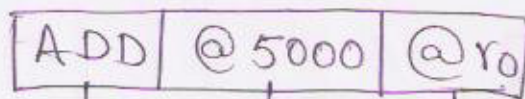
$$\text{data} = [EA] = [[\text{mem Addr}]]$$

$$\text{Mem. refer} = EA$$

$$\frac{1 \text{ mem. ref} = \text{read/write}}{2 \text{ mem ref.}}$$



Eg 1:



↓                      ↓                      ↓

+                      STPD                      S2

↓                      ↓                      ↓

Indirect                      Indirect

2+2                      1

⇒ 5 mem. ref.

$$M[\underset{2}{[5000]}] \leftarrow M[\underset{2}{[5000]}] + M[\underset{1}{[r0]}]$$

### Indexed Addressing Mode

It is used to implement the arrays. Array is an ordered set of homogeneous data elements. Array is always stored in the sequential memory locations.

To access the array element, there is a need of 2 parameters i.e. (i) Base address (ii) Index value.

In the indexed mode, EA is calculated by adding the index value to the base address, i.e.  $EA = BA + \text{index value}$ .

BA  $\rightarrow$  starting or ending address of the array.

Index  $\rightarrow$  Position of the array element.

Processor supports one special purpose register called as Index register ( $R_i$ ). It is used to hold the index value.

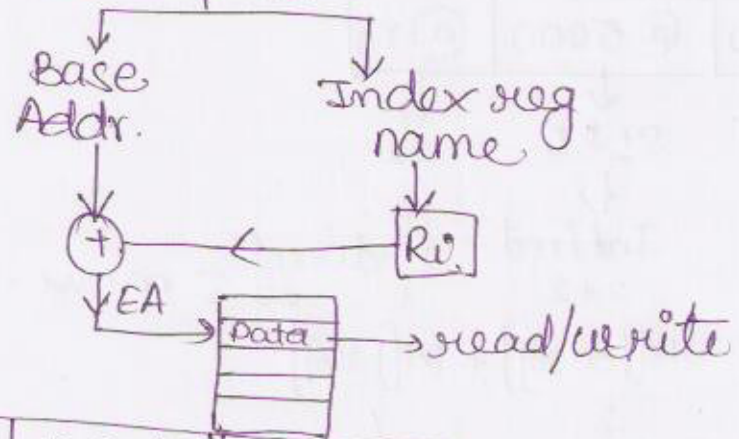
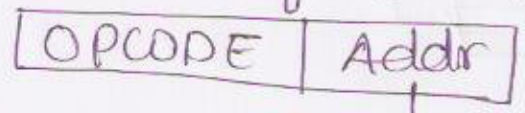
Base address is present in the address field of instr.  $\therefore EA = \text{Addr. field value} + [R_i]$

$$\text{Data} = [EA] = [\text{Addr field value} + [R_i]]$$

1 reg ref = Index

1 arithmetic = EA (Mem. addr)

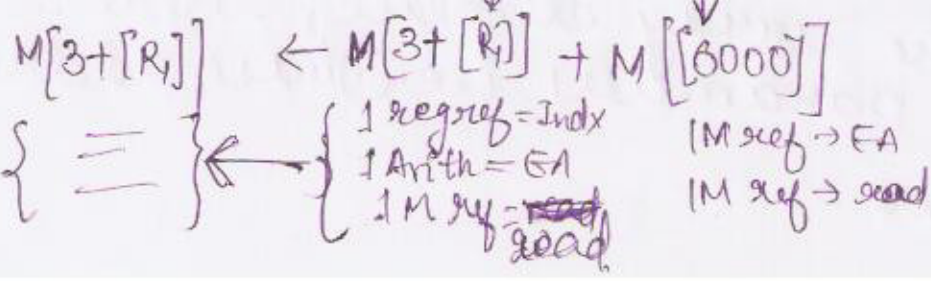
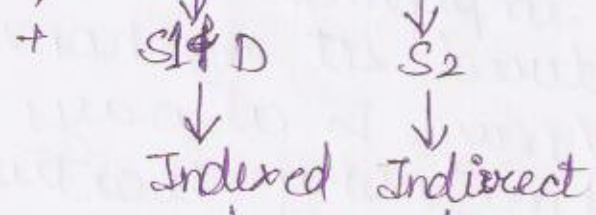
1 mem. ref = read/write.



Eg 

Add	3( $r_1$ ), @6000
-----	-------------------

 $r_1 = \text{Index reg}$



Indexed addressing mode is also called as based indexed addressing mode.

NOTE:- Some of the processors support base and index register used to hold the base and index value respectively.

Eg:- 8086  $\mu$ P supports 2 registers named as  
 BX  $\rightarrow$  Base reg,  
 SI/DI  $\rightarrow$  Index reg. (Source/Destination Index)

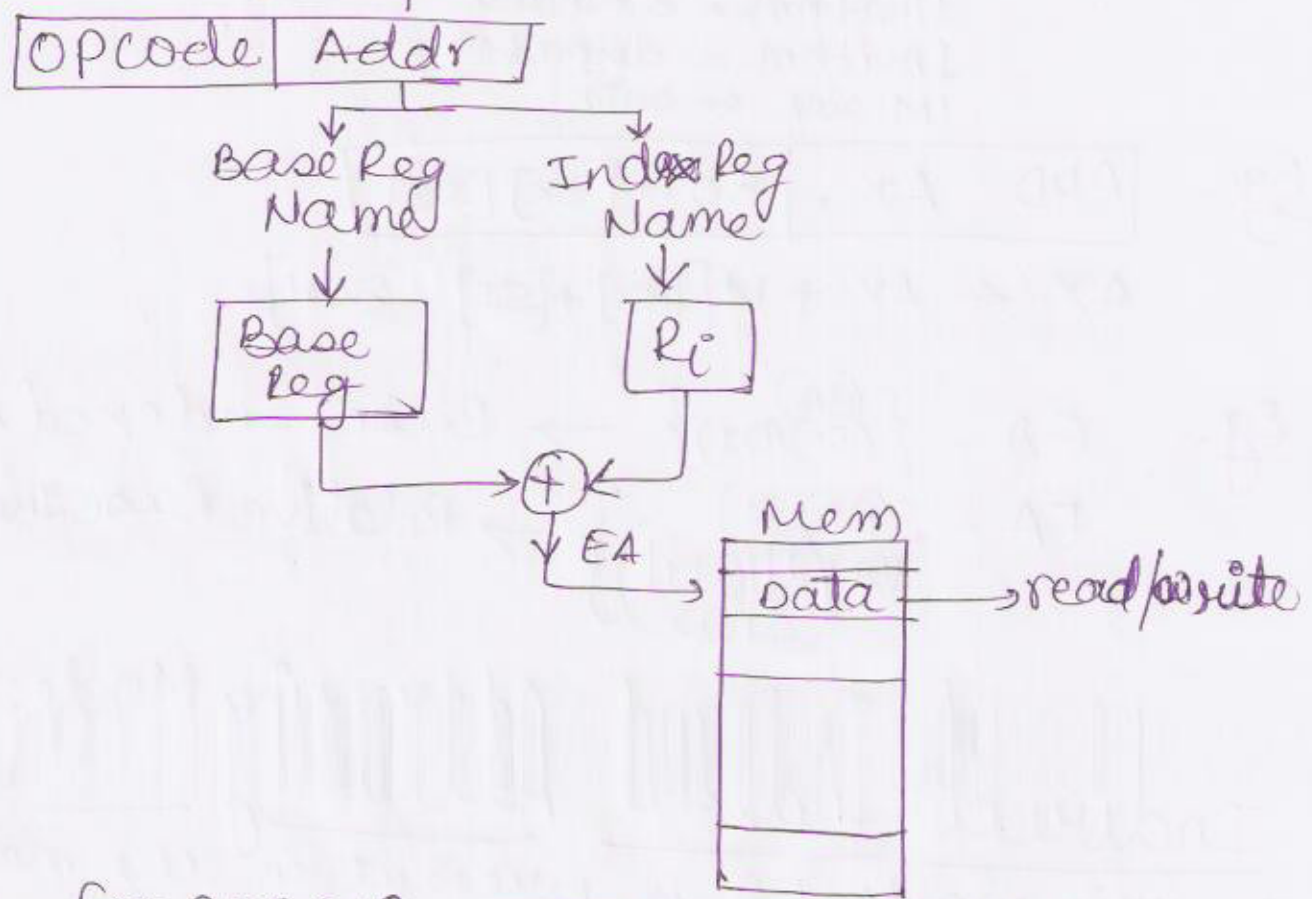
$$EA = [Base\ Reg] + [Index\ reg]$$

$\downarrow$  Base Addr                       $\downarrow$  Index value

$$Data = [EA]$$

$$= [[Base\ Reg] + [Index\ reg]]$$

1 reg ref = Base Addr  
 1 reg ref = Index  
 1 Arithm. = EA  
 1 M ref = read/write



Eg:- For 8086  $\mu$ P:

ADD	AX,	[BX][SI]
$\downarrow$	$\downarrow$	$\downarrow$
+	SI, DI	S <sub>2</sub>



$$\text{Data} = [EA]$$

$$= \left[ [ \text{Addr. field} + [R_i] ] \right] / \left[ [ [ \text{Base Reg} ] + [R_i] ] \right] \quad (\text{or})$$

1 Reg ref = Index

1 Arith. = Mem Addr where  
EA is present

1 Mem. Ref = EA

1 Mem Ref = Read/Write

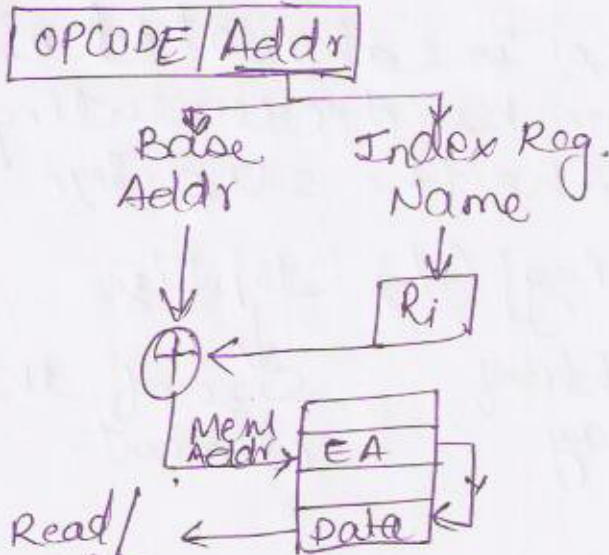
1 reg ref = Base Reg.

1 reg ref = Index Val.

1 Arith = Mem Addr,  
where the EA is  
present

1 Mem. ref = EA

1 Mem ref = Read/Write



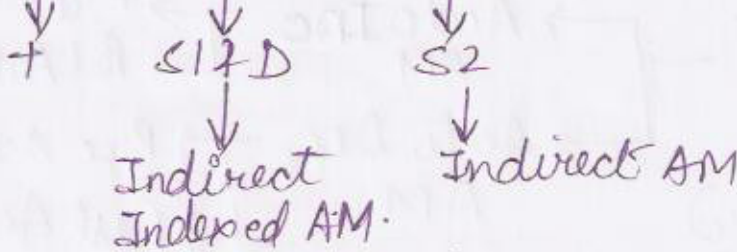
Not in use becoz  
it takes more time

Indirect

Indexed Addressing mode takes more # cycles  
to read/write data to/from memory.

Eg: 

ADD	3(R4),	@6000
-----	--------	-------



$$M[3 + [R_4]] \leftarrow M[3 + [R_4]] + M[6000]$$

1 reg ref = Index    1 reg → R<sub>i</sub> (Index)    1 Mem → EA  
1 Arith = Mem Addr    1 Arith → Mem Addr    1 Mem → read  
1 Mem. ref = EA    1 Mem ref → EA  
1 Mem ref = Write    1 Mem ref → read

NOTE: - Indexed addressing mode is used to  
access the random array element based  
on the index value.

Auto Indexed AM. It is used to access the linear array elements. To access the linear array elements, there is a need of base address.

- Base address ~~is~~ ~~ba~~ indicates the starting or ending of an array.
- Base reg. is used to hold the base addr. In this mode EA is calculated either by incrementing or by decrementing the base address with the step size. i.e.

$$EA = [\text{Base Reg}] (+/-) \text{Step Size}$$

↓  
Base addr of  
an array

↓  
Size of the data  
element.

$$\text{Data} = [EA]$$

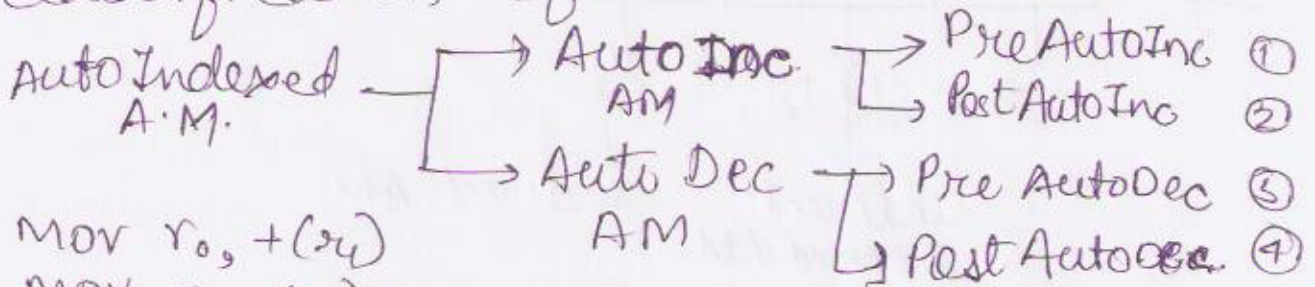
$$= [[\text{Base reg}] (+/-) \text{Step Size}]$$

1 reg. ref. = Base Addr.

1 with m = EA

1 M. ref = read/write

The classification of auto Indexed AM is



1.  $\text{MOV } r_0, +(r_1)$

2.  $\text{MOV } r_0, (r_1) +$

3.  $\text{MOV } r_0, -(r_1)$

4.  $\text{MOV } r_0, (r_1) -$

where R = Base register.

Q. Consider a hypothetical CPU. It uses different operands accessing modes.

operand & accessing mode.	Frequency (%)
Register	30
Immediate	20
Direct	22

Mem. Indirect = 17

Indexed = 11.

Assume the 2 cycles required for memory access, one cycle consumes for arithmetic operations & zero cycles consumes when the operand is coming in the register or instruction itself. What is the average operand fetch rate of the processor?

Ans :-

$$\begin{aligned} \text{Average operand fetch rate} &= (0.3 \times 2) + (0.2 \times 0) + (0.22 \times 2) + (0.17 \times 4) \\ &\quad + (0.11 \times 3) \\ &= 0.44 + 0.68 + 0.33 \\ &= 1.45 \text{ cycles} \end{aligned}$$

Transfer of control flow AM :-

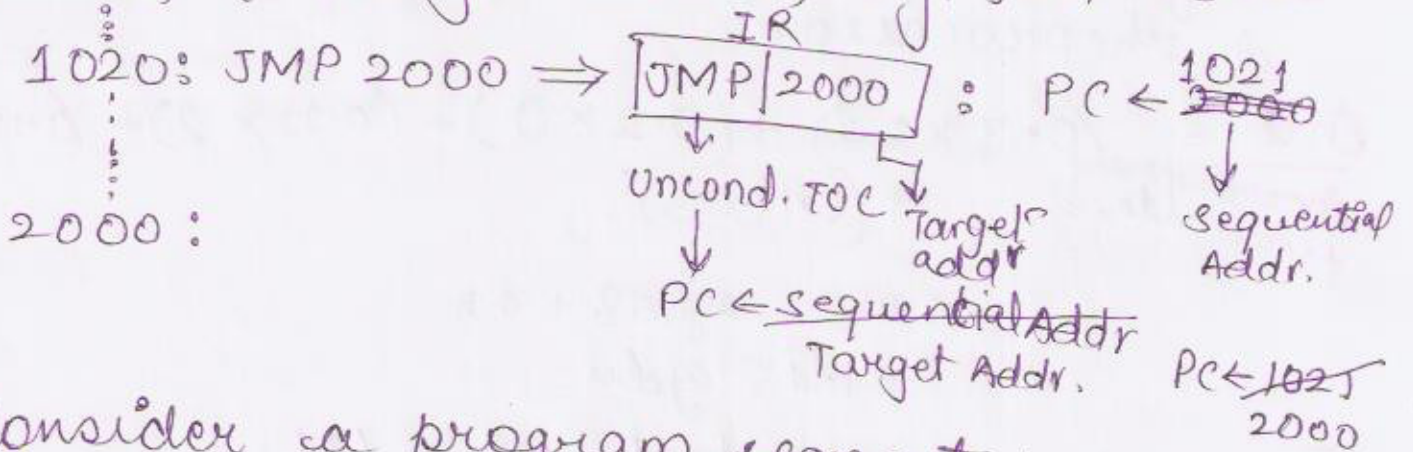
While execution of control structures i.e. selection statements (goto, if-then-else, switch etc), iterative statements (while, do-while, for loop etc) and subprogram concept, the program control is transferred from 1 location to another location.

- To implement the control structures, process or supports special type of instructions called as transfer of control flow.
- Transfer of control is associated with the target address.
- While execution of the TOC operation, the control is transferred from current location to target location.
- There are 3 possible mnemonics used to implement the TOC op<sup>n</sup>.
  - (i) Branch
  - (ii) Jump
  - (iii) Skip

TOC op<sup>n</sup> classified into 2 types

- (i) Unconditional TOC
- (ii) Conditional TOC

Unconditional TOC : While execution of these instr, without checking any cond<sup>n</sup>, the control will be transferred from current location to target location. Eg. JMP 2000

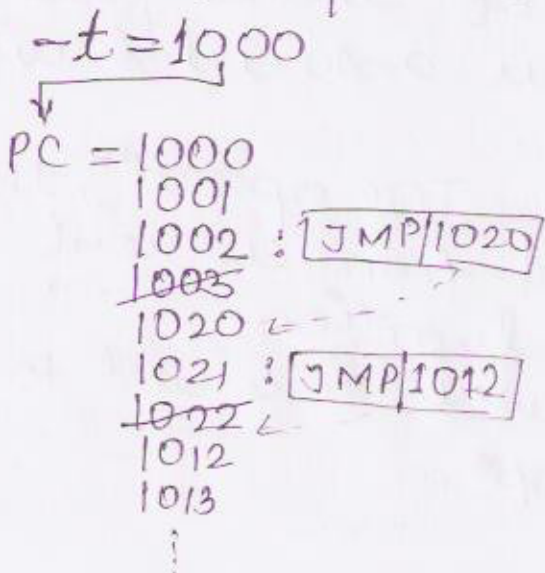


Consider a program segment:

```

1000: I1
1001: I2
1002: I3 (JMP 1020)
1003: I4
    ⋮
1020: B1
1021: B2 (JMP 1012)
1022: B3
    
```

Execution sequence of program is:-



NOTE: Uncond. TOC operations are used to implement the goto statement, halt instr and M/C control instr.



HALT: It is implemented using the unconditional TOC with starting addr as the target addr.

Eg: 1000: I<sub>1</sub>  
 1001: I<sub>2</sub>  
 1002: I<sub>3</sub> (Halt)  
 1003: I<sub>4</sub>

-t = 1000

PC = 1000

1001

1002: Halt: JMP 1002

1003

1002: halt

~~1003~~

1002: halt

M/C control inst<sup>r</sup>:-

Uncond<sup>n</sup> TOC op<sup>n</sup> are used to control the functionality of some specific M/Cs.

Eg: Traffic light controller.

1: G: ON  
 Y: OFF  
 R: OFF

call delay( )

G: OFF  
 Y: ON  
 R: OFF

call delay1( )

G: OFF  
 Y: OFF  
 R: ON

call delay2( )

Imp 4.

Conditional TOC:-

Whole execution of this instruction, the associated cond<sup>n</sup> undergoes evaluation. When it evaluates to true, then the control is transferred to the target location, otherwise the next sequential inst<sup>r</sup> is executed.

Eg: JNZ 2000 PC seq. Addr. cond: NZ

Cond<sup>n</sup>  
TOC

target addr

NZ comp status of previous inst<sup>r</sup>  
 false → true: PC ← sequential  
 true → false: PC ← target Addr

Consider the program segment:

```
1000 : MOV r0, #4
1001 : ADD r1, r2
1002 : DEC r0
1003 : JNZ 1001
1004 : MOV
      :
```

Execut<sup>n</sup> Seq:-

-t = 1000

↳ PC : 1000 : r0 = 4  
1001 : +  
1002 : r0 = 3 (NZ)  
1003 : NZ (True)  
~~1004~~  
1001 : +  
1002 : r0 = 2

1003 : NZ (True)

~~1004~~

1001 = +

1002 = r0 = 1

1003 = NZ (True)

~~1004~~

1001 = +

1002 = r0 = 0

1003 = Zero (false)

1004 = MOV.

1005 :

NOTE:- Cond<sup>nal</sup> TOC oper<sup>ns</sup> are used to implement the cond<sup>nal</sup> select<sup>n</sup> statmnts & iterative ~~select<sup>n</sup>~~ stmnts.

Subprogram:- Subprogram is a reusable program. That means when the same functionality is repeatedly occurred in the main app<sup>n</sup>, develop the functionality related code only once (subprogram) and call the functionality many times in the main app<sup>n</sup>.

The advtage of subprogram is it reduces the memory space, development time & development cost.

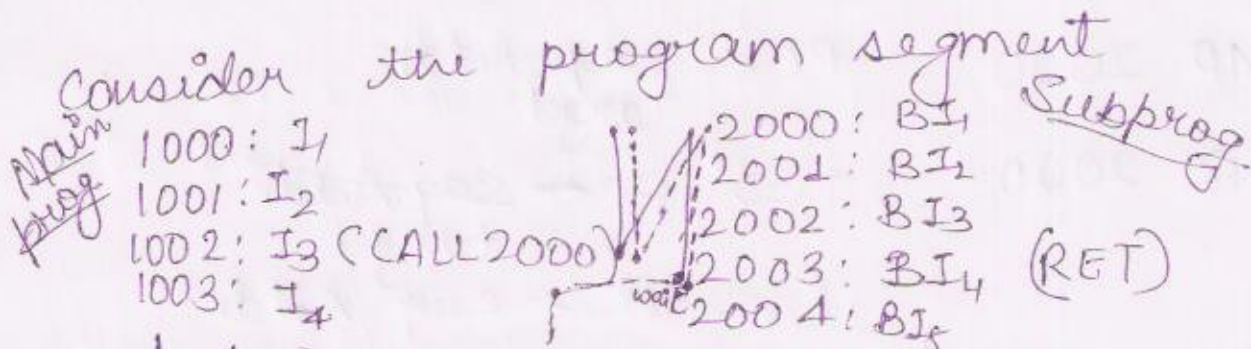
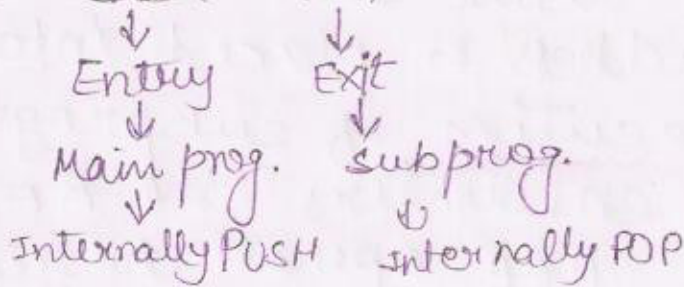
## Characteristics of subprogram:-

- (i) Single entry point - single exit points.
- (ii) Main program is suspended during the execution of the subprogram.
- (iii) Control will be transferred to the main program after completion of the subprogram.

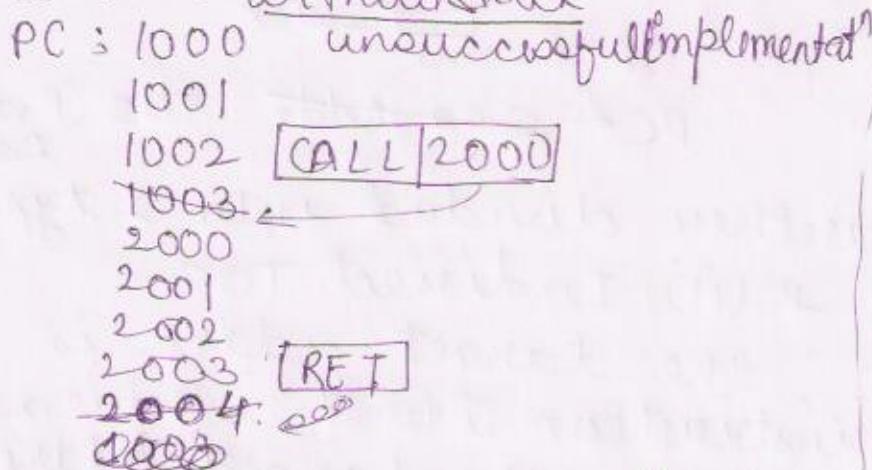
### IMPLEMENTATION:-

To implement the subprogram concept, the processor supports pair instructions.

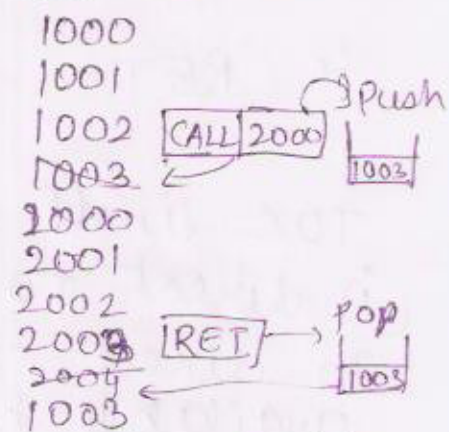
i.e. CALL-RET



$t=1000$  without stack



with stack



In the above execution sequence, due to the lack of return statements, RET address during the execution of RET instr, CPU enters into the WAIT state. ∴ It is not

- a successful implementation of subprogram.
- To make it successful, need of return address.
- Return address is Next inst<sup>r</sup> address after the CALL inst<sup>r</sup>.
- Runtime stack is used to store the return addresses.
- During the execut<sup>n</sup> of program, when the CPU encounters the CALL inst<sup>r</sup>, it pushes the PC value into the stack. Later target address is placed into PC.
- During the execution of subprogram, when the CPU encounters the RET inst<sup>r</sup> it invokes the pop operat<sup>n</sup> to restore PC with return address.

1. JMP 2000       $PC \leftarrow \begin{matrix} \text{Seq. Addr} \\ \text{Target} \end{matrix}$

2. JNZ 2000       $T \Rightarrow PC \leftarrow \begin{matrix} \text{Seq. Addr} \\ \text{target} \end{matrix}$   
 $F \Rightarrow PC \leftarrow \text{Seq. Addr.}$

3. CALL 2000 :  $PC \xrightarrow{\downarrow} \text{stack}$   
 $PC \leftarrow 2000$

4. RET       $PC \leftarrow \text{Seq Addr TOS (Returning Addr)}$

TOC inst<sup>r</sup> further divided into 2 types:-

(i) direct TOC & (ii) Indirect TOC.

In direct TOC, the target address is available in instruction itself, whereas in the indirect TOC, the target address is not present in the inst<sup>r</sup>.

TOC	uncond.	cond.
Direct	JMP, CALL, goto Halt	JNZ, CALLZ, CALLNZ, CALLC, for, while, do-while, if, switch
Indirect	RET	RETZ, RETC, RETNC, RETNZ

In the transfer of control flow AM category, two addressing mode are employed:-

(a) Relative / PC relative AM

(b) Based / Base Reg AM.

→ Inst are always stored in memory. So to access the instruct's efficiently, there is a need of understanding the memory org.

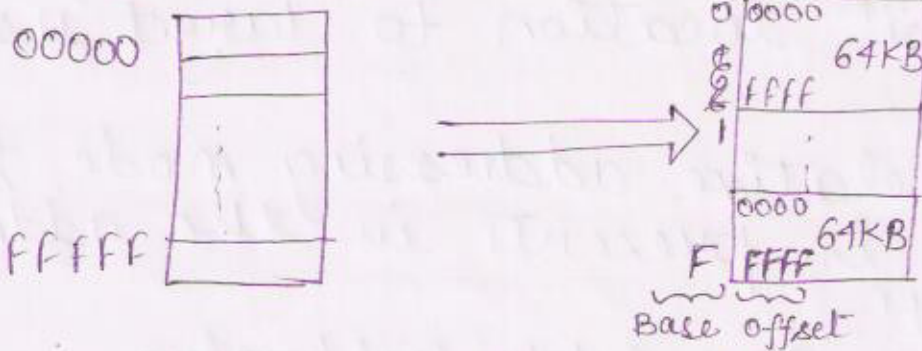
→ Consider 8086 mem org. It contains 1MB physical memory. The physical memory is divided into 16 logical segments with each segment size of 64KB.

(No change in size)

1MB  $\Rightarrow 2^{20}$  cells.

1MB = 16 x 64KB =  $2^{20}$  B.

change in organization



Physical addr = BA + Offset Addr.

In 8086, segment registers are used to hold BA and pointer or index register is used to hold offset address.

Segment Registers  $\rightarrow$  CS  
 $\rightarrow$  DS  
 $\rightarrow$  SS  
 $\rightarrow$  ES

Pointer/Index Reg  $\rightarrow$  IP } Pointer Reg  
 $\rightarrow$  BP }  
 $\rightarrow$  SP }  
 $\rightarrow$  SI } Index Reg  
 $\rightarrow$  DI }

If  $CS=0 \Rightarrow \left. \begin{matrix} 00000 \\ 0FFFF \end{matrix} \right\} \text{Inst}^n$

If  $DS=2 \Rightarrow \left. \begin{matrix} 20000 \\ 2FFFF \end{matrix} \right\} \text{Data}$

To access code  $\Rightarrow CS + IP$

Notation

$CS:IP$

To access data  $\Rightarrow DS + BP/SI/DI$

$DS:BP/SI/DI$

Relative/PC-relative AM:-

When the target inst<sup>r</sup> is present in the same segment, then the control will be transferred within the segment (Intrasegment TOC)

In the intrasegment TOC the base address is common. So only 1 parameter is required to calculate the target inst<sup>r</sup> addr. i.e. relative value.

Relative value is the distance b/w the current inst<sup>r</sup> location to target inst<sup>r</sup> location.

In the PC-relative addressing mode, the relative value is present in the address field of inst<sup>r</sup>.

$$\begin{array}{ccc} \text{Next inst}^r & \text{Current} & \text{relative} \\ \text{Addr} & \text{Inst}^r \text{ Addr} & \text{value} \\ \downarrow & \downarrow & \downarrow \\ \text{EA} = \text{PC} + \text{Addr. field value.} & & \end{array}$$

$$\text{PC} \leftarrow \text{PC} + \text{Addr. field value.}$$

PC relative addressing mode is suitable for program relocation at runtime.



# INSTRUCTION SETS :-

Processor supports 3 categories of ins<sup>n</sup>.

- (i) Data Transfer Ins<sup>n</sup>
- (ii) Data Manipulation Ins<sup>n</sup>
- (iii) Transfer of control Ins<sup>n</sup>

## (i) Data Transfer Instruction

While execution of this instr, data is transferred from 1 location to another location. The limitation of data transfer operation is source may or may not have the storage functionality but the dest<sup>n</sup> always has the storage.

Immediate addressing mode may be for

There may be different data transfer operations possible.

(i) MOV :- General purpose data transfer instruction.

<u>Source</u>	<u>Dest<sup>n</sup></u>
reg	reg
reg	Mem
Mem	reg
Mem	Mem
Immediate	reg
Immediate	Mem.

(2) LOAD } These are memory specific op<sup>n</sup> used to perform the memory read & memory write op<sup>n</sup> resp.

(3) STORE } These op<sup>n</sup> are stack specific op<sup>n</sup> used to perform the insert & delete operations resp.

(4) IN } No specific op<sup>n</sup> used to perform the thread & IO write op<sup>n</sup> resp.

(5) OUT }





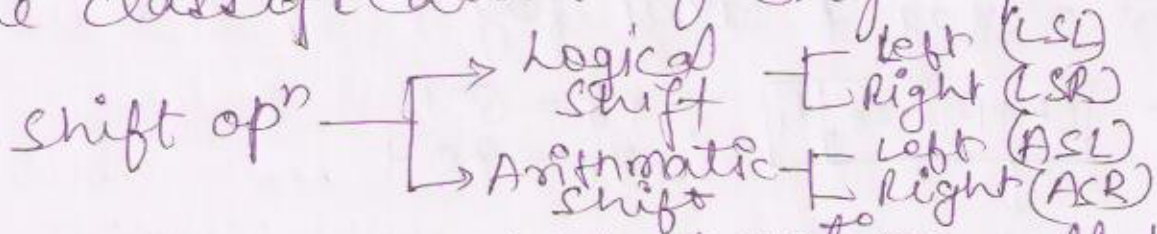
Eg 3:  $r_0 = 40H$   $r_1 = 80H$

$$\begin{array}{r} 01000000 \\ 10000000 \\ \hline 11000000 \end{array}$$
 (0) zero-reset  
 set (1) carry

After CMP execution, if zero flag is reset & carry flag is set, that means 2nd operand is greater than 1st operand.

Shift operation:-

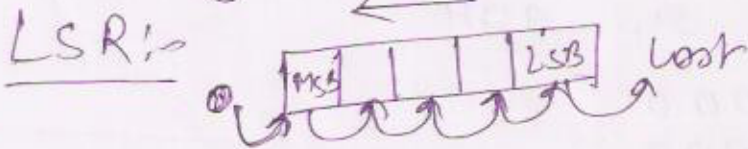
→ while execution of shift operation, data is transferred towards left or right bit by bit with loss of data.  
 → The classification of shift operation is



In the logical shift operation all the bits are moving towards left or right including the sign bit.



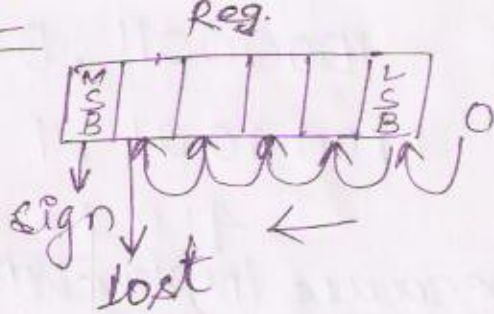
$r_0 = 83H$   
 Eg. LSL  $r_0, \#2$   
 $1000\ 0011$   
 ← 2 times  
 Ans. →  $0000\ 1100$   
 0CH



Eg. LSR  $r_0, r_0, r_0 = 83H$   
 $1000\ 0011$

Appl<sup>n</sup> :- Transfer data into another register. Used in the serial data comm<sup>n</sup>. to transfer the data bit by bit.  
 In the arithmetic shift operations, all the bits are moving towards left or right except the sign bit.

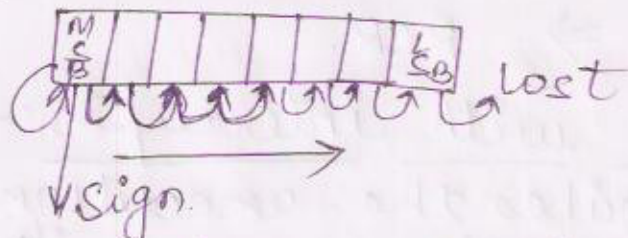
ASL



Eg:-  $\boxed{ACL} \times 2$   $r_0 = 83H$

$1000\ 0011$   
 $\rightarrow 1000\ 0110$   
 $86H$

ASR



Eg:-  $r_0 = 83$

$1000\ 0011$   
 $\rightarrow 1100\ 0001$   
 $(C1)H$

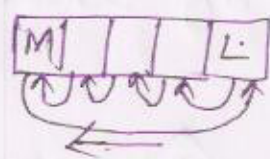
App<sup>n</sup>:- Used in signed arithmetic.

Rotate Instruction:- While execution of rotate Instr, all the bits are moving towards left or right without loss.

The classification of rotate operations is

- without carry —  $\left[ \begin{array}{l} \text{Left (ROL)} \\ \text{Right (ROR)} \end{array} \right.$
- with carry —  $\left[ \begin{array}{l} \text{Left (RCL)} \\ \text{Right (RCR)} \end{array} \right.$

ROL:-

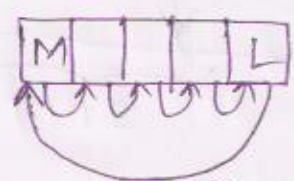


No new data is added  
Existing data is not lost,  
only position is changed.

Eg:-  $83H$

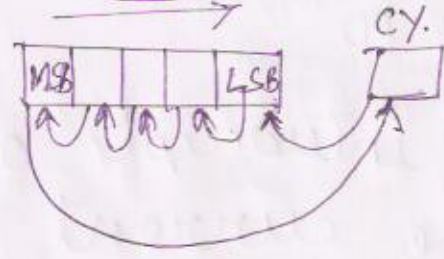
$1000\ 0011 \rightarrow 0000\ 0111$   
 $(07H)$

ROR:-



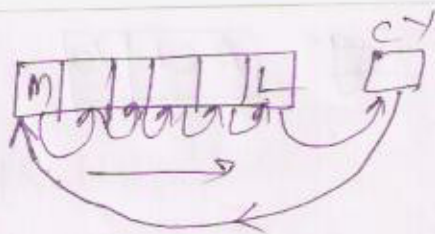
$1000\ 0011 \rightarrow 1100\ 0011$   
 $(C1)H$

RCL:-



$1000\ 0001\ C=0$   
 $\Rightarrow 0000\ 0110\ C=1$   
 $06H\ \&C \rightarrow 1$

ROR:-



10000011, 0  
01000001, 1

41H PC 71.

App<sup>n</sup>:- optimization because implement<sup>n</sup> is simple.  
Motor movement  $\Rightarrow$  ROL or ROR,  
even odd  $\Rightarrow$  RCR.

Instruction cycle with interrupts:-

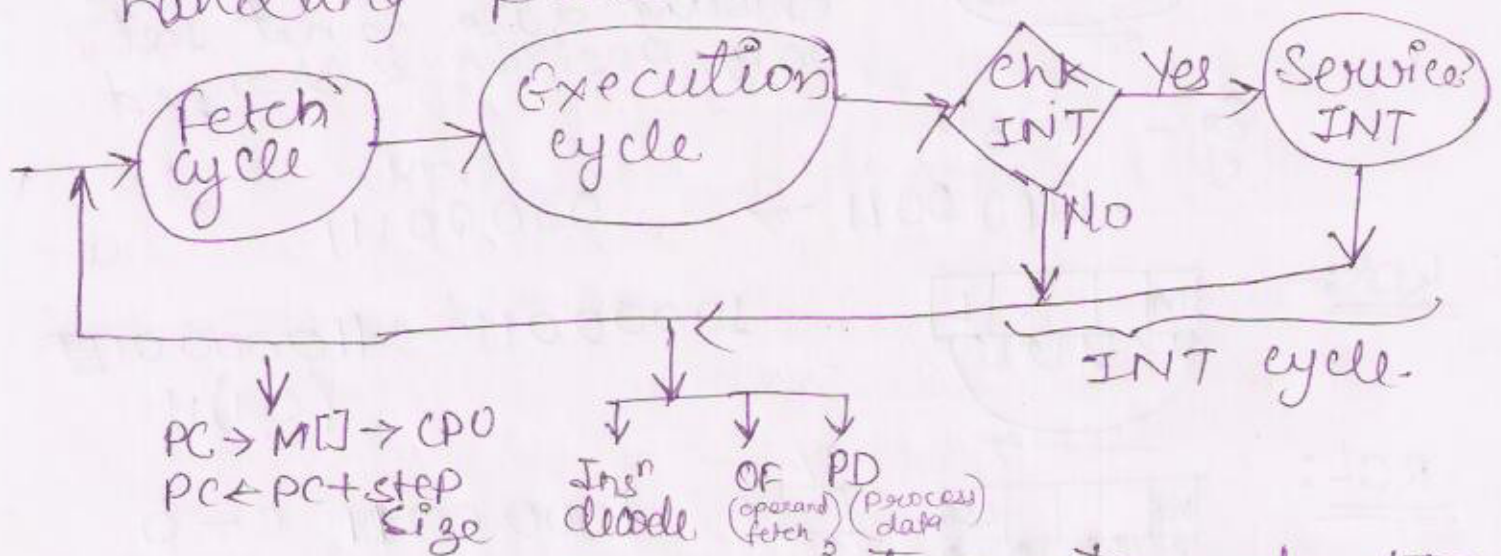
This concept describes the execution sequence of instructions along with the interrupts.

Interrupt is an unusual event to disturb the normal flow of execution.

Instruction cycle with interrupt consists of 3 subcycles:

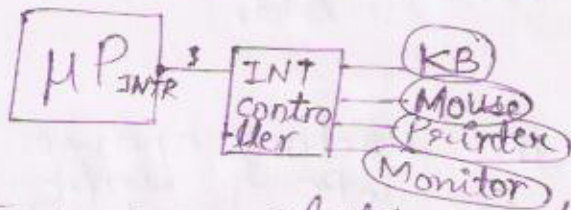
- (i) fetch cycle
- (ii) Execution cycle
- (iii) Interrupt cycle

The sequence diagram of interrupt handling process is



The CPU responds to interrupts only after completion of the current instruction execution.

When the interrupt flag is ENABLED that means interrupt sources are enabled. So interrupt cycle is required to service the interrupts. Otherwise, interrupt sources are disabled, so no need of interrupt cycle.



After completion of every inst<sup>r</sup> execution the CPU reads the status of the interrupt sources to identify whether the interrupt is present or not.

INTR = 1 ; INT is present

INTR = 0 ; INT is not present.

- When the int. are not present, then the CPU fetch the next inst. from the memory.
- When the int. is present, the CPU push the program counter value into the stack and transfers the control to the INTERRUPT VECTOR TABLE.
- IVT is a part of memory. It is used to store the interrupt subprograms.
- There are multiple interrupts possible in the processor. So different INT subprograms are required to service the different INT.
- Each INT subprogram must be end with IRET/RETI instruction.
- During execution of INT subprogram, when the CPU encounters the RETI inst<sup>r</sup> it invokes the POP inst<sup>r</sup> to restore the PC with return address.

∴ After service the interrupt, the CPU fetch the next instr based on the PC:

PC:	Data Transfer	ALU	Uncond TOC	Cond TOC
Fetch cycle	PC ← Seq. Addr.	PC ← Seq. Addr.	PC ← Seq.	PC ← Seq.
Execution cycle	No change in PC	No change in PC	change in PC PC ← target Addr.	change or Not PC ← changed or Not.
Int Service	Yes RA: Seq. Addr. return Addr.	Yes RA: Seq. Addr.	yes RA: target Addr.	Yes RA: Changed or Same.

**4D** Types of interrupts: [Mainly → H/W + S/W INT rest are part of these]

1. Hardware INT: - Hardware int. are present at the H/W pins of the processor. H/W INTs are also called as the ext. int. & internal int.

External INTs are generated by the ext. sources i.e. Power supply, basic IO devices etc.

Internal INTs generated by the internal components of the processor i.e. invalid OPCODE, divide by zero etc.

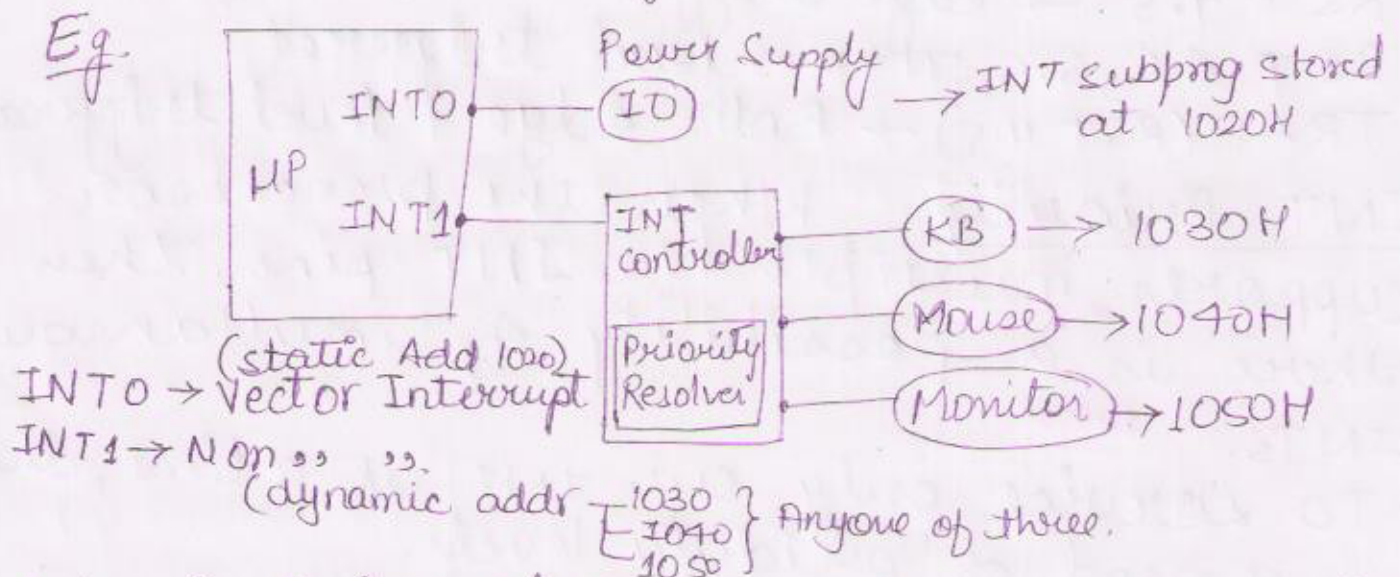
2. Software INT: - S/W INT is an instr. It is defined in the instr set of the processor. Eg: - Supervisor call.

3. Maskable INT: - Maskable INT pins are enabled or disabled based on the signals. All the basic IO devices are connected to maskable INT pins.

4. Non Maskable INT: - Here, the NMI pin is always present in the enabled state. So critical conditions of the processor i.e. power failure & temp sensor etc. are connected to the NMI pin.

5. Vector Interrupt: - Vector INT is associated with fixed vector point (Static vector Addr)

6. Non Vector INT: - It is associated with dynamic vector point (means vector point addr is supplied by the source).



7. Level triggered INT: - This INT is enabled based on the clock signals

8. Edge triggered INT: - These INT are enabled either with raising edge transition or with falling edge transition.

INT structure: - To analyse the implementation of various INTs, let us consider 8085 processor as a reference component.

8085 H/W INT

- |                    |                                  |
|--------------------|----------------------------------|
| 1) RST 7.5 (003CH) | (4) RST 4.5 (TRAP) (NMI) (0024H) |
| 2) RST 6.5 (0034H) | (5) INTR                         |
| 3) RST 5.5 (002CH) |                                  |

1, 2, 3, 5 → Maskable

# S/W INT - RSTO-7.

or S/W INTs, ~~maskable~~ maskable and non-maskable is not applicable because instr<sup>n</sup> cannot be enabled or disabled.

INTR - Non vectored INT.

Other 4 h/w INTs are vectored.

All the S/W INT are vectored.

RSTO - 0000    RST2 - 0010    RST4 - 0020    RST6 - 0030

RST1 - 0008    RST3 - 0018    RST5 - 0028    RST7 - 002F

For S/W INTs, level triggering and edge triggering is not applicable.

RST 7.5 → edge triggered

RST 6.5, 5.5, INTR → level triggered

TRAP (RST 4.5) → Both edge & level triggered

INT Priority: When the processor supports multiple h/w INT pins, then there is a possibility of simultaneous

INTs.

- To service only one INT at a time, there is a need of priority levels.

- |                   |                             |
|-------------------|-----------------------------|
| 1. TRAP (RST 4.5) | ↓ High<br>Priority<br>↓ Low |
| 2. RST 7.5        |                             |
| 3. RST 6.5        |                             |
| 4. RST 5.5        |                             |
| 5. INTR           |                             |

Q. A certain  $\mu P$  supports 4 INTs,  $I_1, I_2, I_3, I_4$ .  $I_1$  → priority decreases. The CPU responds to INT in every 2.5  $\mu sec$ . The service time of INTs are 20  $\mu sec$ , 40  $\mu sec$ , 30  $\mu sec$ , 25  $\mu sec$  resp. The INTs may or maynot occur simultan. What is the possible range of time required to execute the INT4.

⇒ Min ⇒ 27.5  $\mu s$     Max - 105 + 10 = 125  $\mu s$ .



- (\*) If many INTs are present simultaneously, after servicing one INT, an instr is executed. Only after execution of an instr, next INT is serviced.

Ans to Ques:-

min  $\rightarrow$  (without simul)  $\rightarrow$  2.5  $\mu$ s. with simul  $\rightarrow$  2.5  $\mu$ s.  
 max  $\rightarrow$  125  $\mu$ s. with simul  $\rightarrow$  2.5  $\mu$ s.

### RISC and CISC:-

Chars. of RISC  $\rightarrow$

- (i) It supports more no. of registers
- (ii) It supports less no. of addressing modes.
- (iii) It supports fixed length instructions.
- (iv) One instr per cycle - CPI = 1  
(cycles per instr)
- (v) It supports successful pipeline implementation.
- (vi) Motorola, PowerPC, ARM (Advd. risc M/c) etc are eg. of RISC processors.

Characteristics of CISC:-

- (i) It supports less no. of registers.
- (ii) It supports more no. of addressing modes.
- (iii) It supports variable length instrs.
- (iv) There is not one instr per cycle. CPI  $\neq$  1.
- (v) supports unsuccessful pipeline implementation.
- (vi) Eg - Pentium.

Register organization in RISC CPU:-

- In the RISC CPU, the registers are organized as overlapped register windows.
- The risc processor registers are classifi

-ed as

→ (i) Global reg (G) (iii) IN reg (C)

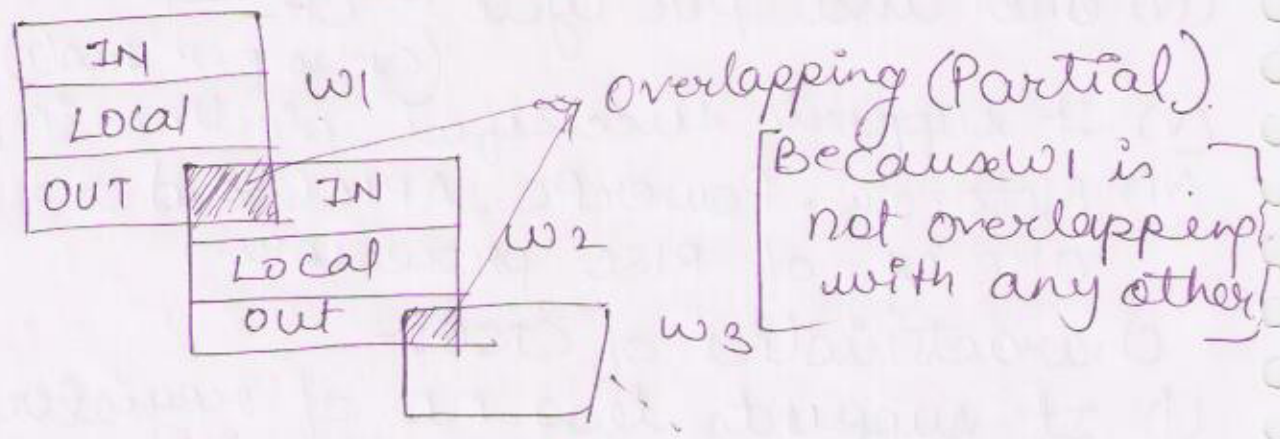
(ii) Local reg. (L) (iv) OUT reg. (C)

- The register window contains 3 set of registers i.e. local, IN and OUT reg.

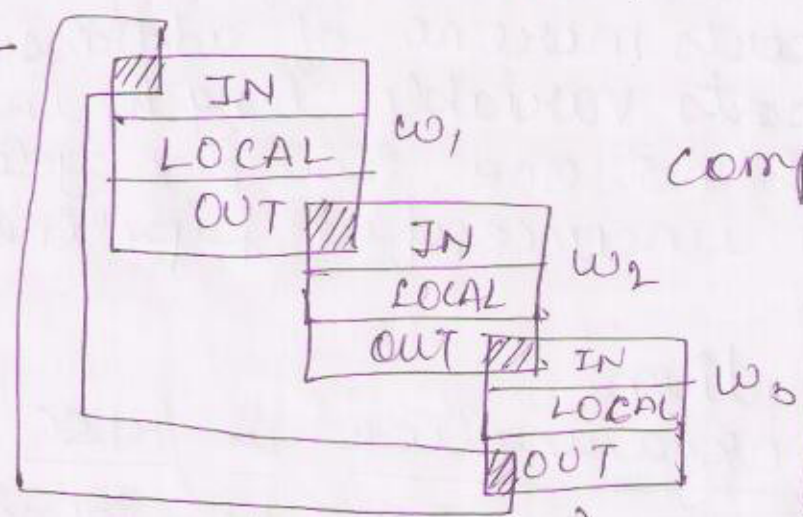
The global registers are accessible by all the register window.

- Window size =  $L + 2C + G$ .

In the RISC CPU, all the system supported registers are organized as overlapped register windows. Overlapping means 1 window OUT registers are used as IN registers in the another window.



Ex<sup>n</sup>:-



complete overlapping

Here,  
 $w_1(\text{out}) = w_2(\text{in})$   
 $w_2(\text{out}) = w_3(\text{in})$   
 $w_3(\text{out}) = w_1(\text{in})$

In overlapped (complete) register file size

$$= \text{No. of window} * (L + C) + G$$

# Components of computer:-

The computer contains 3 basic components.

- (i) CPU (ii) Memory (iii) I/O

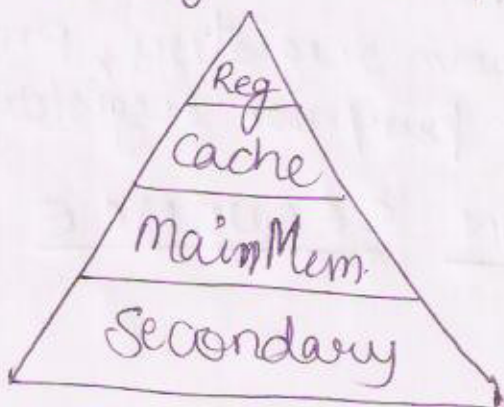
## CPU ORGANIZATION:-

CPU functionality is processing the inst<sup>r</sup>.

It consists 3 basic components (i) Registers (ii) ALU (iii) Control units.

REGISTERS:- Reg. are used to store the frequently referred data.

- During the execution of the program, there is a frequent communication between CPU and memory.
- Due to slow working of memory component, the CPU performance degrades.
- To avoid this problem, memory hierarchy design is used in the system design.
- The objective of memory hierarchy is to balance the bandwidth between the CPU and memory. (or) to minimise the speed gap between CPU + memory. (or) sync the data transfer rate b/w CPU + memory.
- Memory hierarchy describes the way of organization and way of accessing the system supported memories.



Bottom to top approach of memory hierarchy is

- (1) Capacity decreases.
- (2) Access time reduces.
- (3) Cost <sup>(Expensive)</sup> increases.
- (4)  $P \propto \frac{1}{\text{Exec. Time}}$ ,  $(E \propto \text{Acc. Time})$ ,  $(D \propto \frac{1}{AT})$

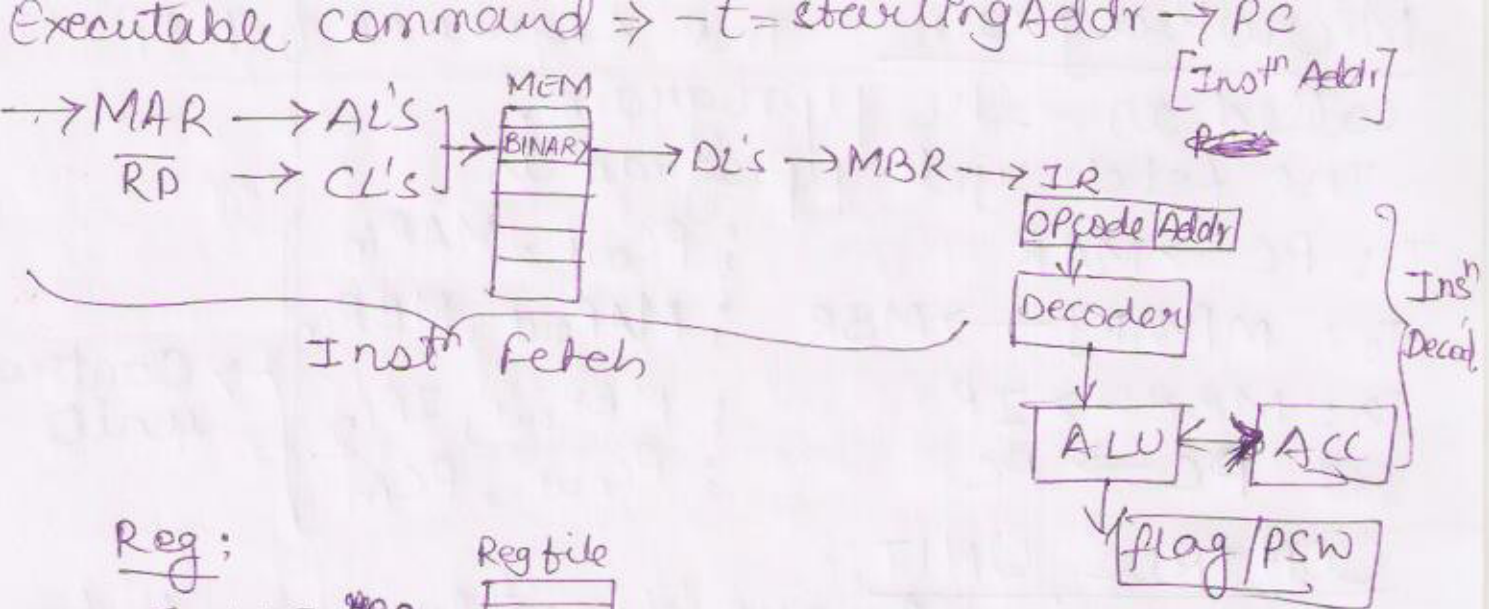
- 5) Performance increases.  
6) ~~Expensive~~ Cost/Bit increases.
- NOTE: During the execution of instruction, data becomes the most imp. So use the smallest storage <sup>(registers)</sup> component to hold the frequently accessed data. ∴ The processor maintains minimum set of registers.

### Mandatory Register:

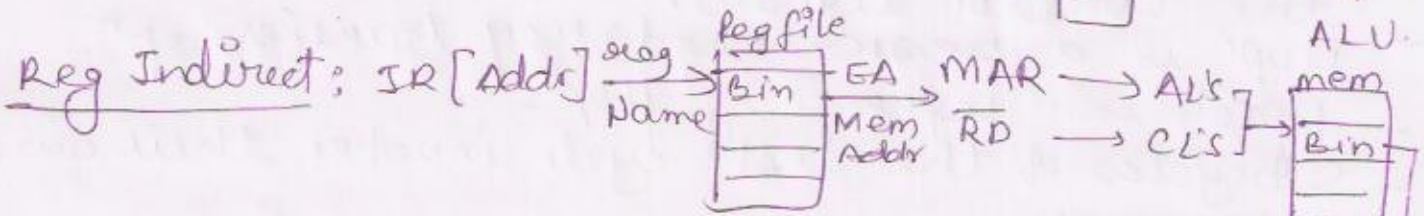
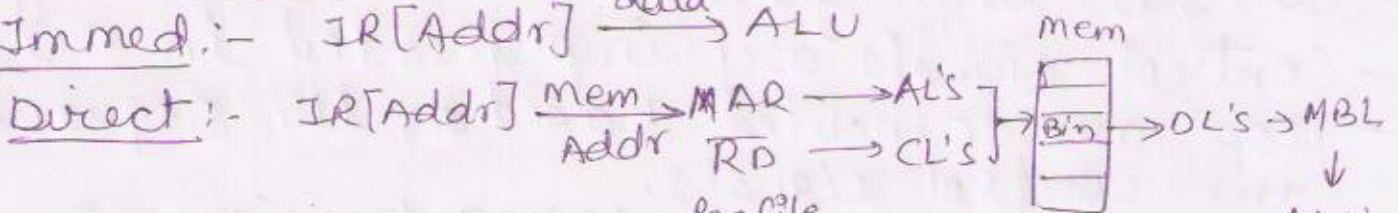
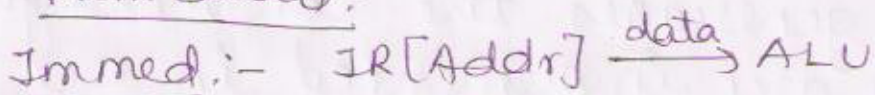
1. PC: It is used to hold the starting address & immediately points to the next inst<sup>n</sup> addr.
2. IR: It is used to hold the currently fetched inst<sup>n</sup> to be decoded because inst<sup>n</sup> format is predefined in this register, so it is not user accessible register.
3. ACC: It is user accessible register, used to hold the ALU input & output.
4. TR: It is used to hold the 2<sup>nd</sup> ALU operand. It is not a user accessible register.
5. MAR: It is used to carry memory address. MAR register is directly connected to address lines of system bus.
6. MBR/MDR: It is used to carry the binary sequence. MDR register is directly connected to data lines of system bus.

NOTE: - Apart from the above registers, processor supports set of the general purpose registers.

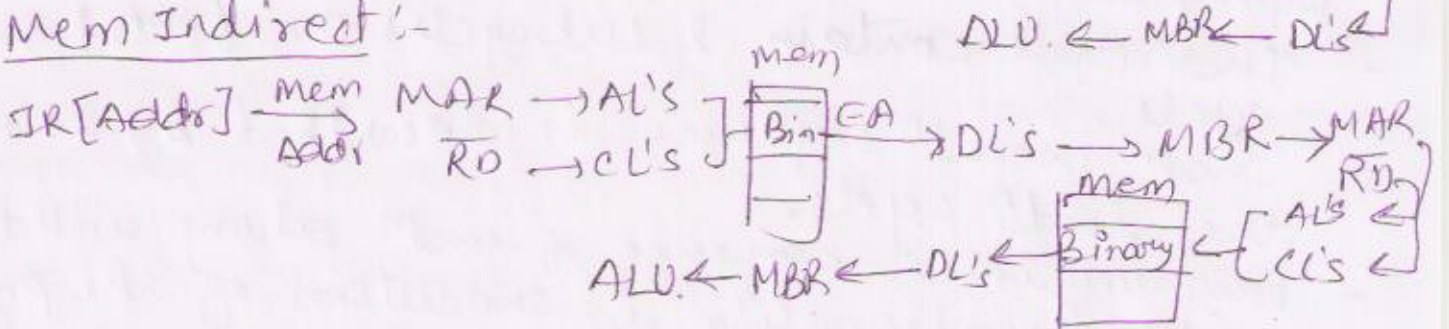
### INSTRUCTION EXECUTION SEQUENCE USING MICROOPERATIONS:-



Mem Based:



Mem Indirect :-



Microoperation is a basic reg. to reg. transfer op<sup>n</sup>. Each  $\mu\text{inst}^n$  can be executed in 1 cycle.  $\mu\text{op}^n$  is also called as  $\mu\text{inst}^n$  or control word or Atomic op<sup>n</sup>.

Control signals are req<sup>d</sup> to execute the  $\mu\text{inst}^n$ .

Eg. (a)  $T_1: R_1 \rightarrow R_2$ ;  $R_{1\text{out}}, R_{2\text{in}}$ .

(b)  $R_{1\text{out}}, R_{2\text{in}}, R_{3\text{in}} \Rightarrow \begin{matrix} R_1 \rightarrow R_2 \\ R_1 \rightarrow R_3 \end{matrix} \left. \vphantom{\begin{matrix} R_1 \rightarrow R_2 \\ R_1 \rightarrow R_3 \end{matrix}} \right\} \text{in same cycle.}$

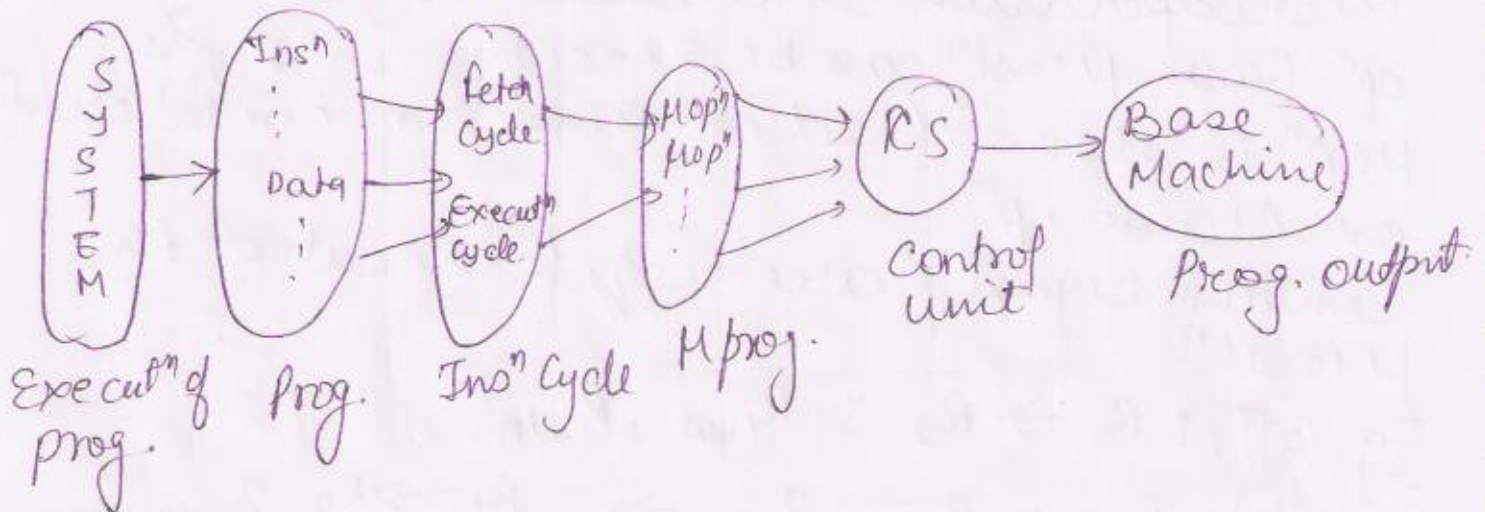
MicroProgram: The sequence of  $\mu op^n$  is called as the  $\mu$ program.

- The fetch cycle  $\mu$ program is

$T_1: PC \rightarrow MAR$	$; PC_{out}, MAR_{in}$	} Control unit
$T_2: M[MAR] \rightarrow MBR$	$; MAR_{out}, MBR_{in}$	
$T_3: MBR \rightarrow IR$	$; MBR_{out}, IR_{in}$	
$PC \rightarrow PC$	$; PC_{out}, PC_{in}$	

CONTROL UNIT:-

- Control signals are directly executed on base machine.
- Control unit generates the control signals.
- Control signals are implemented in CU.
- $\mu op^n$  are executed based on the sequence of the control signals.
- $\mu op^n$  is a basic reg to reg transfer  $op^n$ .
- $\mu$ prog is a sequence of  $\mu op^n$ s.
- Subcycles of the  $inst^n$  cycle invokes their own  $\mu$ programs.
- $Inst^n$  cycle contain 2 subcycles i.e fetch & execute cycle.
- $Inst^n$  execute $^n$  sequence is described by using the  $inst^n$  cycle.
- Program is a sequence of  $inst^n$  along with data.
- System functionality is execution of the program.



# Control unit

To design the control unit the designer must identify following 3 parameters:

- The no. of op<sup>n</sup> (Inst<sup>n</sup>) possible in the processor.
- The no. of  $\mu\text{op}^n$  required to execute each inst<sup>n</sup>.
- The no. of Control signals required to execute each  $\mu\text{op}^n$ .

After identifying the above 3 parameters designer maintains the database.

Eg: #inst<sup>n</sup> { I<sub>1</sub>, I<sub>2</sub>, I<sub>3</sub> } op code = 2 bits  
 { ADD SOB MUL }

- 00 - ADD
- 01 - SOB
- 10 - MUL
- 11 - undef.

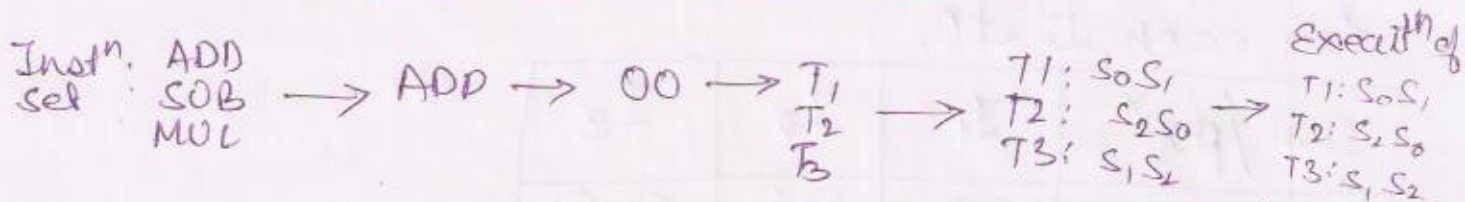
# of control signals recognized by b/w = { S<sub>0</sub>, S<sub>1</sub>, S<sub>2</sub> }

#  $\mu\text{op}^n$  req to execute each inst<sup>n</sup>: { T<sub>1</sub>, T<sub>2</sub>, T<sub>3</sub> }

Database :

Inst <sup>n</sup> / $\mu\text{op}^n$ :	00 I <sub>1</sub> ADD	01 I <sub>2</sub> SOB	10 I <sub>3</sub> MUL
T <sub>1</sub>	S <sub>0</sub> , S <sub>1</sub>	S <sub>0</sub> , S <sub>2</sub>	S <sub>2</sub> , S <sub>1</sub>
T <sub>2</sub>	S <sub>2</sub> , S <sub>0</sub>	S <sub>2</sub> , S <sub>1</sub>	S <sub>0</sub> , S <sub>2</sub>
T <sub>3</sub>	S <sub>1</sub> , S <sub>2</sub>	S <sub>1</sub> , S <sub>0</sub>	S <sub>1</sub> , S <sub>0</sub>

implemented in  
 $\Rightarrow$  CU



Match to previous circle diagram.

- After defining database, designer uses any one of the following approaches to imp. the database

- Hardwired Approach
- microprogrammed Approach

## Hardwired control unit design :-

- In this design, control signals are expressed as sum of product expressions; they are directly implemented on the independent h/w.
- It is the fastest control unit design when control signals are implemented on independent hardware.
- It is used in real time app<sup>n</sup>. Eg:- aircraft simulation & weather forecasting, etc.
- Even a minor modification requires redesign and reconnection of the control signals. ∴ It is not flexible.
- It is not suitable in the designing and testing places.
- Risc control unit is hardwired CU.

Q. Consider a hypothetical control unit. It supports 3 inst<sup>n</sup> ( $I_1, I_2, I_3$ ) and 4 control signals ( $S_0, S_1, S_2, S_3$ ). Each inst<sup>n</sup> requires 4  $\mu$ op<sup>n</sup>s ( $T_1, T_2, T_3, T_4$ ) to complete the execution. The following table indicates what are the control signals required for each  $\mu$ op<sup>n</sup>, for each inst<sup>n</sup>.

Inst <sup>n</sup> / $\mu$ op <sup>n</sup>	$I_1$	$I_2$	$I_3$
$T_1$	$S_0 S_1$	$S_2 S_0$	$S_3 S_2$
$T_2$	$S_1 S_2$	$S_3 S_2$	$S_0 S_2$
$T_3$	$S_3 S_2$	$S_1 S_0$	$S_3 S_0$
$T_4$	$S_0 S_3$	$S_2 S_3$	$S_3 S_0$

Get the control exp. for  $S_0, S_1, S_2, \& S_3$ .



$$S_0 = T_0(I_1 + I_2) + T_2 I_3 + T_3(I_2 + I_3) + T_4(I_1 + I_3)$$

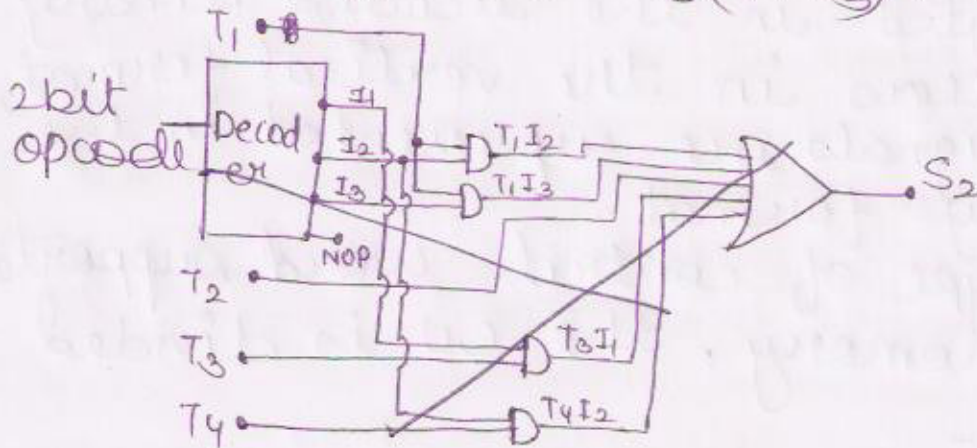
$$S_1 = T_1 I_1 + T_2 I_1 + T_3 I_2$$

$$S_2 = T_1(I_2 + I_3) + T_2(I_1 + I_2 + I_3) + T_3 I_1 + T_4 I_2$$

$$S_3 = T_1 I_3 + T_2 I_2 + T_3(I_1 + I_3) + T_4(I_1 + I_2 + I_3)$$

$$T_1(I_2 + I_3) + T_2 + T_3 I_1 + T_4 I_2$$

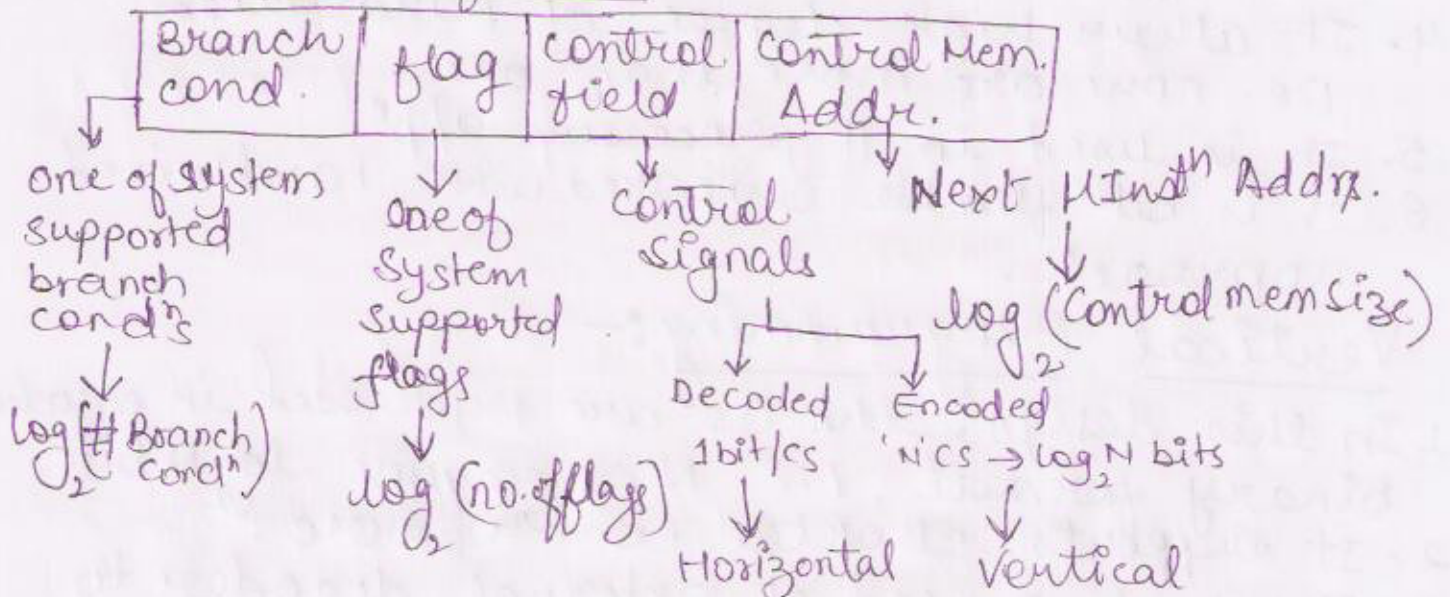
$$T_1 I_3 + T_2 I_2 + T_3(I_1 + I_3) + T_4$$



### MicroProgram Control Unit Design:-

- In this design, control memory is used to store the  $\mu$ prog. The control memory is permanent memory i.e. ROM.
- $\mu$ prog is a sequence of  $\mu$ operations.
- $\mu$ inst<sup>n</sup> are stored in the control memory by using the following format.

$\mu$ inst<sup>n</sup> format:-



If control format is horizontal  $\rightarrow$  Horizontal programming.

Based on way of representing control signals,  $\mu$ inst<sup>n</sup> is divided into 2 types i.e. (i) Horizontal  $\mu$ inst<sup>n</sup> (ii) Vertical  $\mu$ inst<sup>n</sup>.

In Horizontal format, the control signals are represented in the decoded binary format, whereas in the vertical format, the control signals are represented in the encoded binary format.

-Based on type of control word supported by control memory, the CU is divided into 2 types:-

- (i) Horizontal program CU
- (ii) Vertical program CU.

### Horizontal Programming:-

1. In this design, the CS are expressed in decoded binary format i.e. 1 bit/CS.
2. It supports longer control word.
3. No need of the external decoder to generate the CS, so it is faster than vertical.
4. It allows high degree of parallelism. i.e. none ~~are~~ more than one.
5. It is used in II<sup>l</sup> processing app<sup>n</sup>.
6. It is bit flexible compared with hardwired approach.

### Vertical Programming:-

1. In this design, the CS are expressed in encoded binary format, i.e. N CS require  $\log_2 N$  bits.
2. It supports shorter control word.
3. There is a need of external decoder to generate the control signals. So it is slower.

than horizontal.

4. It allows low degree of parallelism. i.e. none or one.

5. It allows easy implementation of new inst<sup>n</sup>.  
↳ it is more flexible.

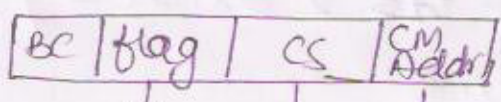
NOTE:-

The ascending order of the CU design in terms of speed is Vertical - Horizontal - Hardwired  
low → → → high

The ascending order of CU design in terms of flexibility  
low → high  
Hardwire - Horizontal - Vertical.

Q. Consider a hypothetical CU. It supports 1024 control word memory. It has 48 CS, and 16 flags. What is the size of Cword in bits & control memory in bytes using the horizontal programming?

Ans:-



↓                      ↓                      ↓  
 4                      48                      10 = 62 bits

CW size = 62 bits  
 = 1024 CW  
 = 1024 × 62 bits  
 =  $\frac{(1024 \times 62)}{8}$  bytes

Q. A hypothetical CU supports 8 mutually exclusive groups of CS that is tabulated below:-

Groups	G <sub>1</sub>	G <sub>2</sub>	G <sub>3</sub>	G <sub>4</sub>	G <sub>5</sub>	G <sub>6</sub>	G <sub>7</sub>	G <sub>8</sub>
CS	1	3	7	6	12	14	20	2

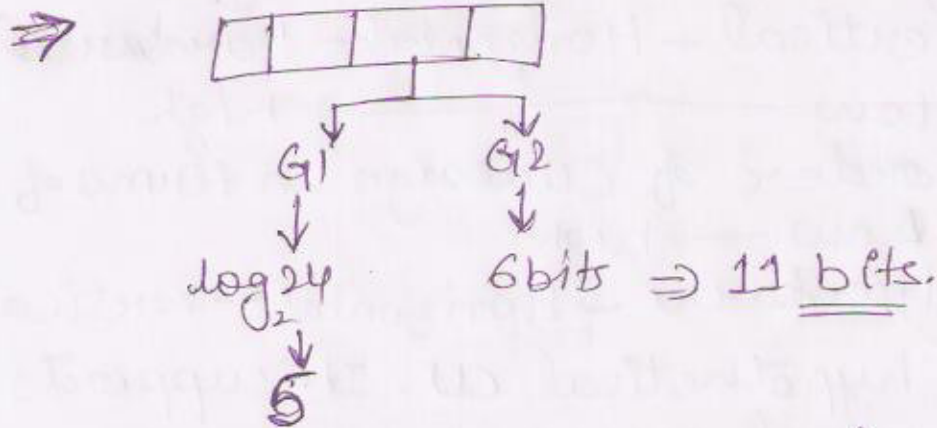
How many #bits saved using the vertical programmed over horizontal programming?

Ans:-    H = 1 + 3 + 7 + 6 + 12 + 14 + 20 + 2 = 65 bits  
           V =  $\log_2 1 + \log_2 3 + \log_2 7 + \log_2 6 + \log_2 12 + \log_2 14 + \log_2 20 + \log_2 2$

= 23 bits

∴ No. of bits saved =  $(65 - 23) = 42$  bits

Q. A C field ~~has~~ of  $\mu$ inst<sup>n</sup> supports 2 groups of CS. Group 1 indicates ~~number~~ of 24 CS and group 2 indicates at most 6 from remaining. What is the size of control field?



## Performance Evaluation of Processing

Performance is an indirect measured parameter

It is depending on ~~the~~ 2 directly measured parameters

a) Execution time

(b) Throughput

- The time required to complete the program is called execution time or elapsed time or response time.

- Consider 2 systems X and Y. Both are used to execute the same task in which X system takes 5 ns + Y system takes 10 ns time to complete the task.

$$P \propto \frac{1}{E.T.}$$

- The no. of tasks completed in a unit of time is called as throughput.

$$P \propto T$$

- The performance gain can be measured by using the speedup factor. ~~Speedup~~ speedup factor compares the 2 different entities.

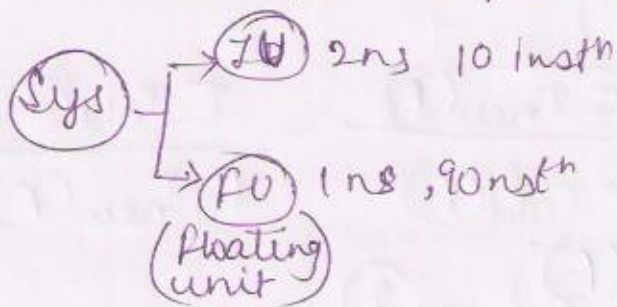
$$\text{i.e. speedup} = \frac{\text{Performance}_x}{\text{Performance}_y}$$

(5)

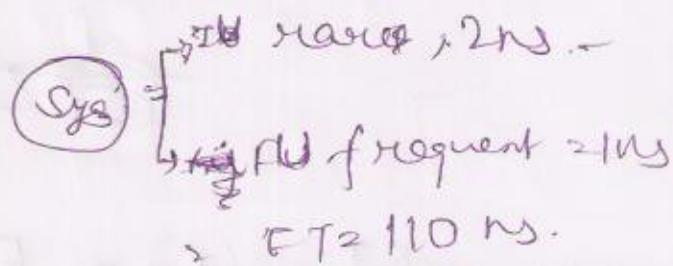
$$S = \frac{ET_y}{ET_x} ; S = 'n'$$

System 'X' runs 'n' times faster than system 'Y'

- Quantitative principles are used to improve the performance.
- These principles states that improve the performance of processor within the cost & price limits while making the common case fast.
- This principle is called as "Amdahl's Law"



$$ET = 2 \times 10 + 1 \times 90 = 100 \text{ ns.}$$



~~$$ET = 2 \times 10 + 1 \times 90 = 100 \text{ ns.}$$~~

Amdahl's law is focussed on performance gain after enhance the frequent components.

i.e.

$$\text{Overall} = \frac{\text{Performance of the Sys with enhancement}}{\text{Performance of the Sys with out enhancement.}}$$

$$\text{Overall} = \frac{1/ET_{\text{new}}}{1/ET_{\text{old}}} = \frac{ET_{\text{old}}}{ET_{\text{new}}} \quad (1)$$

- After enhance the system, ~~it~~ it contains 2 portions i.e. enhanced portion and unenhanced portion.

$$\Rightarrow ET_{\text{new}} = ET_{\text{old}} + ET_{\text{enhanced portion}} \quad (2)$$

2 parameters are needed to calculate the  $ET_{\text{new}}$

- (i) Fraction enhanced (F).
- (ii) Speedup enhanced (S).

$f$  fraction enhanced indicates how much portion of the system is enhanced.

$f =$  enhanced  
 $(1-f) =$  unenhanced

$\Rightarrow$  Speedup enhanced indicates, how many times the enhanced portion is running faster than the old portion.

$S = \frac{\text{Performance of enhanced portion}}{\text{Performance of old portion.}}$

$$= \frac{1}{\cancel{ET_{new}(f)}} \cdot \frac{ET_{old}(f)}{ET_{new}(f)}$$

②<sup>nd</sup> in ③<sup>rd</sup>  $ET_{new}(f) = \frac{ET_{old}(f)}{S} \quad \text{--- (2)}$

$$ET_{new} = \frac{ET_{old}(1-f) + ET_{old}(f)}{S}$$

substitute ② in ①,

$$S_{overall} = \frac{ET_{old}}{\cancel{ET_{old}(f)} + ET_{old}(1-f) + \frac{ET_{old}(f)}{S}}$$

Let  $ET_{old} = 1$ .

$$S_{overall} = \frac{1}{(1-f) + \frac{f}{S}}$$

If  $ET_{old} \geq 1$ , then  $ET_{new}$  is always less than 1 &  $S_{overall}$  is always greater than 1.

NOTE:- When the system contains multiple enhancements then  $S_{overall} = \frac{1}{(1 - \sum f_i) + \sum \frac{f_i}{S_i}}$   
 $S_i = \#$  enhancements.

Q. Consider a computer which is used in the mathematical model simulations. It consists of 2 components: - Integer & floating point units. FPU is enhanced then it runs 2 times faster but only 10% of the instr are floating point. What is the performance gain.

Ans:- 
$$S_{\text{overall}} = \left( (1-f) + \frac{f}{S} \right)^{-1} = \left[ (1-0.1) + \frac{0.1}{2} \right]^{-1}$$

$$= 1.05$$

$$\approx 1$$

According to the problem statement, instead of frequent case the rare case goes on enhancement.  $\therefore$  there is no considerable performance gain.

The corrective measurement is enhance the frequent case i.e. Integer unit.

$$S_{\text{overall}} = 1.82 \approx 2 \quad \therefore \text{considerable gain}$$

Q. Suppose the cache memory is 10 times faster than main memory & it is used 70% of the time. How much speedup do we gain by using this cache.

$$\left[ (1-0.7) + \frac{0.7}{10} \right]^{-1} = \left( 0.3 + \frac{0.7}{10} \right)^{-1}$$

$$= \left( \frac{3.7}{10} \right)^{-1} = \frac{10}{3.7} = 2.7$$

$\Rightarrow$  Introducing cache improves performance by 2.7 times.

Q. 3 enhancements are proposed for a new architecture with the following speedups  $\Rightarrow S_1=30, S_2=20, S_3=15$ .

If the enhancement 1 & 2 both will

be used 25% of the time, what fraction of the time enhancement 3 will be used to gain the overall of 10

$$10 = \left( \frac{1}{0.5} \right)$$

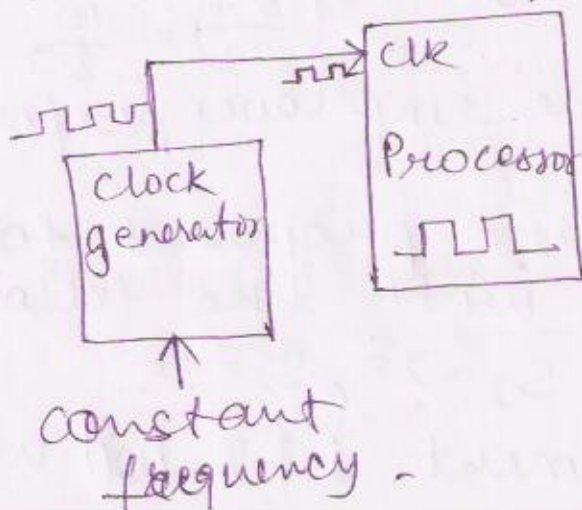
$$10 = \left[ \left[ 1 - (0.5 + x) \right] + \frac{0.25}{30} + \frac{0.25}{20} + \frac{x}{15} \right]^{-1}$$

$$10 = \left[ 0.5 + x + \frac{0.5 + 0.75}{60} + \frac{x}{15} \right]^{-1}$$

$$10 = \left( 0.5 + x + \frac{0.25}{60} + \frac{x}{15} \right)^{-1}$$

### Calculation of CPU time:-

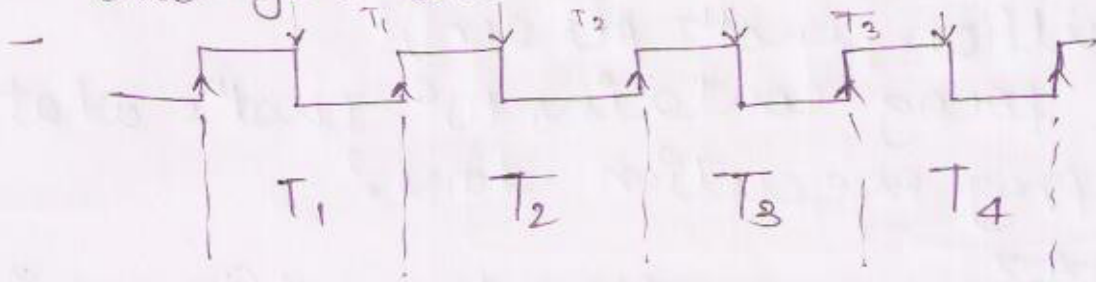
- CPU time means program execution time.
- Processor is having clock pin as an input pin.
- The clock pin is externally connected with the clock generator.
- Clock generator is operating with the constant frequency to generate the clock pulses.
- These clock pulses are carried into the CPU through the CLK pins. The CPU operations are controlled by the clock.



The program execution time can be calculated based on cycle & cycle time.



- Cycle is defined as raising edge to raising edge transition or falling edge to falling edge transition.
- It is also called as the 'tick' / ck tick / cycle / clock cycle etc.



- The time required to transfer the pulse either from raising edge to raising edge or from falling edge to falling edge is called as cycle time.
- The cycle time is depending on the clock frequency.

$$1 \text{ GHz} = \text{frequency}$$

$$\Rightarrow t = \frac{1}{1 \text{ GHz}} = 10^{-9} \text{ sec.} = 1 \text{ ns.}$$

The program execution time = # sec per

$$= \frac{\# \text{ inst}^n / \text{prog} \times \# \text{ cycles} / \text{inst}^n}{\# \text{ seconds} / \text{cycle}} \Rightarrow \frac{\# \text{ inst}^n / \text{prog}}{\# \text{ cycles} / \text{inst}}$$

$$ET = \text{Inst}^n \text{ Count} * \underset{\text{cycles/inst}}{\text{CPI}} * \text{cycletime}$$

Program is a mixing of ~~3~~ 3 categories of inst<sup>n</sup>. So different inst<sup>n</sup> take different # cycles.  $\therefore$  The Prog. execut<sup>n</sup> time / cost time

$$= \sum_{i = \text{type of inst}^n} (I_i * CPI_i) * \text{Cycletime}$$

Q. Consider 1ns clock cycle processor. It consumes 4 cycles for load + store op<sup>n</sup>, 6 cycles for ALU op<sup>n</sup> and 5 cycles for branch op<sup>n</sup>.

The relative frequencies of these inst<sup>n</sup> are 40%, 40% & 20% resp.

(a) What is the average inst<sup>n</sup> execution time?

(b) What is the performance in terms of MIPS. (Million inst<sup>n</sup> per sec)?

(c) If the prog contains  $10^6$  Inst<sup>n</sup>s what is the prog execution time?

Ans:-

$$\begin{aligned} \text{(a) Avg inst}^n \text{ execution time} &= \sum (I_i \times CPI_i) \times \text{cycle} \\ &= [(0.4 \times 4) + (0.4 \times 6) + (0.2 \times 5)] \times 1 \text{ ns} \\ &= 5 \text{ ns.} \end{aligned}$$

$$\begin{aligned} \text{(b) } 1 \text{ inst}^n &\rightarrow 5 \text{ ns.} \\ \# ? &\rightarrow 1 \text{ sec.} \end{aligned}$$

$$\frac{1}{5 \times 10^{-9}} = 2 \times 10^8 \text{ Inst/sec} = 200 \text{ MIPS.}$$

$$\begin{aligned} \text{(c) } 10^6 &\rightarrow 5 \times 10^{-9} \times 10^6 \\ &= 5 \times 10^{-3} \\ &= \underline{\underline{5 \text{ ms.}}} \end{aligned}$$

## (5D) HIGH PERFORMANCE ARCHITECTURE

High performance architecture exhibits the concurrency.

- Concurrency means more than one event execution.
- Event may be program, subprogram, inst<sup>n</sup> or stage (Intrinsic)
- Concurrency implies parallelism, simultaneous, pipelining.
- Parallelism means 2 or more event execut<sup>n</sup>

at the same time period.

→ Simultaneous means 2/more event execution at the same time instance.

→ Pipelining means 2 or more event executed in overlapping span.

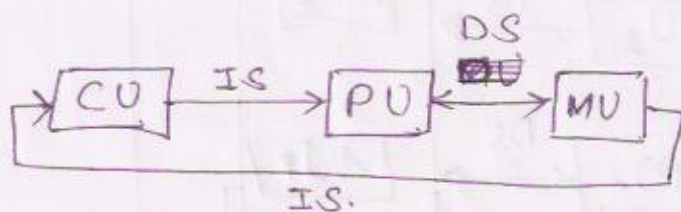
According to the Flynn's archi classification, the parallel architecture is divided into 4 catg:

- (i) SISD → Single inst<sup>r</sup> single data~~ing~~ stream
- (ii) SIMD → Single " multiple "
- (iii) MISD → Multiple " single "
- (iv) MIMD → ~~Multi~~ " " multiple "

PU: Processing Unit  
CU: Control Unit  
MU: Memory Unit

IS: Inst<sup>r</sup> stream  
DS: Data stream

SISD:

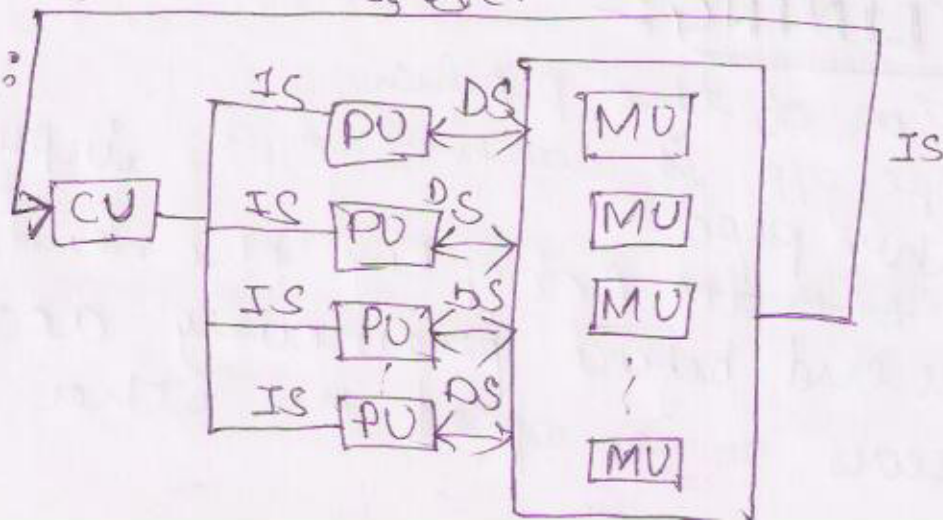


Inst<sup>n</sup> & data both are stored in same memory. After performing the operation, the result is also stored in same memory. This concept is known as stored program structure.

→ Von Neumann Archi. was designed based on stored program concept.

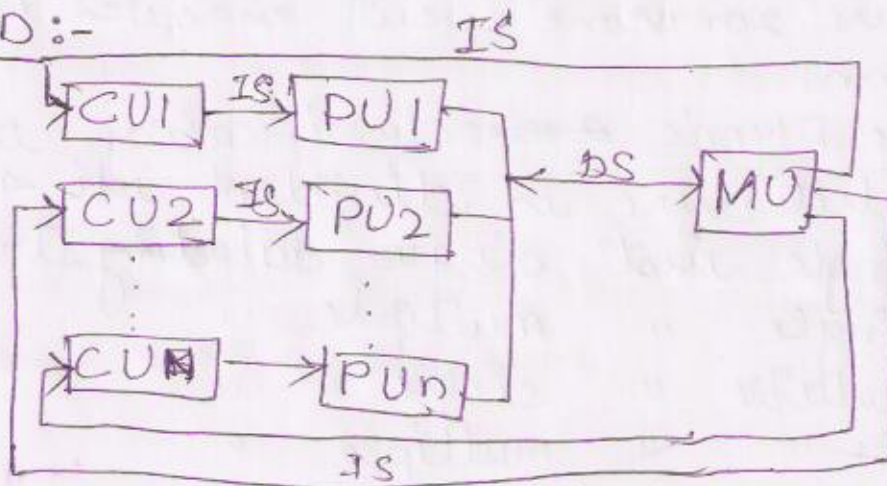
→ Uniprocessor is using the SISD architecture.  
Eg → VAX, IBM 360, etc.

SIMD:



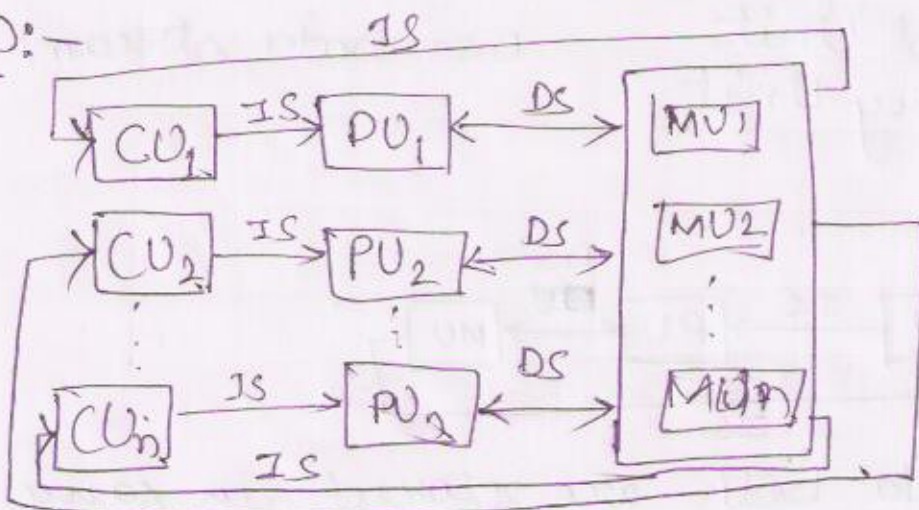
→ Array processors use the SIMD arch.  
 eg → STERAN, PEPE, MPP etc.

MISD:-



Not yet implemented.

MIMD:-



Multiprogramming <sup>sys</sup>  
 Multiprocessor system uses MIMD architect  
 use. By Cray, Cyber etc.

→ To improve the performance of uniprocessor system pipelining is used.

### PIPELINING:-

Function of the pipeline

→ 1 pipe out is connected as input to the another pipe.

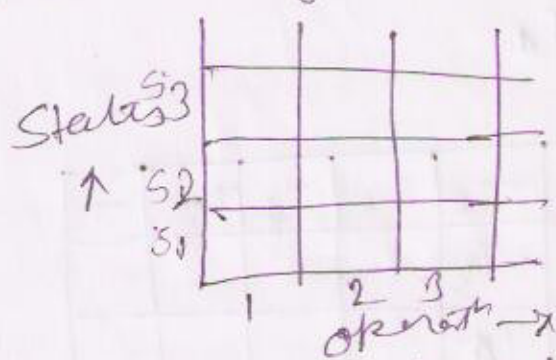
Def<sup>n</sup>:- ~~Accessing~~ Accepting the new inputs at one end before previously accepted input appears as an output at the other end.

Def<sup>n</sup> states that insert the new IP before completion of old inputs. That means new IPs are executed along with the old inputs. So pipeline allows the overlapping execution.

Non pipeline allows the non-overlapping execution. i.e. new inputs are inserted only after completion of old inputs.

The overlapping execution sequence can be represented using space time diagram.

- In this diagram  $x \rightarrow$  cycles,  $y \rightarrow$  stages.



The successful char. of pipelines is that in every new cycle, new input must be inserted into pipeline.

$\therefore \text{CPI} = 1.$

Designing of the pipeline - IP end  
 to Pipe has 2 ends.  $\left[ \begin{array}{l} \text{IP end} \\ \text{OP end} \end{array} \right.$   
 Between the inp & outputs ends, the input & output pipe, is interconnected.

→ Interface registers are used b/w the stages - to hold the intermediate results. They are also known as buffer / Latch, Pipelining register.

→ All the stages ~~are~~ along with interface registers are controlled by the common clock.



It takes  $k$  cycles to complete the operation.  
 The remaining  $n-1$  tasks are emerged from the pipe at the rate of 1 cycle per task.  $\therefore$   $n-1$  tasks consumes  $n-1$  cycles to complete the operation.  
 $\rightarrow$  So the total time required to process the  $n$  tasks using the  $k$  segment pipeline is

$$ET_{\text{pipe}} = (k+n+1) \text{ cycles.}$$

$$= (k+n+1) t_p.$$

Consider non pipeline processor used to execute the  $n$ -tasks in which each task takes  $t_n$  time to complete the operation.

$\therefore$  The total time required to complete the  $n$  tasks in non pipeline is

$$ET_{\text{non-pipe}} = n t_n$$

The performance gain of the pipeline processor over the non pipeline is

$$\text{Speedup (S)} = \frac{\text{Performance}_{\text{pipe}}}{\text{Performance}_{\text{non-pipe}}}$$

$$S = \frac{1/ET_{\text{pipe}}}{1/ET_{\text{non-pipe}}} = \frac{ET_{\text{non-pipe}}}{ET_{\text{pipe}}}$$

$$= \frac{n t_n}{(k+n+1) t_p}$$

When the no. of tasks increases, then  $n$  becomes too large than  $k-n$ . Then  $k+n+1$  value is approached to  $n$  when  $n \gg k$ . then

$$k+n+1 \approx n.$$

Now

$$S = \frac{n t_p}{(k+n+1) t_p} \approx \frac{n t_p}{n t_p}$$

$$S = \frac{t_n}{t_p}$$

When all the tasks take same no. of cycles then 1 task execution time is also equal to the no. of stages in the pipeline.

$$t_n = k \text{ cycles.}$$

$$t_n = k \cdot t_p$$

Under this cond<sup>n</sup>,  $S = \frac{k \cdot t_p}{t_p}$

$$\Rightarrow \boxed{S = k}$$

$k \rightarrow$  no. of stages in pipeline

Also called as pipeline depth.

When the processor is operating with 100% efficiency then maximum speedup depth the pipeline depth.

$$100\% \eta \rightarrow S_{\max}$$

$$?, \eta \rightarrow \frac{S}{k}$$

$$\boxed{\eta_{\text{pipe}} = \frac{S}{S_{\max}} = \frac{S}{k}}$$

The throughput of the pipeline is # task processed per total time taken to process the task.

$$T_{\text{pipe}} = \frac{\# \text{ tasks processed}}{\text{total time taken to process the tasks.}}$$

$$= \frac{n}{(k+n-1) t_p}$$

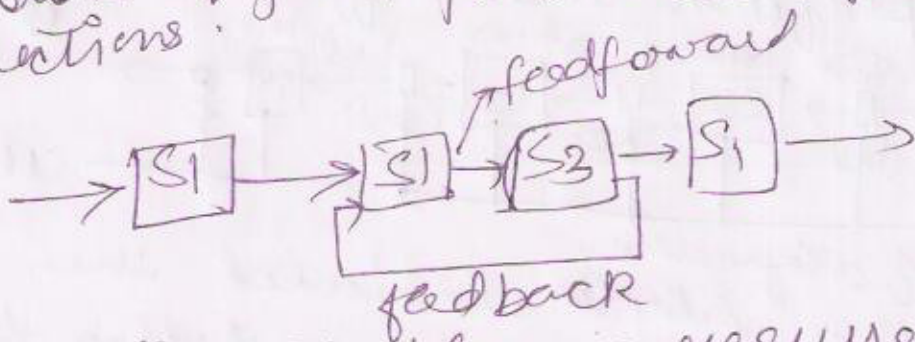
Types of the pipeline:-

(1) Linear pipeline:- This pipeline performs only 1 specific functionality. It consists only the feed forward connect<sup>n</sup>.





Non-linear pipeline: This pipeline can perform multiple functionality. It consists of feed forward & feed backward connections. (Asynchronous)



Reservation table is required to use the pipeline stages.

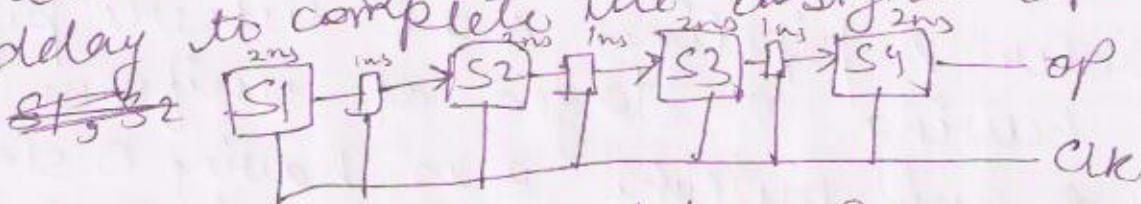
- Synchronous pipeline:-

In this pipeline operations are initiated based on the clock.

- Asynchronous pipeline:-

→ In this pipeline, operations are initiated based on the hand shaking signals.

Uniform delay pipeline:- In this pipeline all the stages maintains the same amt of delay to complete the assigned operation.



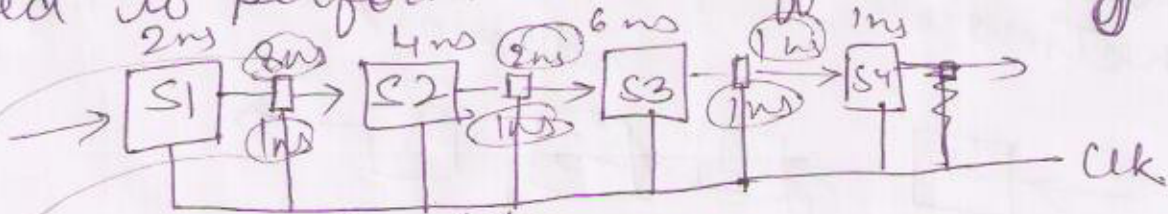
(i)  $t_p = \text{stage delay} = 2\text{ns}$ .

(ii) If buffer delay is included, then  
 $t_p = \text{Stage delay} + \text{buffer delay}$   
 $= 3\text{ns}$

(iii) Skew or setup time overhead is included, then  
 $t_p = t_p + \text{overhead time}$ .

## Non-uniform delay:-

In this pipeline, different stage delays are used to perform the different stage operations.



1)  $t_p = \max(\text{stage delay})$

2)  $t_p$  - if buffer delay is included then  
 $t_p = \max(\text{stage delay}) + \text{Buffer delay}$

3) In nonuniform buffer delay is included,  
 $t_p = \max(\text{stage delay} + \text{buffer delay})$

Q. Consider an  $n$  stage pipeline which has the speedup factor ~~10~~ while operating with 80% efficiency. What would be # stages present in the pipeline?

$$\Rightarrow \eta = 80\% \quad \eta = \frac{S}{K} \quad K = \frac{10}{0.8} = 12.5 = \underline{13}$$

$S = 10$

Q. Consider 2 pipelines A & B where pipeline A is having 8 stages of uniform delay of 2ns. pipeline B is having 5 stages with the respective stage delays of 2ns, 3ns, 4ns & 1ns, 2ns. How much time is saved when 100 tasks are pipelined using A instead of B.

$$\Rightarrow \text{Pipe A: } K = 8$$

$$\eta = 100$$

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

$$ETA = (K + n - 1) t_p$$

$$= (8 + 100 - 1) 2 = 214 \text{ ns}$$

$$\text{Pipe b: } (5+100-1) \cdot t_p = 104 \times 4 = 416 \text{ ns.}$$

$$\text{Time saved} = 416 - 214 = \underline{\underline{202 \text{ ns.}}}$$

Q Consider 4 segment pipeline with respective stage delays of 10, 20, 5, 15 ns. What is the approx speedup when very large no. of instr are pipelined?

$$\Rightarrow S = \frac{t_n}{t_p}$$

$$t_p = \max(10, 20, 5, 15)$$

$$t_p = 20 \text{ ns.}$$

$$\text{(Nonoverlapping execution time)} \quad t_n = S_1 + S_2 + S_3 + S_4 = 10 + 20 + 5 + 15 = 50 \text{ ns}$$

$$S = \frac{50}{20} = 2.5$$

$$\eta = \frac{S}{K} = \frac{2.5}{4} = 0.625$$

Q: Consider 4 stage pipeline with the following time delays:

$$t_1 = 20 \text{ ns}$$

$$t_2 = 40 \text{ ns}$$

$$t_3 = 15 \text{ ns}$$

$$t_4 = 25 \text{ ns.}$$

The interface regis. used b/w the stages have the delay of 5 ns, what is the speedup?

$$\Rightarrow S = \frac{t_n}{t_p} = \frac{100}{45} = \underline{\underline{2.2}}$$

No need of adding buffered delay because there is no need to save intermediate results in nonoverlapping execution.

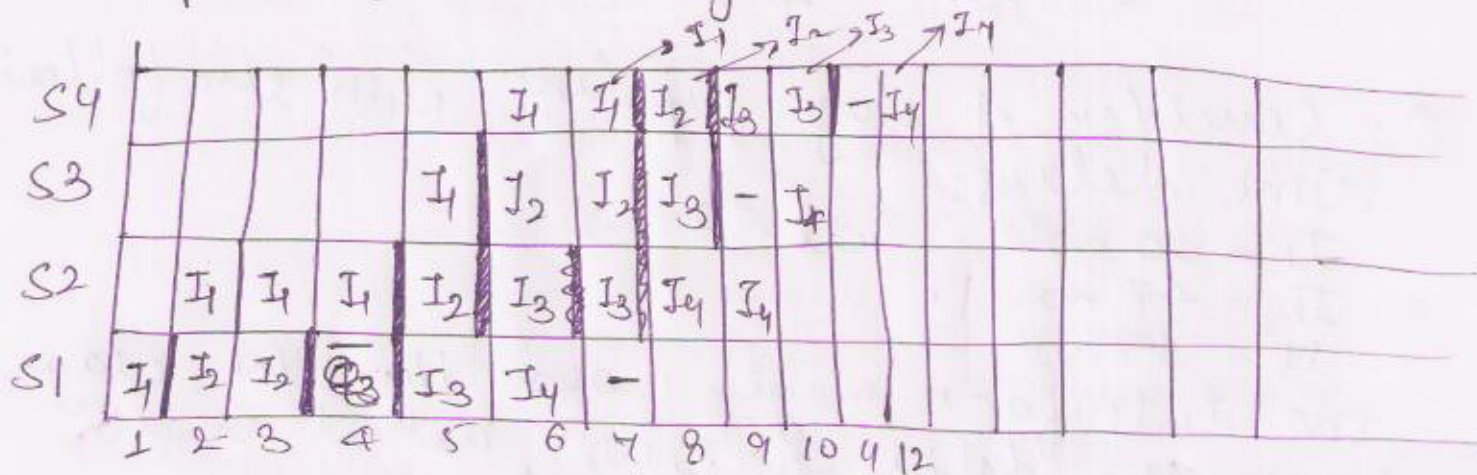
Q. Consider a 4 stage instruction pipeline where different instructions are spending different amount of time at different stages that is shown below: -

	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	S <sub>4</sub>	
I <sub>1</sub>	1	3	1	2	⇒ 7
I <sub>2</sub>	2	1	2	1	⇒ 6
I <sub>3</sub>	1	2	1	2	= 6
I <sub>4</sub>	1	2	2	1	= 6
					<u>25</u>

(i) How many cycles required to complete the above instr<sup>n</sup>?

(ii) what is the speedup?

⇒ ~~Not~~ Not uniform delay, not non-uniform delay. Special case. Solved using space time diagram.



(i) ⇒ 12 cycles.

$$(ii) \Rightarrow S = \frac{ET_{\text{nonpipe}}}{ET_{\text{pipe}}} = \frac{I_1 + I_2 + I_3 + I_4}{12} = \frac{25}{12} \approx 2.08$$



pipeline. It is used to execute the inst<sup>n</sup>. RISC processor supports 3 categories of instructions

(i) Data transfer inst<sup>n</sup>. LOAD & STORE inst<sup>n</sup> are used to transfer the data between CPU & memory, in which to access the memory Indexed addressing mode is used.

The inst<sup>n</sup> format is ~~LOAD r0, r1~~, ~~STORE r1, r2~~

LOAD r<sub>0</sub>, r<sub>1</sub> ; r<sub>0</sub> ← M[r<sub>1</sub>]

STORE r<sub>1</sub>, r<sub>2</sub> ; M[r<sub>2</sub>] ← r<sub>1</sub>

(ii) ALU operations: ALU operations are performed only on registers.

Inst<sup>n</sup> format is

Add r<sub>0</sub> r<sub>1</sub> r<sub>2</sub> ; r<sub>0</sub> ← r<sub>1</sub> + r<sub>2</sub>

(iii) Branch op<sup>n</sup>:-

unconditional jump: JMP 2000 ; PC ← ~~Seq Addr~~ TA

conditional jump: JNZ r<sub>0</sub>, 2000

NZ ← CMP r<sub>0</sub>

T: PC ← ~~Seq Addr~~ TA  
F: PC ← Seq Addr.

To execute the above inst<sup>n</sup>, RISC processor supports 5 stage inst<sup>n</sup> pipeline.

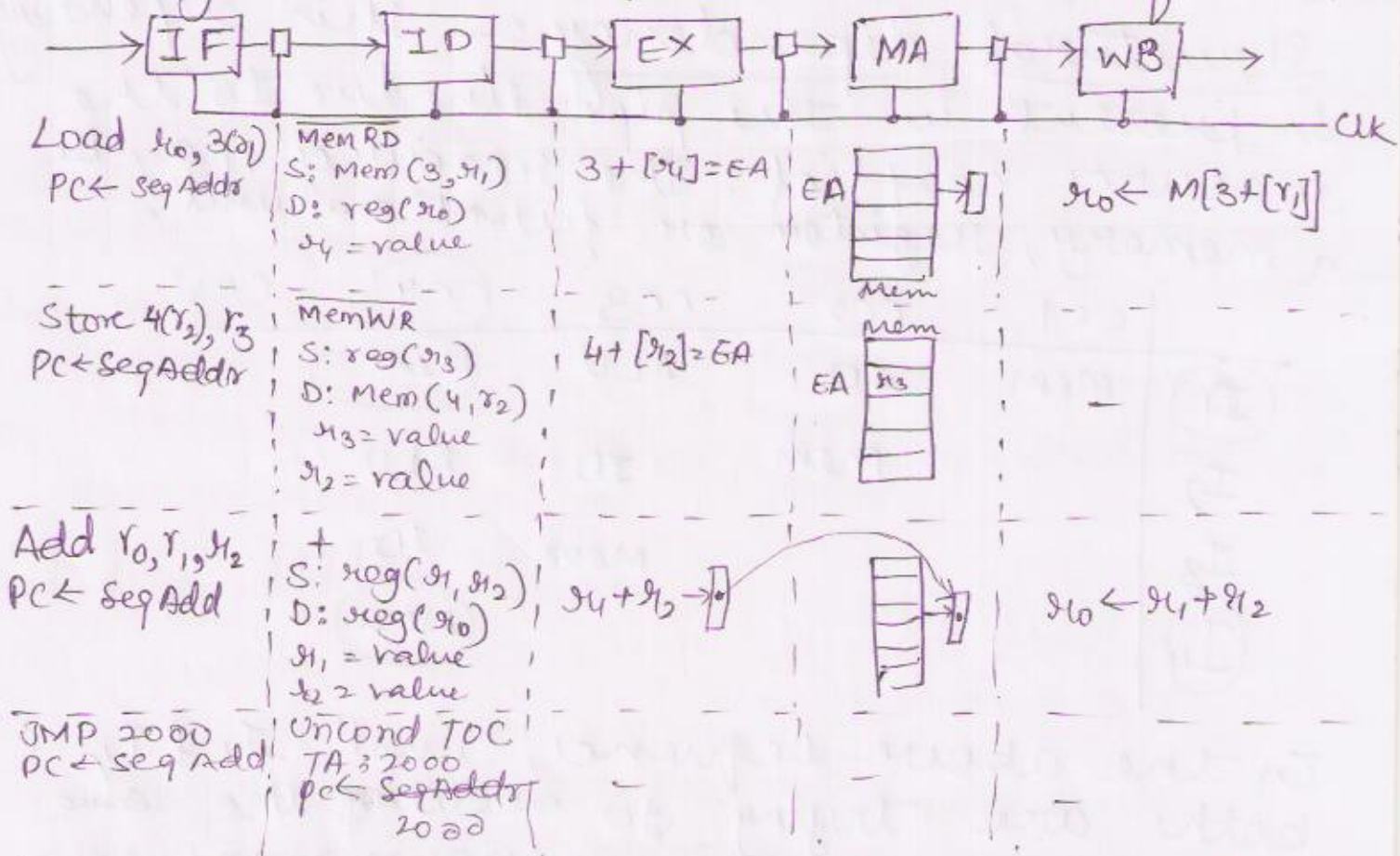
STAGE 1 (Inst<sup>n</sup> fetch):- In this stage inst<sup>n</sup> are fetched from memory based on the program counter. At the end of this stage, the program counter value is incremented to fixed constant.

STAGE 2 (INST DECODE) :- In this stage 2 operations are performed. i.e decode inst<sup>n</sup> & accessing the register file. This stage also contains the comparator logic to evaluate the branch conditions.

STAGE 3 (EXECUTION): In this stage ALU operations are performed.

STAGE 4 (Memory Access): In this, memory read & memory write operations are performed.

STAGE 5 (Write Back): This stage op<sup>n</sup> is divided into 2 halves. In the first half, register write operation is performed, in the 2<sup>nd</sup> half, register read operation is performed.



JNZ 20, 2000  
PC ← Reg Addr



Branch inst<sup>n</sup> → 2<sup>nd</sup> stage  
ALU + reg. load → 5<sup>th</sup> stage  
Reg ← reg store → 4<sup>th</sup> stage

Depend  
There are 3 types of dependencies possible in pipeline.

- (i) structural dependency.
- (ii) data dependency
- (iii) control dependency.

Dependency creates the extra cycles in the pipeline. Extra cycle without operation (New i/p initiation) called as stall.

Structural dependency: - This dependency is present in the pipe-line due to the resource conflict. The resource may be a memory, register or functional unit.

	CC1	CC2	CC3	CC4	CC5
I <sub>1</sub>	Mem	ID	ALU	Mem	
I <sub>2</sub>		Mem	ID	ALU	
I <sub>3</sub>			Mem	ID	
I <sub>4</sub>				Mem	

In the above sequence, inst<sup>n</sup> I<sub>1</sub> & I<sub>4</sub> both are trying to access the same resource memory. This is known as memory conflicts.



Conflict is an unsuccessful operation. So to make it successful, keep insts 4 into the wait until the resource becomes available. i.e., this waiting creates the stalls in the pipeline.

	CC1	CC2	CC3	CC4	CC5	CC6	CC7	CC8
I <sub>1</sub>	mem	ID	ALU	<del>mem</del> mem	WB			
I <sub>2</sub>		mem	ID	<del>mem</del> ALU	mem	WB		
I <sub>3</sub>			mem	ID	ALU	mem	WB	
I <sub>4</sub>			<del>mem</del>	x	x	x	x	mem

To minimise the no. of stalls in pipeline, due to the structural dependency, renaming concept is used.

Renaming states that use the independent memory modules to store the insts & data separate called as code memory & data memory.

Use the code memory in stage 1 & data memory in stage 4. ∴ without conflict, inst 1 fetch accessing the data from data memory and insts 4 & 5 fetched from code memory.

	CC1	CC2	CC3	CC4	CC5	CC6
I <sub>1</sub>	CM	ID	ALU	DM	<del>WB</del>	
I <sub>2</sub>	-	CM	ID	ALU	DM	WB
I <sub>3</sub>	-	-	CM	ID	ALU	DM
I <sub>4</sub>	-	-	-	CM	ID	ALU
I <sub>5</sub>	-	-	-	-	CM	ID
I <sub>6</sub>	-	-	-	-	-	CM



$$A[R_0] \leftarrow A[R_0] + (R_1)B$$

$$M[A+R_0] \leftarrow M[A+R_0] + M[(R_1)B]$$

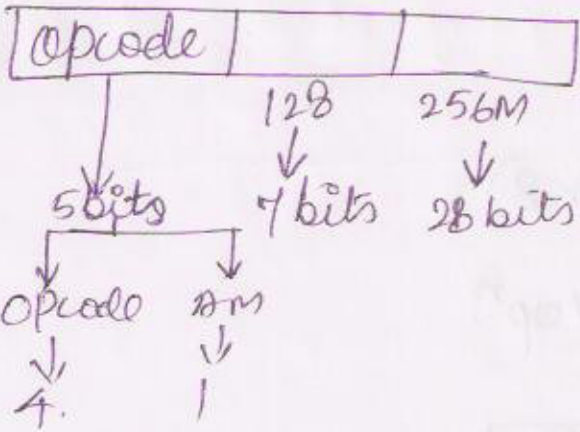
1 reg ref - Index  
1 arith EA  
1 mem ref -> write

1 reg. ref = Index  
1 arith - EA  
1 mem ref -> read

1 mem ref -> EA  
1 mem ref -> read

1 mem ref -> A } 6 mem  
1 mem ref -> B } 9 mem  
ref.

20



7 + 28 = 35  
⑤

23

$$R_s \leftarrow 1$$

$$[R_d \leftarrow 1000 + R_s = 1001 = 1]$$

$$R_d \leftarrow 1$$

$$R_d \leftarrow 1001$$

$$[0 + 1001] \leftarrow 20$$

20  
1001

31

Before replacing the register, take the existing value and compare the remaining program instruction to check whether it is used as i/p or not.

- If it is not used at input, replace register content otherwise verifies the i/p as old i/p or new i/p.
- replace register whenever it is a new i/p and donot replace register when it is old output.

31

$$a = 1$$

$$b = 10$$

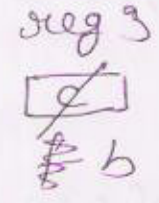
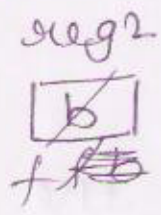
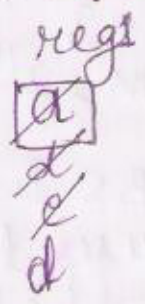
$$c = 20$$

$$d = a + b = 11$$

$$e = c + d = 31$$

$$f = c + e$$

$$b = c + e$$



CISC

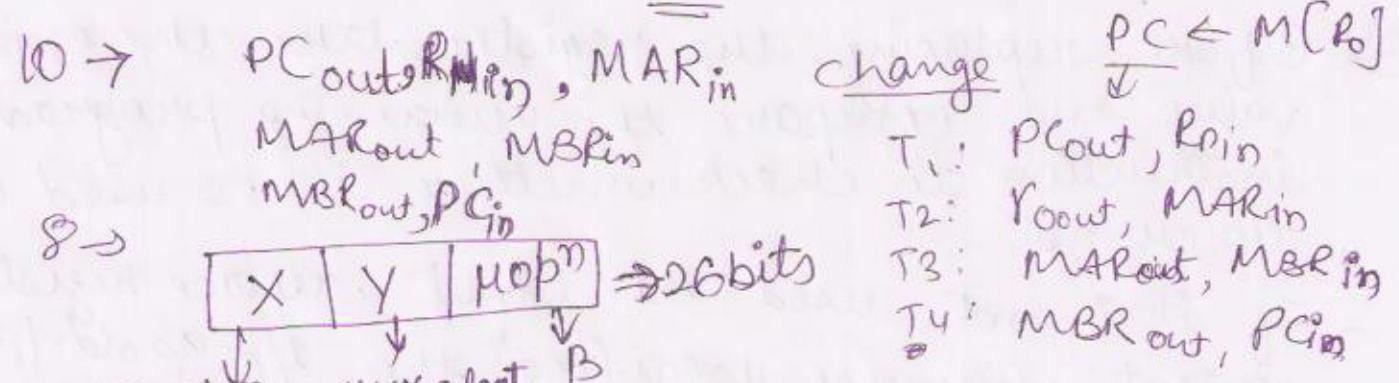
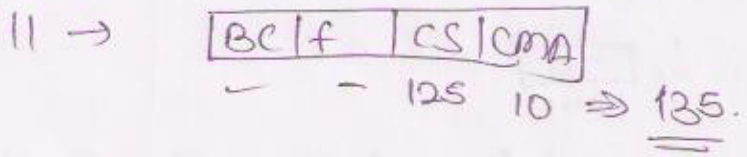
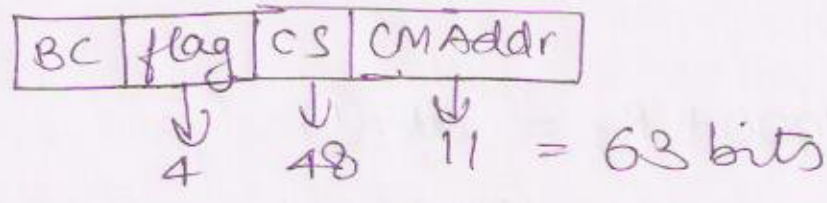
RISC  $\rightarrow$  1 CPI cycles per inst<sup>n</sup>

1  
1  
40  
MC  $\uparrow$

1  
1  
1  
1  
6 times  
MC  $\uparrow$

$$S = \frac{E_{T_{CISC}}}{E_{T_{RISC}}} = \frac{42}{15} = 2.8$$

12. #inst<sup>n</sup> = 256  
 #HOP<sup>n</sup> = 8  
 Total = 2048 (256 x 8)  
 = CM  
 CM size = 2048 HOP<sup>n</sup>  
 HCW



$\Rightarrow X + 3 + 13 = 26 \Rightarrow X = 10$   
 $Y = 3$   
 CM = 2<sup>10</sup> cells

= 1024 CW  
 = 1024 x 26 bits  
 =  $\left(\frac{1024 \times 26}{8}\right)$  bytes

To implement the BBS inst<sup>n</sup> by using the alternative stmt, make sure that mask must contain only 1 bit as 1 that itself based on the position of next bits must be zero. Swap AND OP<sup>n</sup>, correspond bit value is maintained. interm. inst<sup>n</sup>

## 6D) Data Dependency:-

→ Consider the program segment where the instruction  $j$  follows inst<sup>n</sup>  $i$  in the program order i.e.

$i$ : inst<sup>n</sup>  
 $j$ : inst<sup>n</sup>

- Data dependency exists when inst<sup>n</sup>  $j$  tries to read data before inst<sup>n</sup>  $i$  writes it.

→ Data dependency existing between adjacent inst<sup>n</sup> is called true data dependency.

Eg:  $I_1$ : Add  $R_0, R_1, R_2$   
 $I_2$ : Mul  $R_3, R_0, R_4$

- When the above inst<sup>n</sup> are executed in non overlapping order, there is no data loss prob-

- lem because  $I_2$  is inserted after completion of  $I_1$ .

- When above inst<sup>n</sup> are executed in pipeline, inst<sup>n</sup>  $I_2$  is inserted before completion of  $I_1$ . So  $I_2$  tries to read  $R_0$  before  $I_1$  writes it.  
 ∴ old data is injected into the process.

- To make it correct, keep  $I_2$  in waiting until completion of  $I_1$ .

- This waiting creates "stalls" in the ~~program~~ pipeline.

- To detect the data dependency condition, some kind of database will be maintained in the ~~reading~~ stage.

- The dB contains following fields:-

(i) Serial No	Func <sup>n</sup> Unit	Dest <sup>n</sup>	Indep <sup>n</sup> sources	Indep <sup>n</sup> S2	Indep <sup>n</sup> sources	Indep <sup>n</sup> S2	Dep <sup>n</sup> S2	Dep <sup>n</sup> S1	Dep <sup>n</sup> S2	Dep <sup>n</sup> S1
SNO	Fun <sup>n</sup> Unit	Dest <sup>n</sup>	Independent (I)	Ind. S2 (I)	dependent (D)	Dep S2 (D)	(I)	(D)	(D)	(D)
$I_1$	ADD	$R_0$	$R_1$	$R_2$	-	-	-	-	-	0 ⇒ EX Stage
$I_2$	MUL	$R_3$	-	$R_4$	$R_0$ (I)	-	-	-	-	0 ⇒ NOEX Stage

No further fetch are performed because  $I_2$  is not passed to EX stage & there is a stall. The above algo is called "TOMASULO Algo."

Status 0 means the execution is not yet completed. Execution completed  $\Rightarrow$  Status 1.

	CC1	CC2	CC3	CC4	CC5	CC6
$I_1$	IF	ID + S: $r_1, r_2$ (I) (I) D: $r_0$ $r_1 = \text{value}$ $r_2 = \text{value}$	EX	MA	WB	
$I_2$		IF	ID * S: $r_0, r_4$ (O) (I) D: $r_3$ $r_4 = \text{value}$	X	X	EX
$I_3$			IF	X	X	ID
$I_4$						IF

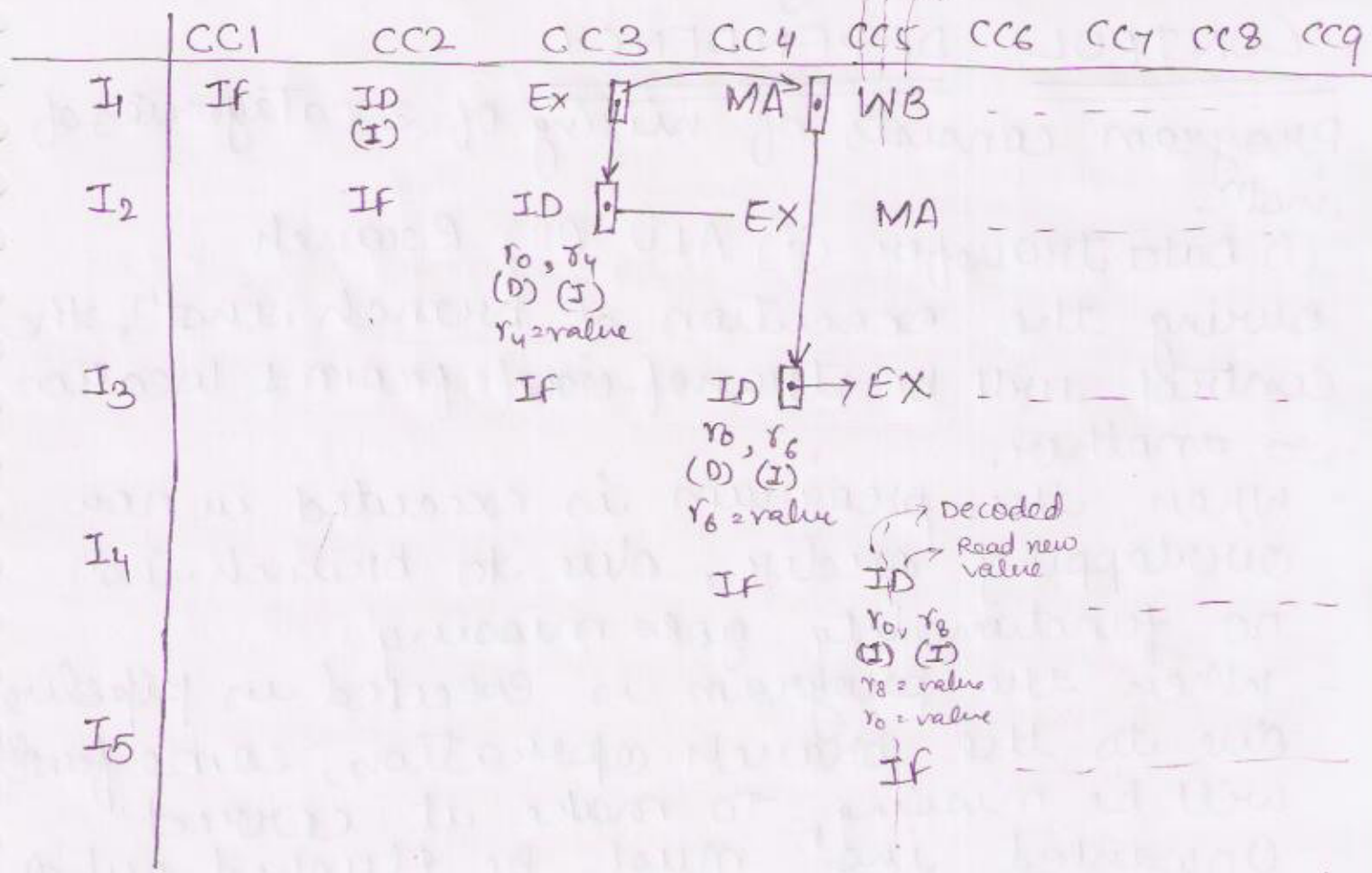
CC4 & CC5 are called "stall cycles" because no new input is accepted due to presence of  $I_3$  in IF stage.

- To minimise the no. of stalls, due to data dependencies, one hw technique is used called as operand forwarding. It is also called as bypassing or short-circuiting.
- Operand forwarding states that access the new data from interface register before updating the register file.

Eg. I<sub>1</sub>: ADD r10, r4, r12  
 I<sub>2</sub>: SUB r3, r0, r4  
 I<sub>3</sub>: MUL r5, r0, r6  
 I<sub>4</sub>: DIV r7, r0, r8

← True data dependency  
 ← dependency but not special name.

→ updation of register files



CC5 divided into two halves. 1st half → Decoding of I<sub>4</sub> & updation of register file as I<sub>4</sub>.  
 2nd half I<sub>4</sub> → reading updated value from register file.

- In the above sequence diagram due to operand forwarding technique, the dependent inst<sup>n</sup> (I<sub>2</sub> & I<sub>3</sub>) are accessing the data from interface registers before updating the register file.  
 ∴ No wait No stall.

- During execution of I<sub>4</sub>, it reads r0 value from register file. At that time, r0 contains the new value because CC5 cycle is divided into 2 halves. In the first half of CC5

I<sub>1</sub> updates R<sub>0</sub> & I<sub>4</sub> decodes ins<sup>n</sup>.  
 In the second half of C5 I<sub>4</sub> reads the register file. At that time R<sub>0</sub> contains new value only.

## CONTROL DEPENDENCY

Program consists of mixing of 3 categories of instr<sup>n</sup>.

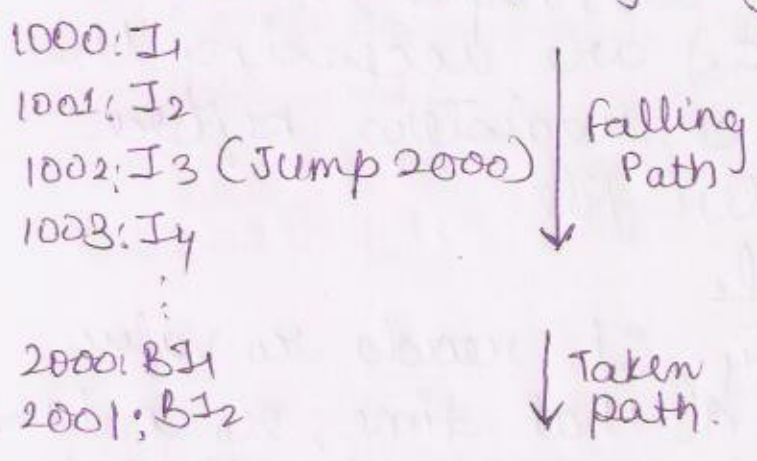
- (i) Data transfer (ii) ALU (iii) Branch.

During the execution of branch instr<sup>n</sup>, the control will be transferred from 1 location to another.

- When the program is executed in non overlapping order, due to branch instr<sup>n</sup>, no functionality goes missing.
- When the program is executed in pipeline due to the branch operation, some funct<sup>n</sup> will be missing. To make it correct.
- Unwanted instr<sup>n</sup> must be flushed out of the pipeline. This operation creates the stalls in the memory pipeline.

## UNCONDITIONAL BRANCH:

Consider the prog segment.



Execution seq.  
 I<sub>1</sub>-I<sub>2</sub>-I<sub>3</sub>-BI<sub>1</sub>-BI<sub>2</sub>-...



	CC1	CC2	CC3	CC4	CC5	CC6
PC=1000 I <sub>1</sub>	IF PC=1001	ID	EX	MA	WB	
I <sub>2</sub>		IF PC=1002	ID	EX	MA	WB
I <sub>3</sub>			IF PC=1003	ID Uncond. TOC TA=2000 PC=1004 2000	EX	MA
needs to be flush/ freeze ↑ I <sub>4</sub> unwanted				<del>IF</del> <del>PC=1004</del> stall	<del>ID</del>	
BI <sub>1</sub>					IF PC=2001	ID
BI <sub>2</sub>						IF PC=2002

I<sub>3</sub> completes at CC4 because uncond. TOC  
completes at IO stage.

1. Inst<sup>n</sup> fetch stage is overlapping with Inst<sup>n</sup> decodes stage. Before decoding fetched inst<sup>n</sup> the next sequential inst<sup>n</sup> is already inserted into the pipeline.

If decoding stage decodes the fetched inst<sup>n</sup> as data transfer or data manipulation then next sequential addr becomes the wanted instruction.

When decoding stage decodes the inst<sup>n</sup> as uncond. TOC then the next sequential inst<sup>n</sup> becomes the "unwanted" inst<sup>n</sup>.

To maintain the correct functionality while executing the program, there is a need of removing the unwanted inst<sup>n</sup> from the pipeline.

The process of removing the unwanted inst<sup>n</sup> from the pipeline is called as flush.

or flush.  
 • The flush operation creates the stalls in the pipeline.

• The no. of stalls cycles created due to branch op<sup>n</sup> is called as "branch penalty".

$$\text{Branch penalty} = \text{At what stage the target Addr - 1 is available}$$

NOTE: The risk pipelining branch penalty is always 2 because the target addr is available at 2<sup>nd</sup> stage.

The no. of stall cycles created during the execution of program, due to branch operation is calculated as

$$\# \text{ stalls from branches} = \text{branch frequency} \times \text{branch penalty}$$

(During Prog. Ex. T).

CONDITIONAL BRANCH:-

Consider the program segment.

1000 : I<sub>1</sub>  
 1001 : I<sub>2</sub> (JNZ r<sub>0</sub>, 2000)  
 1002 : I<sub>3</sub>  
 1003 : I<sub>4</sub>  
  
 2000 : B<sub>1</sub>  
 2001 : B<sub>2</sub>

with true  
 I<sub>1</sub> I<sub>2</sub> B<sub>1</sub> B<sub>2</sub>

with false  
 I<sub>1</sub> I<sub>2</sub> I<sub>3</sub> I<sub>4</sub> ...

	CC1	CC2	CC3	CC4	CC5	CC6
I <sub>1</sub>	IF PC=1001	ID	EX	MA	WB	
I <sub>2</sub>		IF PC=1002	ID Cond. TOC NZCMP <sub>0</sub> ↓ (T) PC: <del>1003</del> 2000		EX	MA
I <sub>3</sub>			<del>IF PC=1003</del> Stall		ID	EX
BI <sub>1</sub>				IF PC: 2001		ID
BI <sub>2</sub>					IF PC: 2002	

- When decoding stage decodes fetch<sup>inst<sup>n</sup></sup> as cond. TOC with true condition then next sequential inst<sup>n</sup> becomes unwanted inst<sup>n</sup>. So, flushout the "unwanted" inst<sup>n</sup>. It creates stalls.

- If cond<sup>n</sup> evaluates to false, the next <sup>seq</sup> inst<sup>n</sup> becomes "wanted" inst<sup>n</sup>, so no flush no stall.

To minimise the no. of stalls in pipeline due to branch operation, one h/w technique is used i.e. branch prediction buffer or Branch target buffer or loop buffer. (BBB/BTB/LB).

BTB:- BTB is a high speed buffer maintained at the inst<sup>n</sup> fetch stage to predict the target addresses. Branch target buffer consists of following ~~fields~~ fields.

PC entry	Expected PC
----------	-------------

During execution of the program when CPU identifies stalls in the pipeline due to the true condition, (first addr is target Addr). Then the flush out Inst<sup>n</sup> address and target address both are inserted into the PC entry and expected PC fields respectively.

→ From the next fetch onwards the PC value is compared with the PC entry field.

→ If none of the values matches, fetch the inst<sup>n</sup> based on current PC, if any one entry is matching then fetch the inst<sup>n</sup> based on the corresponding expected PC value.

```

1000: MOV r0, #4
1001: ADD r1, r2
1002: DJNZ r0, 1001
1003: MOV

```

```

PC: 1000: r0 = 4
    1001: +
    1002: r0 = 3 (N2) ⇒ T
    flush 1003 1001: +
    1002: r0 = 2 (N2) ⇒ T
    No flush 1003
    No stall 1001: +
    No flush 1003 1001: +
    1002: r0 = 1 (N2) ⇒ T
    No flush 1003 1001: +
    1002: r0 = 0 (N2) ⇒ F
    No flush 1003
    flush 1003 1003

```

PC entry	PC Expected
1003	1001
<del>1001</del>	

Useful only in loops.

1st flush → insertion } in expected - entry table.  
 last flush → deletion

When there is a stall due to the false cond<sup>n</sup> delete the ~~entry~~ expected PC value from the buffers & fetch the next inst<sup>n</sup> based on the ~~current~~

corresponding PC value.

## DELAYED BRANCH :-

- It is a compiler technique. It is used to preserve the execution part in the pipeline.
- When there is an unwanted inst<sup>n</sup> present in the pipeline that slot is fixed as the delayed slot.
- When the slot is fixed as delayed then it substitutes the NOP operation.

	CC1	CC2	CC3	CC4	CC5	CC6
I <sub>1</sub>	IF	ID *	EX	MA	WB	
I <sub>2</sub>		IF	ID load	EX	MA	WB
I <sub>3</sub>			IF	ID Uncond TOC	EX	MA
I <sub>4</sub>				IF(NOP)	ID	EX
BT <sub>1</sub>					IF	ID
BT <sub>2</sub>						

CC1 is never delayed slot (1<sup>st</sup> inst<sup>n</sup> in the falling path) <sup>IF</sup>

CC2 is a delayed slot

CC3 " " " " → Forwarded because of I<sub>1</sub>, I<sub>2</sub>

CC4 " " " " → (Fixed) → Uncond TOC

CC5 " never delayed slot (1<sup>st</sup> inst<sup>n</sup> in taken path)

CC6 " delayed slot (depends on decoding of I<sub>0</sub>)  
(wanted or unwanted depends on decoding unit)

## INSTRUCTION SCHEDULE :-

The processor executes the program in sequential order called as inorder execution.

In the inorder execution sequence, if any inst<sup>n</sup> is data dependent on other inst<sup>n</sup>, the

remaining inst<sup>n</sup> are also sharing the stall cycles of the dependent inst<sup>n</sup>.

Eq:  $I_1$ : Add  $r_0, r_1, r_2$       Add  $r_0, r_1, r_2$   
 $I_2$ : Mul  $r_3, r_0, r_4$       Mul  $r_3, r_0, r_4$   
 $I_3$ : Sub  $r_4, r_5, r_6$       Sub  $r_5, r_5, r_6$   
 $I_4$ : Div  $r_3, r_7, r_8$       Div  $r_7, r_7, r_8$

In order Execution sequence: -  $I_1 - I_2 - I_3 - I_4$

- In the above example  $I_2$  is data dependent on  $I_1$ , so it is undergoes wait until completion of  $I_1$ . Due to this,  $I_3$  &  $I_4$  are also sharing stall cycles even they are independent.

- To avoid the above problem, there is a need of scheduler instructions.

Scheduling causes the out of order execution (Reorder exe<sup>n</sup>).  $I_1 - I_3 - I_4 - I_2$  (Out of order exe<sup>n</sup>).

- The out of order exe<sup>n</sup> creates 2 more dependencies in the pipeline

(i) Antidependency

(ii) Output dependency

- Antidependency is existed when the inst<sup>n</sup>  $j$  tries to write the data before inst<sup>n</sup>  $i$  reads it.

eg:-  $I_3$  modifies  $r_4$  before  $I_2$  reads it.

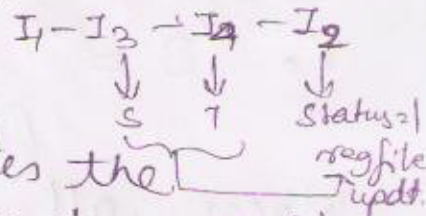
- Output dependency exists when inst<sup>n</sup>  $j$  tries to write the data before inst<sup>n</sup>  $i$  writes it.

eg:-  $I_4$  executed before  $I_2$ , so  $I_4$  writes  $r_3$  before  $I_2$  writes it.

To handle the anti and output dependencies register renaming is used. Register renaming states that use some temp. regis storage (reorder buffer) to store the out of order exe<sup>n</sup> outputs. After receive the accept signal from the dependent inst<sup>n</sup>, update the reg. file with reorder buffer contents.

### HAZARDS:-

- Hazard is a delay. Delay creates the extra cycles. Extra cycle without operation is called stall.



~~But~~ There are 3 types of data hazards.

- (i) RAW Hazard (Read after write)
- (ii) WAR Hazard (write after read)
- (iii) WAW Hazard (write after write)

RAW hazard is created when the inst<sup>n</sup> 'j' tries to read the data before inst<sup>n</sup> 'i' writes it (true data dependency).

WARH is created when inst<sup>n</sup> j tries to write the data before inst<sup>n</sup> i reads it (anti-dependency).

WAWH is created when inst<sup>n</sup> j tries to write data before inst<sup>n</sup> i write it (output dependency).

### Performance evaluation of pipelining with stalls:-

$$S = \frac{\text{Avg ins}^n \text{ exe}^n \text{ of non-pipe}}{\text{Avg ins}^n \text{ exe}^n \text{ of pipe}} = \frac{\text{Avg Ins}^n E_{T_{np}}}{\text{Avg Ins}^n E_{T_p}}$$

$$= \frac{CPI_{np} * \text{cycle time}_{np}}{CPI_p * \text{cycle time}_p}$$

The ideal CPI is almost 1 (always).

$$S = \frac{CPI_{np} * cycle_{time}_{np}}{(1 + \# \text{ stalls/inst}^n) * C_{tp}}$$

↓                    ↓  
Ideal due to dependencies

When all stages are perfectly balanced, then both cycle times are equal. ←

$$\therefore S = \frac{CPI_{np}}{(1 + \# \text{ stalls/inst}^n)} \quad [\because C_{np} = C_p]$$

When all inst<sup>n</sup> take same # cycles, then 1 inst<sup>n</sup> exe<sup>n</sup> is also equal to # stages in pipeline.

$$\therefore S = \frac{\text{pipeline depth}}{1 + \# \text{ stalls/inst}^n}$$

When system operates with 100% efficiency then no stalls exist.

$$\Rightarrow \boxed{S = \text{pipeline depth}}$$

Pg 17 Q2.

cycle time = 10 ns.

Branch freq = 20%.

Branch penalty = 5 - 1 = 4 cycles.

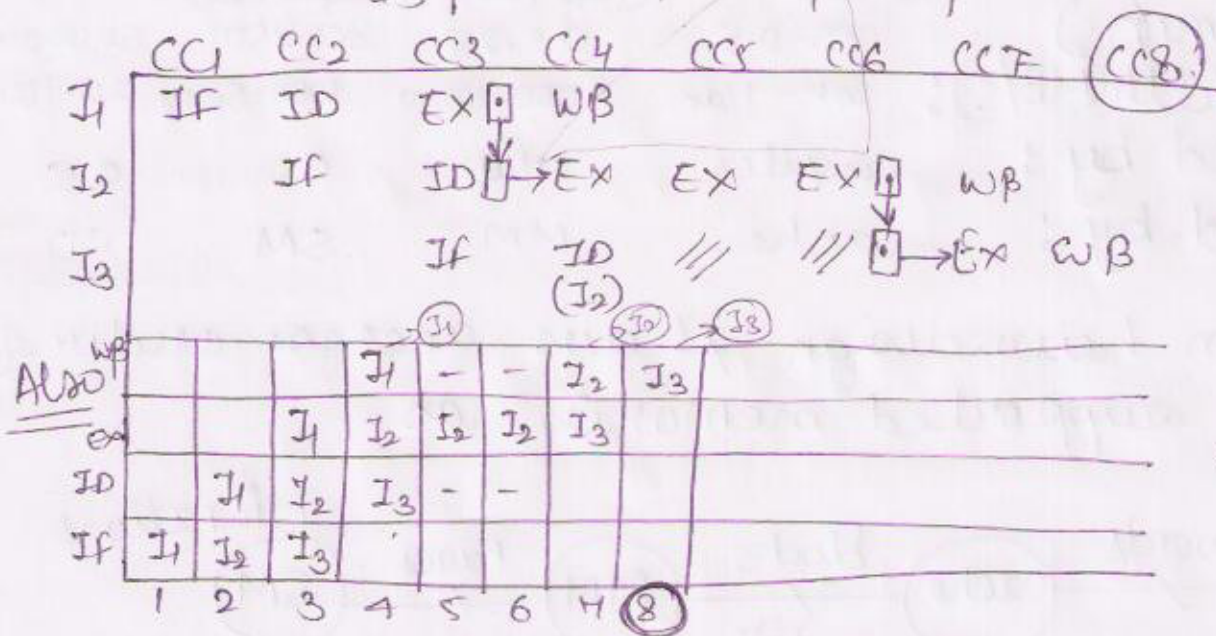
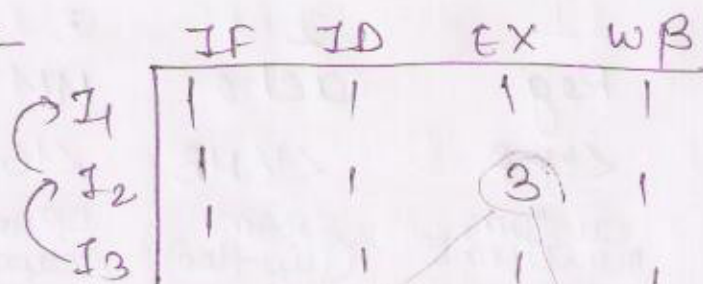
$$\begin{aligned} \text{Avg inst}^n \text{ ET} &= (1 + \# \text{ stalls/inst}^n) * \text{cycle time} \\ &= \left[ 1 + (\text{Branch freq} * \text{Branch penalty}) \right] * \text{cycle time} \\ &= (1 + (0.2 * 4)) * 10 \text{ ns} \\ &= 18 \text{ ns}. \end{aligned}$$

1 inst<sup>n</sup> → 18 ns

$$\begin{aligned} 1 \text{ S} &= \frac{1}{18} * 10^9 \text{ 1ns/sec} \\ &= 55 \text{ MIPS} \end{aligned}$$



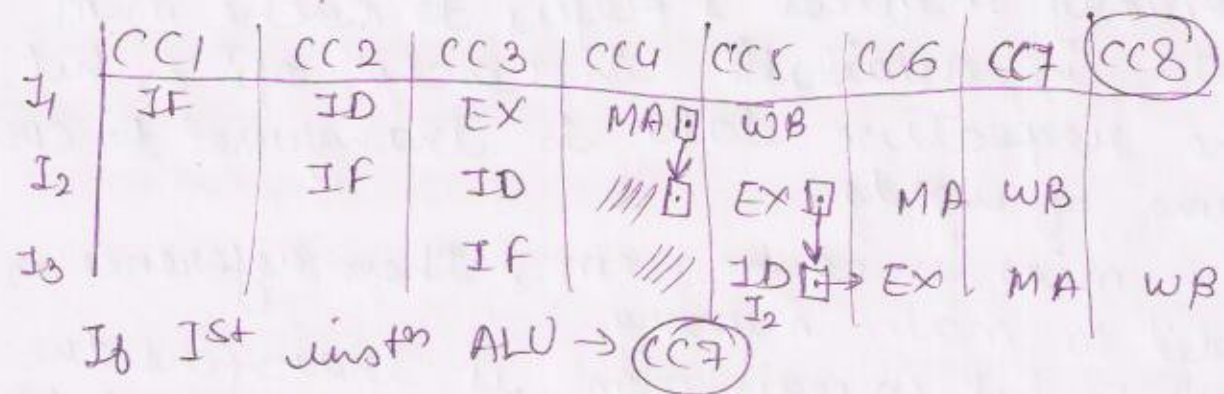
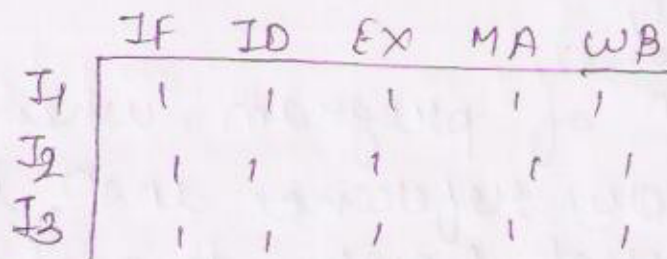
Q3 Pg 17 :-



← Cycle diagram

← space-time diagram

Q7 Pg 18 :-



MEMORY ORGANIZATION :-

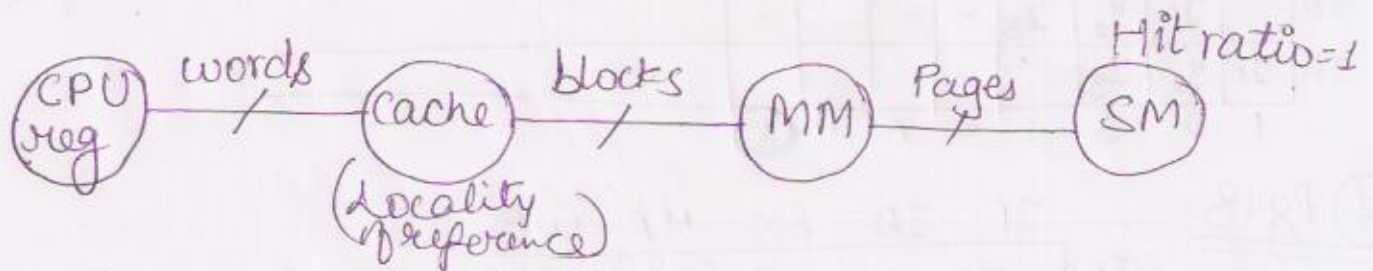
Memory org<sup>n</sup> shows effective utilization of system supported memories.

To satisfy that obj. memory hierarchy design is used.

Memory Hierarchy uses the following specifications :-

level:	1	2	3	4
Name:	Reg.	Cache	MM	SM
Typical Size:	<1KB	<16MB	<16GB	>100GB
Implementation:	custom multiport	SRAM (flip-flop)	DRAM (capacitor)	magnetic
Access time (ns):	0.25-0.5	0.5-25	80-250	50,000,000
Bandwidth (MB/s):	20K-11ac	5000-10000	1000-5000	20-150
Managed by:	compiler	H/W	OS	OS
Backed by:	cache	MM	SM	CD

The mem hiera. design shows access order of system supported memories as



- During the exe<sup>n</sup> of program, when CPU encounters memory reference ins<sup>n</sup>, it generates memory request & prefers to cache mem.
- If data is available in cache, op<sup>n</sup> is hit. So, the respective data is transferred to CPU in terms of words.
- If op<sup>n</sup> is miss in cache mem., then reference is forwarded to main memory.
- When op<sup>n</sup> is hit in main mem., it transfers the data to the cache memory in the form of blocks and cache to CPU in form of words. Otherwise, ref. is forwarded to SM.
- The operation is always hit in SM because it is the last level of memory in hierarchy structure. ∴ hit ratio is always 1 in SM.
- ∴ data is transferred from SM to MM in terms of pages, MM to cache in blocks & from cache to CPU in the form of words.
- According to hierarchy design, the data is Xfered from higher to lower. ∴ lower level data is always

a subset of higher level data.

→ This property is also called as inclusion.

According to heiran design, the CPU can perform the read & write op<sup>n</sup> only the cache memory.

- cache memory is a <sup>reusable</sup> storage which holds the image of main memory so CPU performs op<sup>n</sup> on reusable space. This property is called as locality of reference.

### Types of memory organizations:-

Based on way of accessing the system supported memories, the memory org is divided into 2 types:-

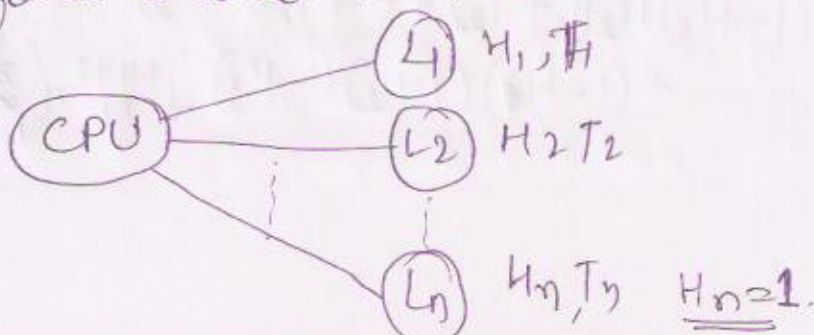
- (i) Simultaneous Access memory org.
- (ii) Hierarchical access " " "

### Simultaneous access:-

In this org. the CPU is able to access all the levels of memories. In this org. level to level comm. is not allowed.

- CPU refers the level 1 memory to read or write the data.

If there is a miss operation, the reference is forwarded to level 2. If op<sup>n</sup> is hit in level 2 i.e. directly transferred to the CPU without the involvement of level 1. otherwise reference is forwarded to next level & so on.



Where  $H_1, H_2, \dots, H_n$  are hit ratios of respec. memories.

$$\text{Hit ratio} = \frac{\text{no. of hits}}{\text{Total \# accesses}}$$

$T_1, T_2 \dots T_n$  are access time of respec. mem.

- The time required to access 1 operand from memory is called as average access time.

$$T_{\text{avg}} = H_1 T_1 + (1-H_1) H_2 T_2 + (1-H_1)(1-H_2) H_3 T_3 + \dots + (1-H_1)(1-H_2)(1-H_3) \dots (1-H_{n-1}) H_n T_n.$$

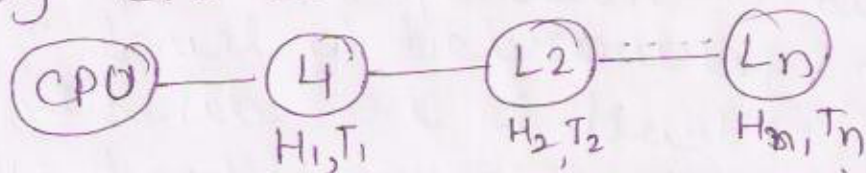
1 word =  $T_{\text{avg}}$ .

? words = 1 sec.

$$\Rightarrow \eta, \text{ efficiency} = \frac{1}{T_{\text{avg}}} \text{ words/sec}$$

(Always when cache is used)  
Hierarchical Access :-

In this org. the CPU can perform read & write op<sup>n</sup> only on level 1 memory. If there is a miss in level 1, transfer the data from higher level to level 1, then only CPU can ~~ac~~ read or write.



The avg access time of hierarchical access org. is

$$T_{\text{avg}} = H_1 T_1 + (1-H_1) H_2 (T_1 + T_2) + (1-H_1)(1-H_2) H_3 (T_1 + T_2 + T_3) + \dots + (1-H_1)(1-H_2)(1-H_3) \dots (1-H_{n-1}) H_n \left( \sum_{i=1}^n T_i \right)$$

Q. In a 2-level memory, level 1 memory is 5 times faster than level 2 memory. And its access time is 10ns less than avg access time. Let level 1 access time be 20ns. What is the hit ratio?

⇒ (Simultaneous)

$L_1$  &  $L_2$

④ ⇒ speed up factor is given

$$S = \frac{T_2}{T_1} = 5$$

← faster

$$\Rightarrow T_2 = 5T_1$$

$$\textcircled{!} \Rightarrow T_1 = T_{avg} - 10$$

$$= T_{avg} = T_1 + 10$$

$$T_1 = 20 \text{ ns}$$

$$\therefore T_{avg} = 30 \text{ ns}$$

$$\& T_2 = 100 \text{ ns}$$

$$T_{avg} = H_1 T_1 + (1-H_1) T_2$$

$$= H_1 T_1 + (1-H_1) T_2$$

$$30 = H_1 * 20 + (1-H_1) 100$$

$$H_1 = 0.87 \text{ Ans}$$

last level hit ratio = 1 (always)

Q. 3 level memory has following specifics

level	Access time/word	Block size (word)	hit ratio
1	20	—	0.8
2	50	2	0.9
3	100	4	1

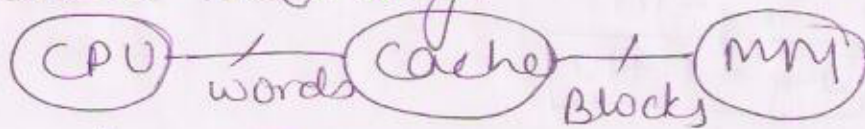
If the required block is not present in  $L_1$ , transfer it from  $L_2$  to  $L_1$ , if not available in  $L_2$ , transfer from  $L_3$  to  $L_2$  &  $L_2$  to  $L_1$ . How long will it take to transfer the blocks

Ans ⇒  $T_1 = 20 \times 1 = 20$  &  $T_3 = 100 \times 4 = 400$   
 $T_2 = 50 \times 2 = 100$

$$\begin{aligned}
 T_{avg} &= H_1 T_1 + (1-H_1) H_2 (T_1 + T_2) + (1-H_1)(1-H_2) H_3 (T_1 + T_2 + T_3) \\
 &= 0.8 \times 20 + 0.2 \times 0.9 (120) + (0.2)(0.1) 1 (520) \\
 &= 16 + 21.6 + 10.40 \\
 &= \underline{\underline{48 \text{ ns}}}
 \end{aligned}$$

## CACHE MEMORY:-

Cache memory is used as the intermediate memory between CPU + main memory.  
 $\therefore$  CPU perform read & write op only on the cache memory.



- Basic elements are
- (i) Memory organization
  - (ii) Mapping techniques
  - (iii) Replacement algorithms
  - (iv) Updating techniques
  - (v) multilevel caches

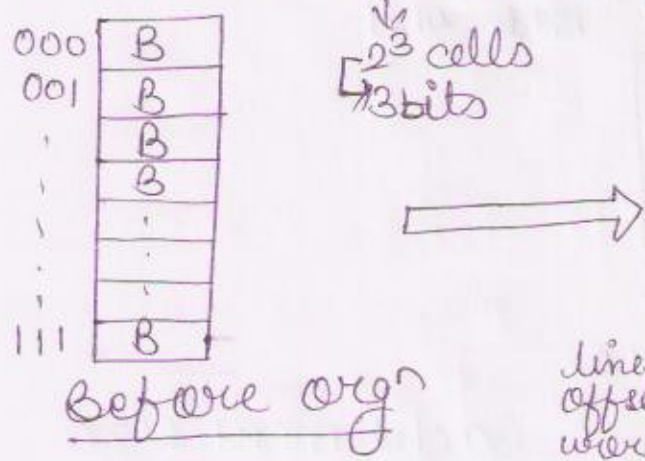
Memory organization: The data is transferred from ~~one~~ main memory to cache memory in form of blocks. So both memories are divided into equal parts based on block size. Each part is called as cache mem. block + main memory block respectively.

$$\# \text{ cache mem. blocks (lines)} = \frac{\text{cache mem size}}{\text{block size}}$$

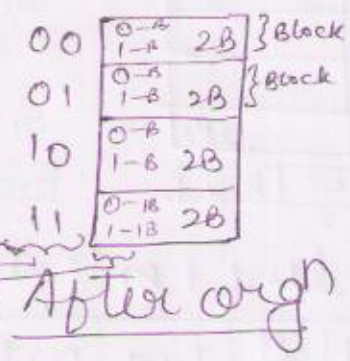
$$\# \text{ main memory blocks} = \frac{\text{MM size}}{\text{Block size}}$$

Consider 8 byte cache memory & 2 byte block size.

$CM = 8B = 8 \times B$

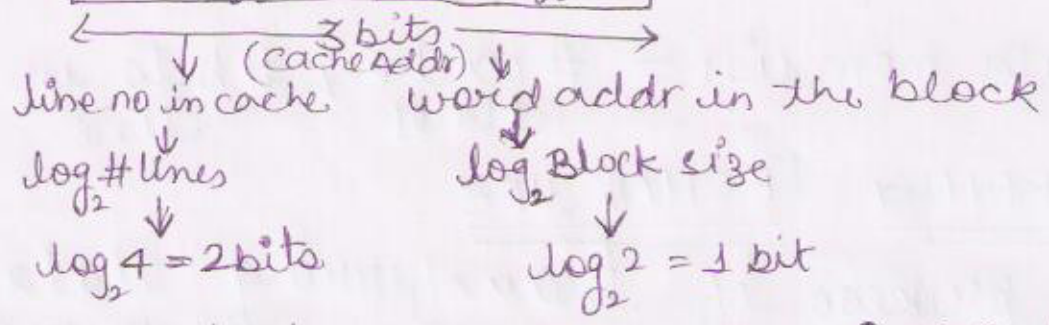
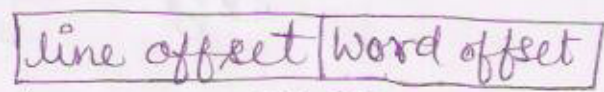


# lines =  $\frac{8B}{2B} = 4$



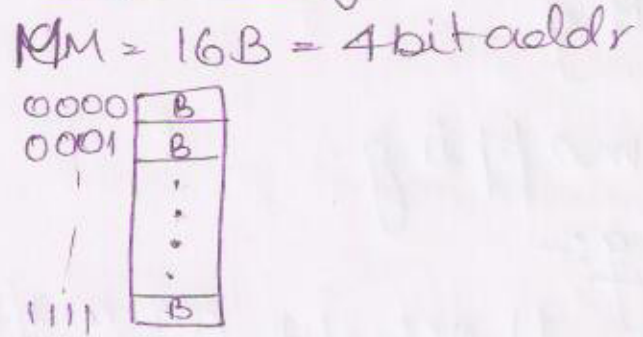
line offset  
word offset

Before and after organization of cache memory there is no change in capacity but the internal structure differs. After cache mem. org., address can be interpreted as

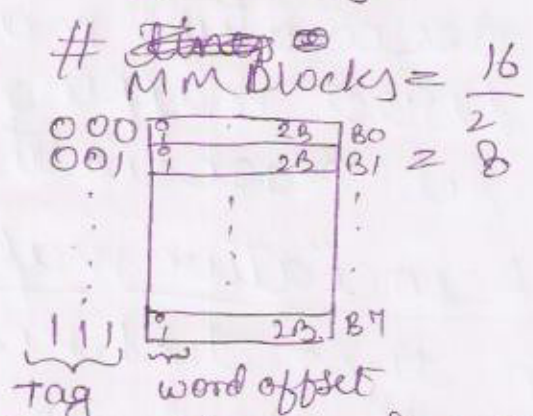


Consider 16 byte main memory & 2B block size.

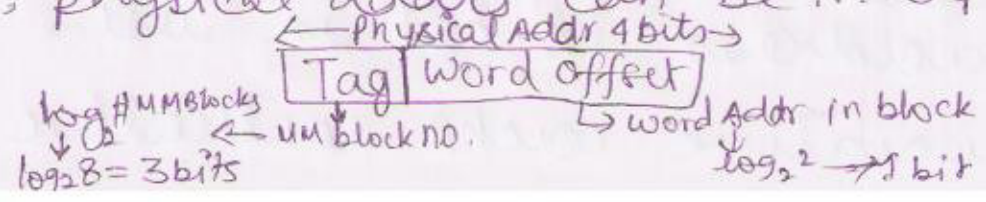
Before org<sup>n</sup>

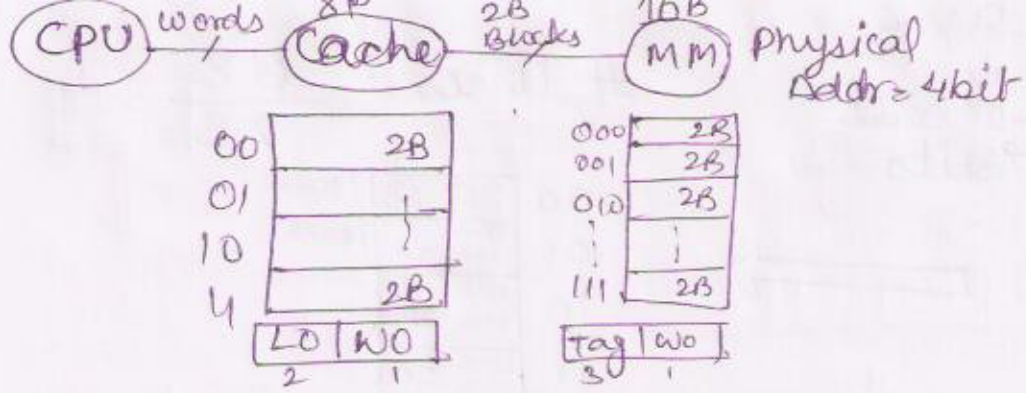


After org<sup>n</sup>

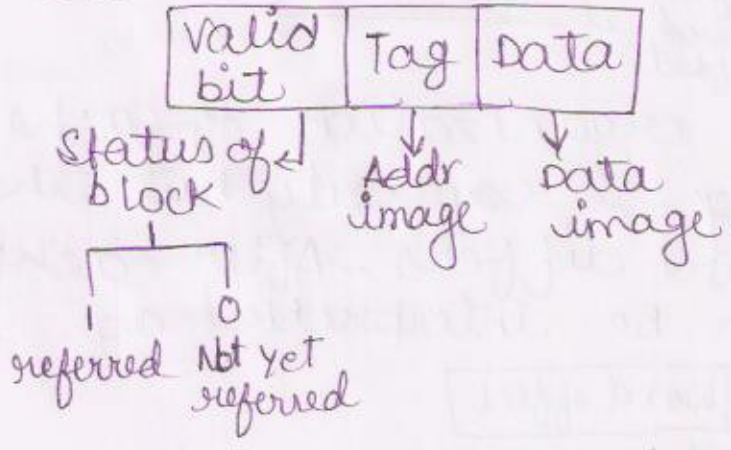


Before & after org<sup>n</sup> of the MM, there is no change in capacity but the internal structure differs so, physical address can be interpreted as





The cache line contents are



Cache mem size = tag mem size + data mem size

Tag mem size = # cache lines \* # tag bits in each line

Data mem size = # cache lines \* # bits in each line

### MAPPING TECHNIQUES

The process of transferring data from MM to cache memory is known as mapping.

There are 3 mapping techniques used:-

- (i) Associative mapping
- (ii) Direct mapping
- (iii) Set associative mapping.

#### Associative mapping:-

In this technique, there is no mapping funct<sup>n</sup> used to transfer the data. That means any MM block can be mapped to any cache line.

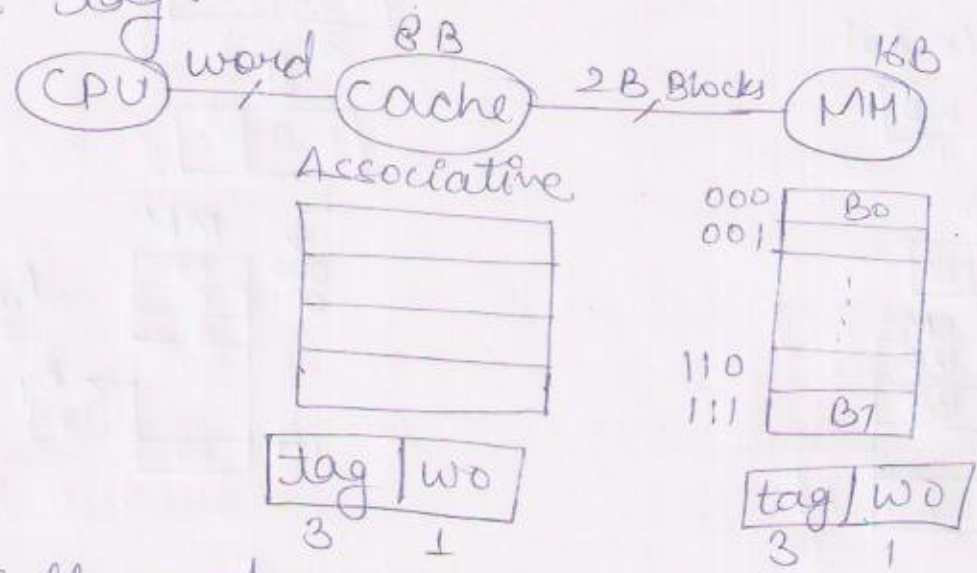
∴ Cache address is not in use.

- The associative cache controller intro-

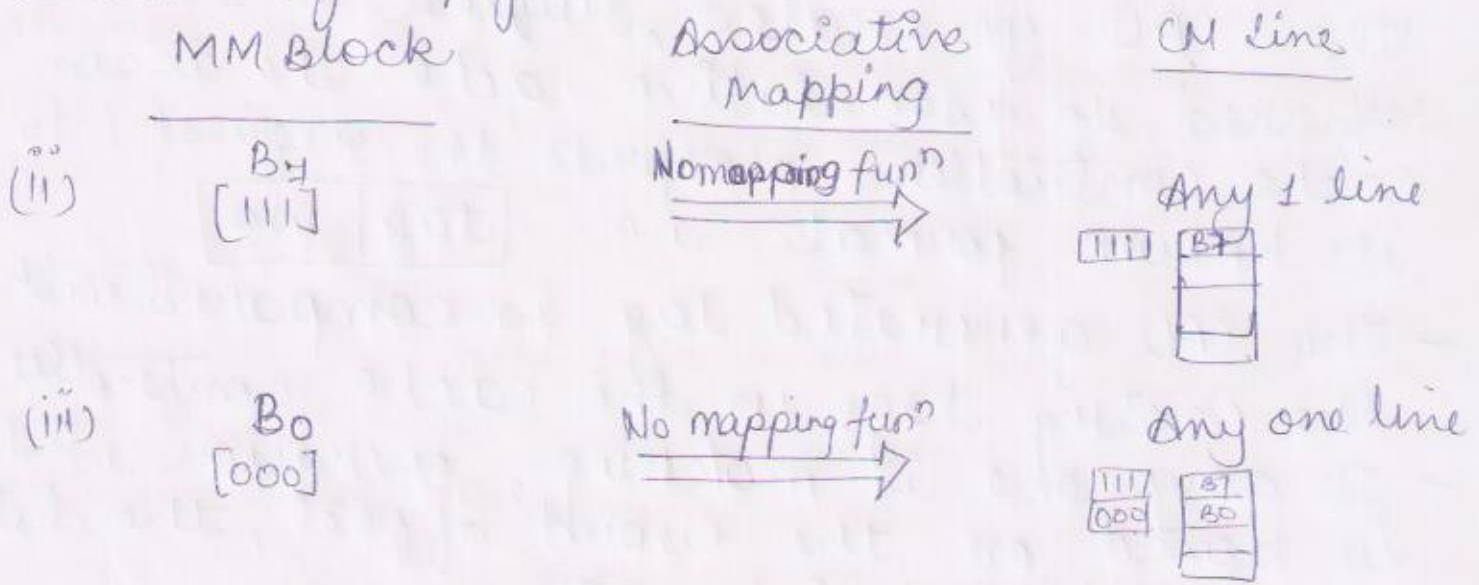


-prets the CPU requests based on MM address format.

-During the mapping, the complete data block is transferred to CM along with tag.



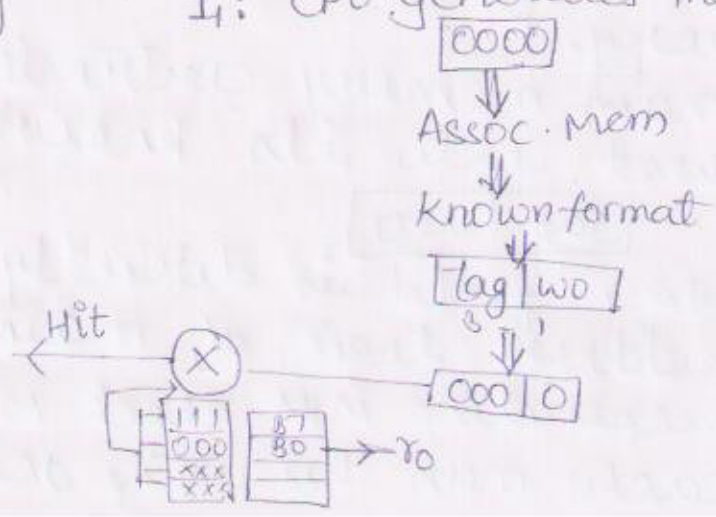
(i) initially empty cache.



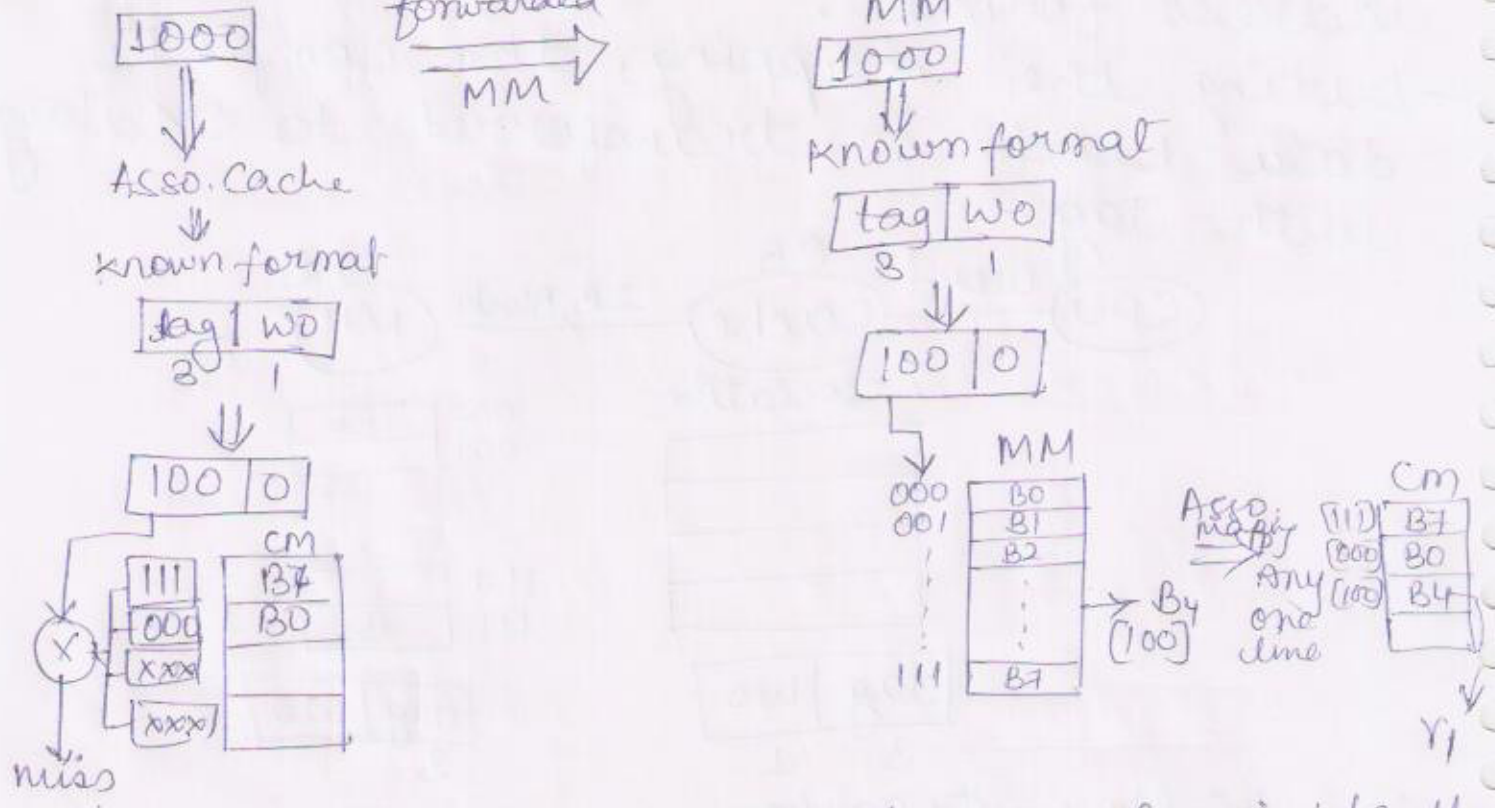
Consider the program segment :-

- (i) MOV r0, [0000]
- (ii) MOV r4, [1000]
- (iii) Add r0, r4

I<sub>1</sub>: CPU generates mem req.



I<sub>2</sub> CPU generates the mem req.



The CPU generated <sup>mem</sup> request is initially referred to associative cache. The asso. cache controller interprets the request into its known format i.e. tag | wo

- The CPU generated tag is compared with the existing tags in the cache controller.
- If any one is matching, operation is hit. So based on the word offset, the data is transferred into CPU.
- If none of them is matching op<sup>n</sup> is miss. So the reference is forwarded to main memory.
- The main memory controller interprets the request into its known format i.e. tag | wo
- The tag field is directly connected with the address logic of main memory. So corresponding MM block is enabled & transferred to cache mem. by using associative mapping.

Technique later data is transferred to CPU. The tag memory size is equal to # of cache lines \* # tag bits in each line.

Tag mem size =  $4 * 3 = 12 \text{ bits}$

DIRECT Mapping:

In this technique memory funct<sup>n</sup> is used to transfer the data. The mapping function is

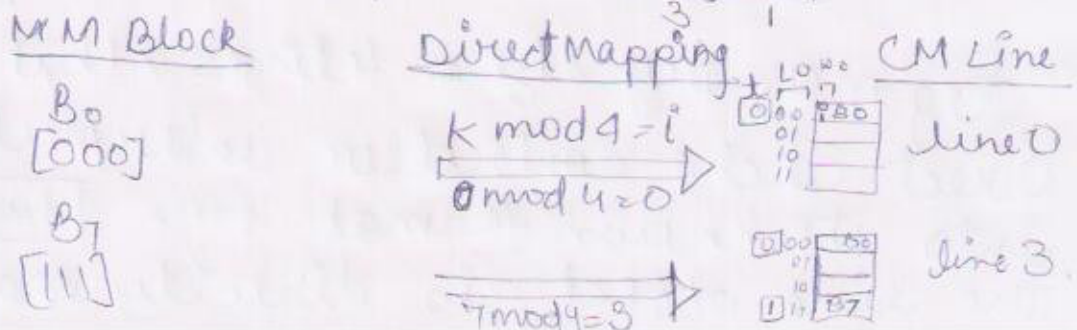
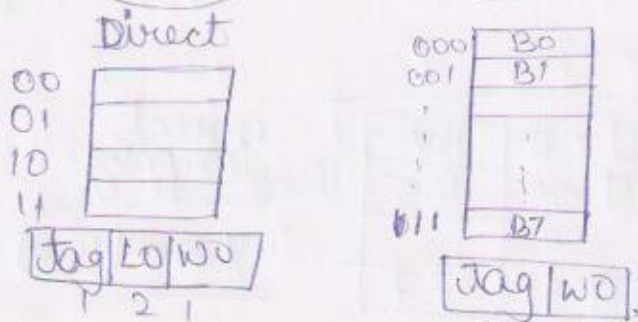
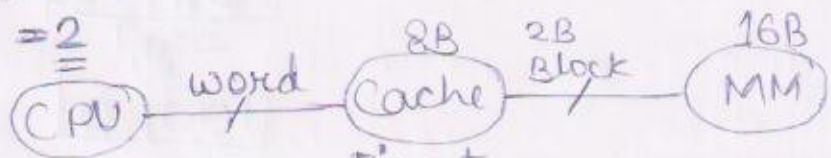
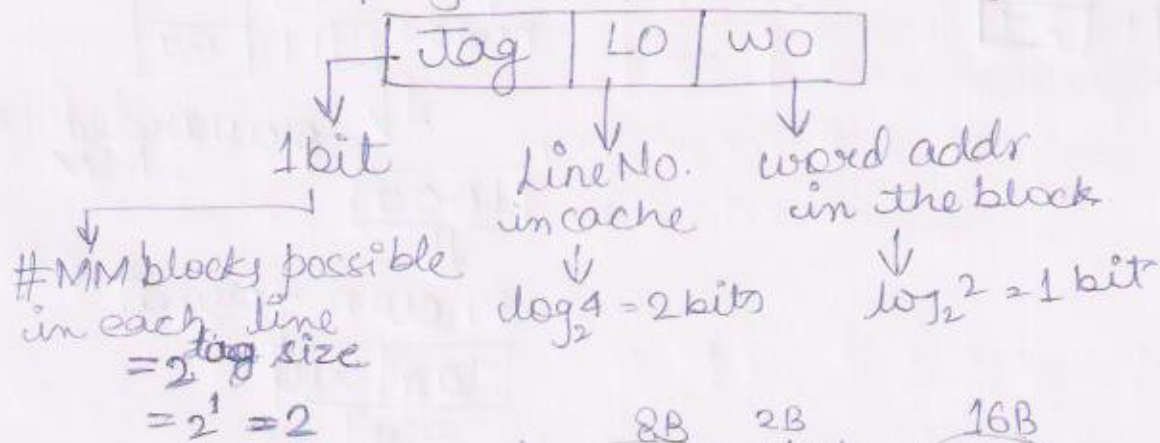
$$K \text{ mod } N = i$$

$\downarrow$                        $\downarrow$                        $\downarrow$   
 MM Block#            # cache lines            CM line number  
 $(0-7) \text{ mod } 4 = (0-3 \text{ \& } 0-3)$

The above mapping fun<sup>n</sup> shows the relationship b/w mm block & cache mem line.

∴ cache mem. address is in use.

The direct cache controller interprets the CPU generated requests as physical Addr (4 bits)



Consider the program segment :-

I<sub>1</sub>: MOV R<sub>0</sub>, [0000]

I<sub>2</sub>: MOV R<sub>4</sub>, [1001]

I<sub>3</sub>: Add R<sub>0</sub>, R<sub>4</sub>

I<sub>1</sub>: CPU generates mem req

I<sub>2</sub>: CPU generates mem req.

0000

Direct cache

known format

tag	LO	WO
1	2	1

0	00	0
---	----	---

only tag compared  
X  
Hit

0	00	B0
1	01	B5
	10	
1	11	B7

→ R<sub>0</sub>

1001

Direct cache

known format

tag	LO	WO
1	2	1

1	00	1
---	----	---

0	00	B0
1	01	B5
	10	
1	11	B7

miss

forwarded to MM

1000

known format

tag	WO
3	1

100	0
-----	---

000	B0
001	B1
	⋮
111	B7

1	00	B4
1	01	B5
	10	
1	11	B7

Direct mapping  
line 0 ← B4  
 $4 \div 4 = 0 [100]$

tag memory size = 4 lines × 1 bit = 4 bits.

- Direct cache controller interprets the req. into its known format i.e., tag|LO|WO
- The line offset is directly mapping with

the address logic of cache memory. ∴ ~~correct~~  
 - pending line is enabled.

- The existing tag in enabled line is compared with CPU generated tag.

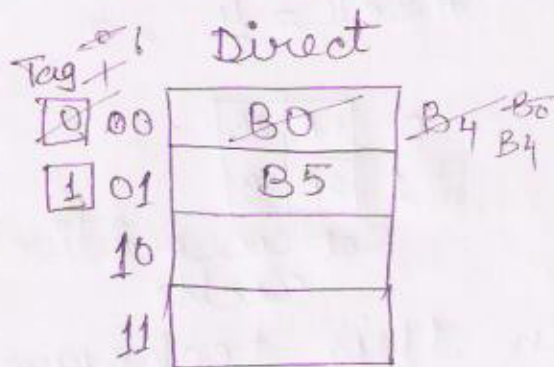
- when both matches, the op<sup>n</sup> is hit. So based on the word offset data is transferred to CPU. otherwise it is a miss.

Reference is forwarded to MM.

- The required block is transferred from MM to cache memory based on the direct mapping function.

(70)

Disadvantage in direct mapping is that each cache line is able to hold one block at a time. ∴ Conflict misses will be increased.



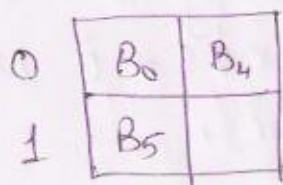
I<sub>1</sub>: MOV R<sub>0</sub>, [0000] ⇒ B<sub>0</sub>, 0<sup>th</sup> Byte, Hit

I<sub>2</sub>: MOV R<sub>0</sub>, [1000] ⇒ B<sub>4</sub>, 0<sup>th</sup>, Miss

I<sub>3</sub>: MOV R<sub>0</sub>, [0001] ⇒ B<sub>0</sub>, 1<sup>st</sup> Byte, Miss

I<sub>4</sub>: MOV R<sub>0</sub>, [1001] ⇒ B<sub>4</sub>, 1<sup>st</sup>, Miss

### Alternative cache organization



I<sub>1</sub>: Hit  
 I<sub>2</sub>: Miss  
 I<sub>3</sub>: H  
 I<sub>4</sub>: H

SET-ASSOCIATIVE  
 CACHE  
 ORGANIZATION

To overcome the disadvantage of direct mapping, alternative cache organization is used. In this, each cache line is able to hold more than 1 block at a time. The alternative cache organization is called a SET ASSOCIATIVE CACHE.

# SET ASSOCIATIVE CACHE ORGANIZATION

In this organization, cache memory lines are again divided into equal parts based on no. of main memory blocks can be placed at a time at each cache line.

Each part is called set. No. of sets,  $S = \frac{\text{no. of cache lines}}{\text{\# blocks placed at a time in each line}}$

$$\frac{N}{P\text{-ways}} = S$$

# blocks placed at a time in each line.

Consider 2-way set associative cache, with block size 2 bytes.

Before org<sup>n</sup> (BB)

000	B
001	B
010	.
011	.
100	.
101	.
110	.
111	B

After org<sup>n</sup> into the lines

000	2B
001	2B
010	2B
011	2B

# lines =  $\frac{8}{2} = 4$   
Direct cache

After org<sup>n</sup> into the sets

$$\# \text{ sets} = \frac{4}{2} = 2$$

0	2B	2B
1	2B	2B

Set associative cache.

Set Associative Mapping: In this technique,

mapping function is used to transfer the data from main memory to cache.

The mapping function is known as  $k \bmod S = i$  where

$k$  = main memory block number.

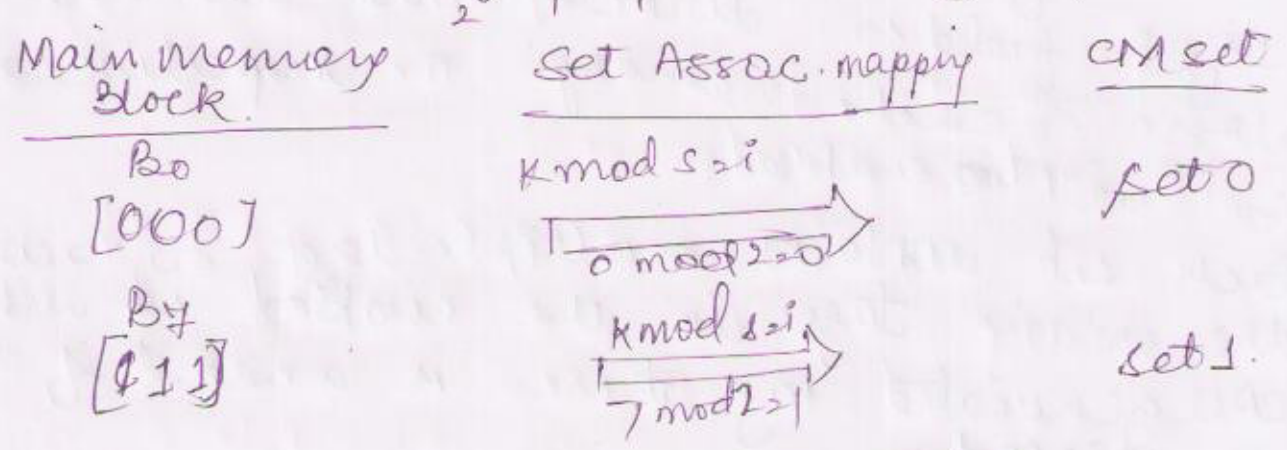
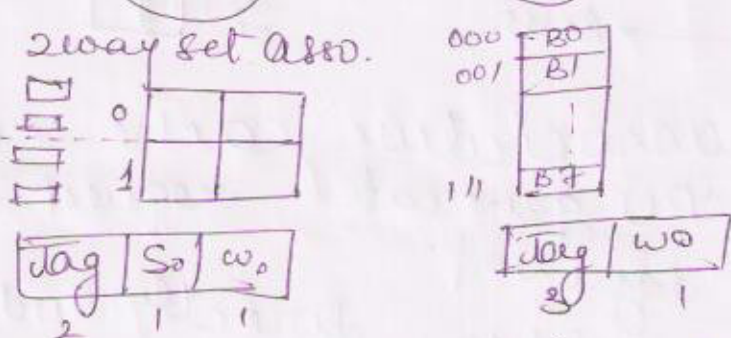
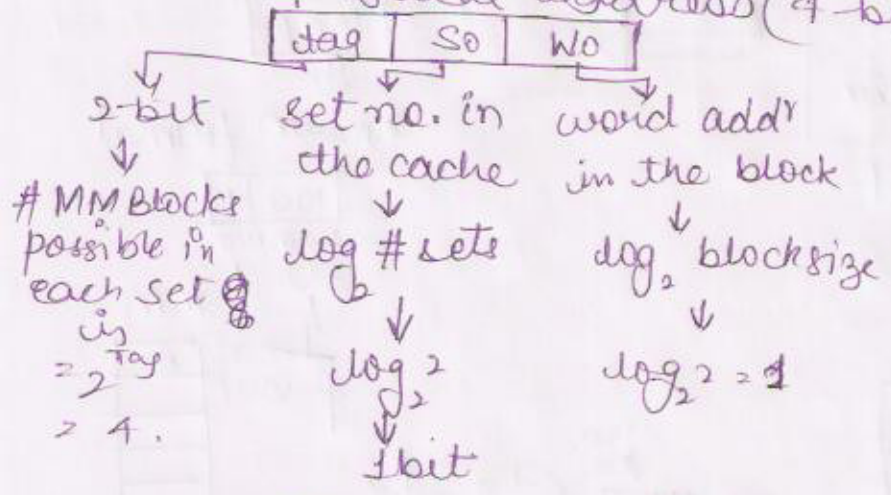
$S$  = # cache sets

$i$  = cache memory set no.

The above function shows the relationship b/w MM block & Cache memory address. Cache memory address is in the use.

The set associative cache controller interprets the CPU generated ~~requests~~ address request

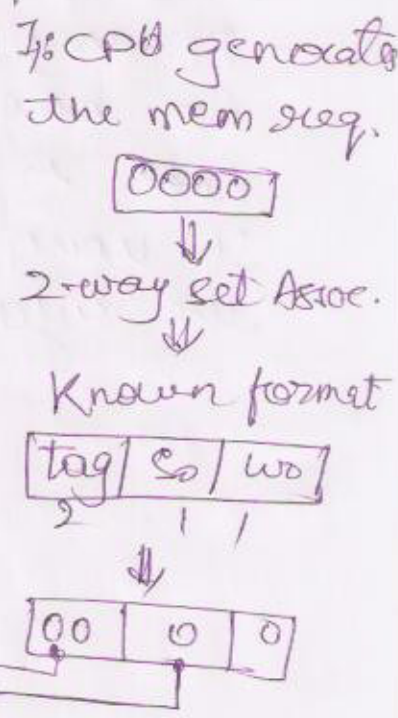
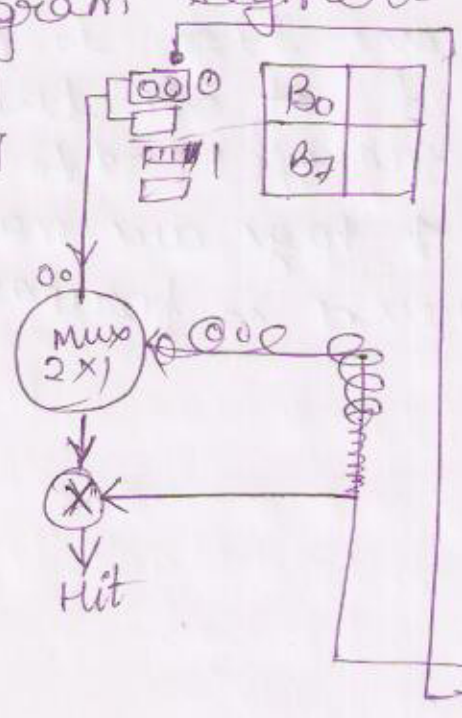
into the following format :-  
 physical address (7-bit)



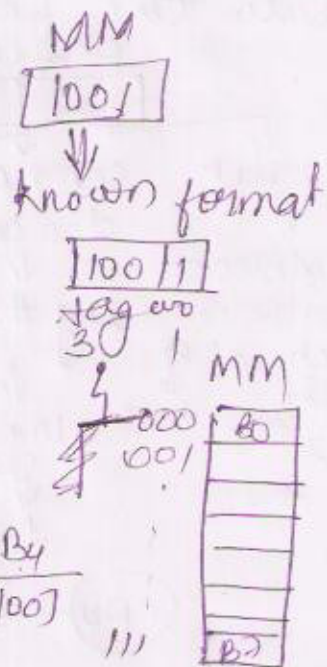
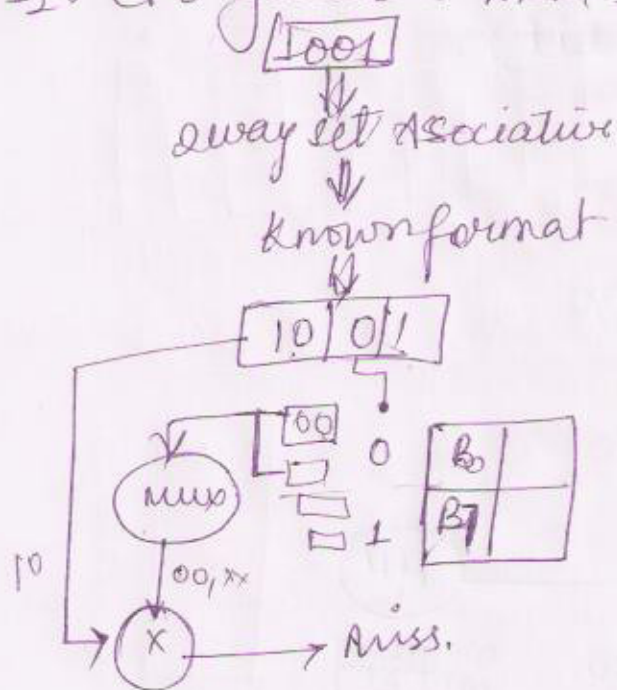
Consider the program segment:

```

I1: MOV r0, [0000]
I2: MOV r4, [1000]
I3: Add r0, r4
  
```



I<sub>2</sub>: CPU generates mem req. forwarded



Set associative cache controller interprets the CPU generated request as follows:-

tag | way |

offset field is directly with the address logic of cache memory. The respective cache set will be enabled.

Each set contains multiple tags. So to compare the existing tags in the enabled set with the CPU generated tag there is a need of multiplexer.

- If any one tag is matching, op<sup>n</sup> is hit. So based on the word offset, the data is transferred to the CPU.

If none of tags are matching, op<sup>n</sup> is miss so the reference is forwarded to main memory.



Based on the tag value, the respective block is transferred to the cache memory by using the set associative mapping technique.

$$\text{tag memory size} = \# \text{ sets} * \# \text{ Blocks in each set} * \# \text{ tag bits.}$$

$$= 5 * 2 * 2 = 20 \text{ bits}$$

$$= 2 * 2 * 2 = 8 \text{ bits}$$

## REPLACEMENT ALGO

When the cache is full, there is a need of replacement algorithms to replace the existing cache blocks with new blocks.

There are 2 replacement algos used.

(i) FIFO (ii) LRU.

In first in first out, replace the cache block with new block which is having the longest time stamp.

In least recently used algo, replace the cache block with new block which is having the less no. of references with longest time stamp.

Q. Consider 4 block cache memory (initially empty) with the following MM block references.

4, 5, 7, 12, 4, 5, 13, 4, 5, 7

Identify the hit ratio with:-

(i) FIFO (ii) LRU (iii) Direct mapped cache

(iv) 2 way set assoc. with LRU.

(M)	(M)	(M)	(M)	(M)	(M)	(M)	(M)	(M)	(M)
4	4	4	4	4	4	13	13	13	13
	5	5	5	5	5	4	4	4	4
		7	7	7	7	7	7	7	7
			12	12	12	12	12	12	7
4	5	7	12	4	5	13	4	5	7

$$\text{Hit ratio} = \frac{2}{10} = 0.2$$

(ii) LRU

4				4				4						
	5				5				5					
		7				13								7
			12											
4	5	7	12	4	5	13	4	5	7					

$$\text{hit ratio} = \frac{4}{10} = 0.4$$

(iii) Direct mapped

0	4/12/4
1	5/13/5
2	
3	7

m, m, m, m, m, h, m, h, m, h

$$\text{hit ratio} = \frac{3}{10} = 0.3$$

(iv) 2-way set asso with LRU :-

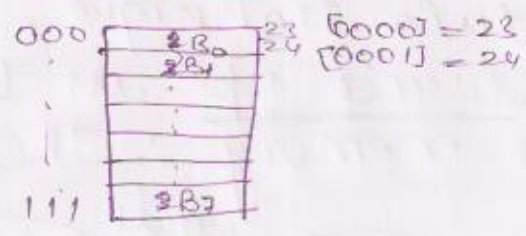
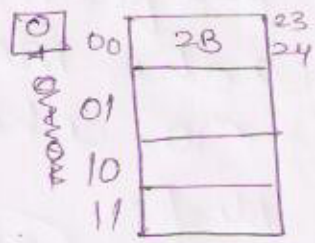
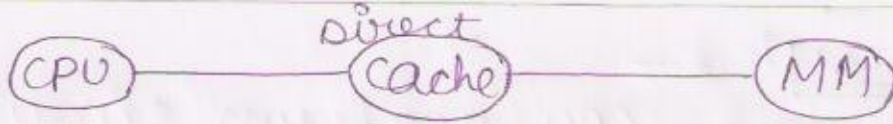
4 5 7 12 4 5 13 4 5 7  
m, m, m, m, h, h, m, h, h, m

0	4	12
1	5	7/13/7

$$\text{hit ratio} = \frac{4}{10} = 0.4$$

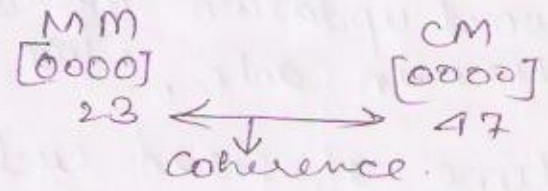
### UPDATING TECHNIQUES :-

CPU performs read & write operation only on the cache ~~memory~~. Before operation, it checks for ~~the~~ availability of block in cache. If the block is available, operation is called as read hit or write hit. If the block is not available, operation is called as read miss/write miss. When there is a miss operation, the required block is transferred from MM to cache memory. This process is called read allocate/write allocate. After the allocate, CPU performs read/write operation.



```

I0: MOV R0, [0000] ; B0, 0th byte → R0, M R0 ← 23
I1: Add R0, #24 ; R0 ← R0 + 23 ; R0 = 47
I2: MOV [0000], R0 ; B0, 0th Byte → R0 ⇒ CM, [0000] ← 47
  
```



During the execution of above program, inst<sup>n</sup> 3 indicates mem. wrt op<sup>n</sup>. So, CPU perform write op<sup>n</sup> only on the cache mem. Therefore, different values are maintained in cache memory + main memory with the same address.

- The same address contains different values at different places is called is Coherence

→ Coherence creates the loss of data problems, i.e.

```

I3: MOV R0, [1000] ; B4-0th byte ⇒ M: CM = B0 B4
I4: —
I5: —
I6: MOV R4, [0000] ; B0-0th Byte ⇒ M: CM = 0 mod 4 = 0
  
```

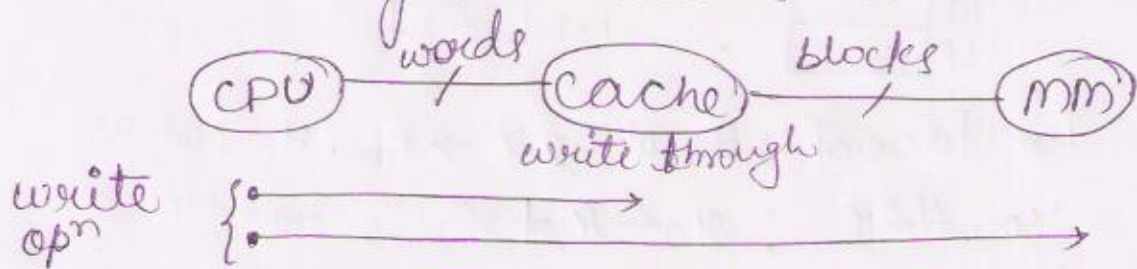
$4 \bmod 4 = 0$   
 $CM = B0$   
 ↓  
 updated data lost  
 $\Rightarrow B4 B0$   
 $23$  (old)  
 $24$

During the execution of I<sub>4</sub>, due to the miss operation, the updated cache block is replaced with new block. So the updated information is lost. To avoid this problem, updating techniques are used:-

- (i) write through
- (ii) write back

(i) Write Through :-

In this technique, CPU perform simultaneous write op<sup>n</sup> in the cache memory and main memory. ∴ There is no coherence.



$$\text{word updation time } (T_w) = \max(\text{word updation time in cache, word updation time in MM})$$

The average access time of read cycle is

$$T_{\text{avg read}} = H_r T_c + (1 - H_r) (T_m + T_c)$$

↓   ↓   ↓   ↓   ↓  
 read hit   read   read miss   read allocated   read.

The average access time of write cycle is

$$T_{\text{avg write}} = H_w * T_w + (1 - H_w) (T_m + T_w)$$

↓   ↓   ↓   ↓   ↓  
 write hit   word updation time   write miss   write allocate   word updation time.

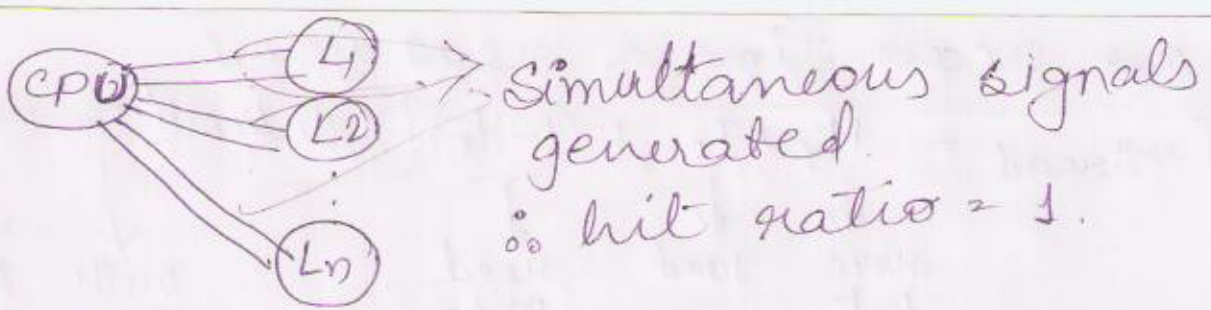
The total average access time when considering both read & write is

$$T_{\text{avg WT}} = f_r * T_{\text{avg read}} + f_w * T_{\text{avg write}}$$

↓   ↓  
 frequency of read op<sup>n</sup>   frequency of write op<sup>n</sup>.

$$n_{\text{WT}} = \frac{1}{T_{\text{avg WT}}} \text{ words/sec}$$

NOTE:- Write through is also ~~write~~ applicable to simultaneous access memory organization.



$$T_{avg_w} = H_w T_m \Rightarrow \boxed{T_{avg_w} = T_m}$$

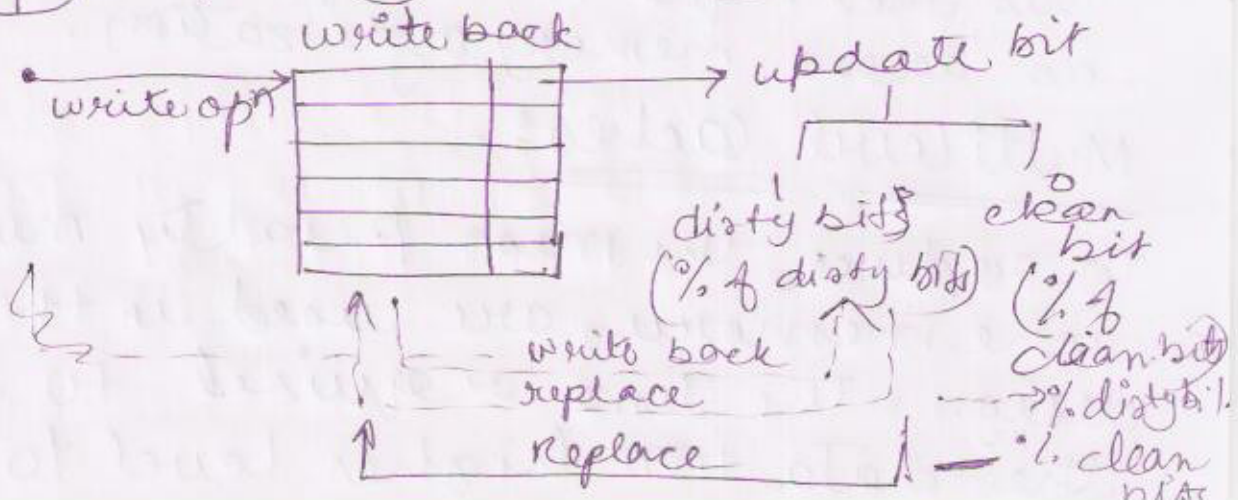
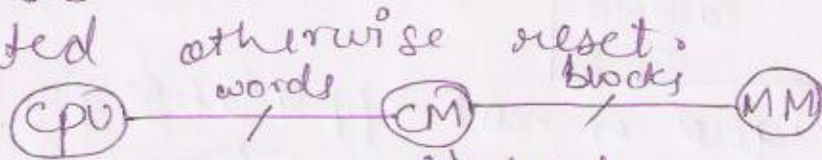
Note: If  $H_w = 1 \Rightarrow$  (simultaneous mem.org.)  
 $H_w \neq 1 \Rightarrow$  (Hierarchical " " " " )

Write back:

In this technique, the CPU performs write operations only in the cache memory so still coherence is present but this coherence does not create the data loss problem because the CPU reads status of cache block before replace.

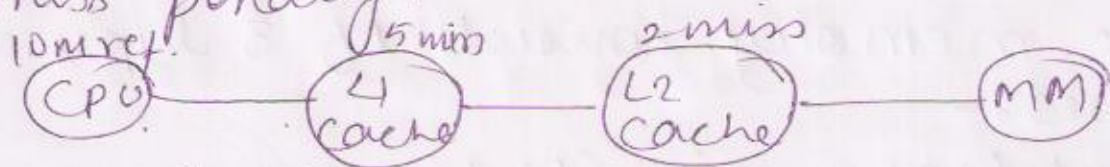
If the cache block is updated block then, it transfers the information back to main memory. Later replaces the cache block. Otherwise the new block will be placed into the cache.

In write back cache, each line maintains one extra bit called as update bit. This bit will be set when the block is updated otherwise reset.





memory due to the miss operation is called as miss penalty.



In multilevel cache org<sup>n</sup> we can calculate two types of miss-rate

- (i) local miss rate
- (ii) Global miss rate

$$\text{Local miss rate} = \frac{\# \text{ misses in cache}}{\text{Total no. of accesses to that cache.}}$$

$$\text{Global miss rate} = \frac{\# \text{ misses in cache}}{\text{Total \# CPU generated references.}}$$

$$\text{LMR}_{L_1} = 5/10 = 0.5$$

$$\text{LMR}_{L_2} = 2/5 = 0.4$$

$$\text{GMR}_{L_1} = 5/10 = 0.5$$

$$\text{GMR}_{L_2} = 2/10 = 0.2$$

always same.

In multilevel cache org. the avg. access time calculated as.

$$T_{avg} = \text{Hit time}_{L_1} + \text{miss rate}_{L_1} * \text{miss penalty}_{L_1}$$

$$\text{Miss penalty}_{L_1} = \text{Hit time}_{L_2} + \text{miss rate}_{L_2} * \text{miss penalty}_{L_2}$$

$$\text{Miss penalty}_{L_2} = \text{main memory access time}$$

In the multilevel cache org<sup>n</sup>, avg stall cycles per inst<sup>n</sup> is calculated as

$$\text{Avg stalls/inst}^n = \# \text{ misses in } L_1 / \text{inst}^n * \text{Hit time}_{L_2} + \# \text{ misses in } L_2 / \text{inst}^n * \text{miss penalty}_{L_2}$$

## Types of MISSES:-

Cache memory consists of 3 types of misses:-

(i) Compulsory miss (First-reference or cold start miss)

This miss will occur when the very first reference to cache is miss (cache is initially empty).

(ii) Capacity miss:- Due to the small size of cache, it cannot hold all memory blocks required for program, so some blocks are discarded & later retrieved.

(iii) Conflict miss (collision miss/Inferencing miss)

This miss will occur when too many blocks are mapped into the same cache line or cache set.

## Types of Locality of reference:-

(1) Temporal locality:-

When the CPU is referring same word in same block in near future is called as temporal locality.

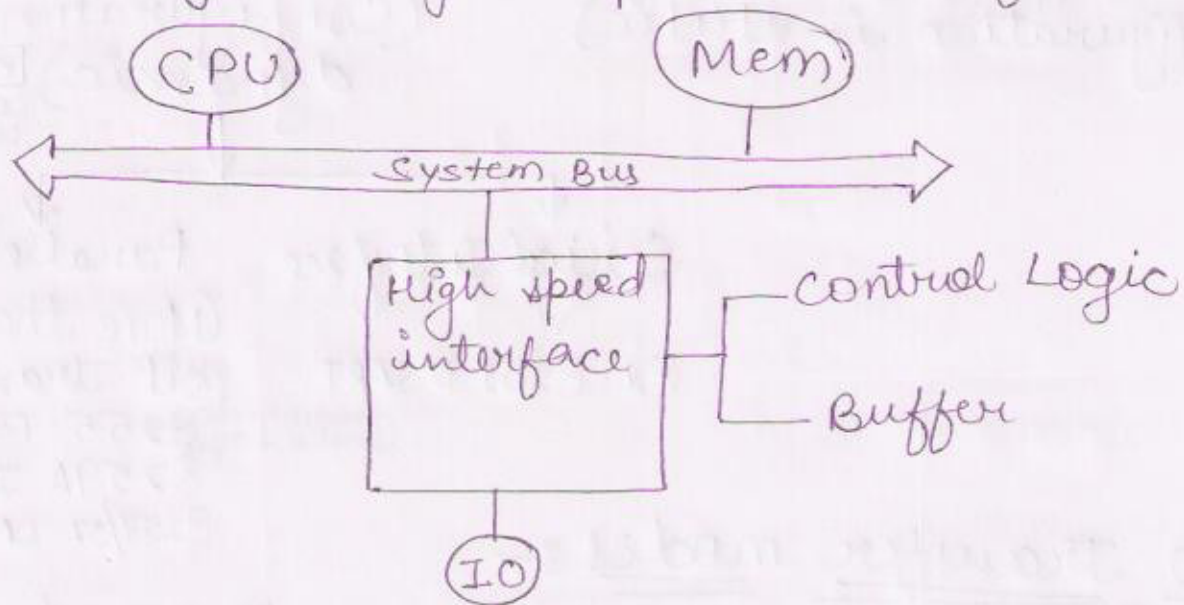
(2) Spatial locality:-

When CPU refers the adjacent words in same block in near future is called as spatial locality.



## IO ORGANIZATION :-

- IO devices are very slow devices. ∴ IO devices are never directly connected with system buses.
- IO devices are always communicating with CPU through high speed interface (IO module).
- High speed interface is required because IO devices are electromagnetic devices and CPU is an electronic component. ∴ there is a difference in operating modes, difference in word format & data transfer rate.
- ∴ To avoid the above problems, IO devices are connected with CPU through high speed interface.



When CPU executes IO operations, it generates IO requests and transfers to high speed interface. Later CPU continues its own some other task. During that process, CPU periodically reads the status of interface.

- High speed interface control logic interprets the request & enables the

corresponding port for respective operation.  
Based on the speed of IO device, it prepares data & later transfers to buffer.

- When buffer has the data, high speed interface generates the signals to CPU & waits for ack.

- After receiving ack, the buffer contents are placed on data lines. ∴ speed gap will be minimised.

### Types of interfaces :-

#### Interface

(i) Non Programmable  
→ Buffer, latch etc  
(Configuration is static)

(ii) Programmable  
USART  
(Configuration is dynamic) CW  
Control word

Serial Interface

Ex: 8251 USART

Parallel Interface

(More than 1 bit transmission)

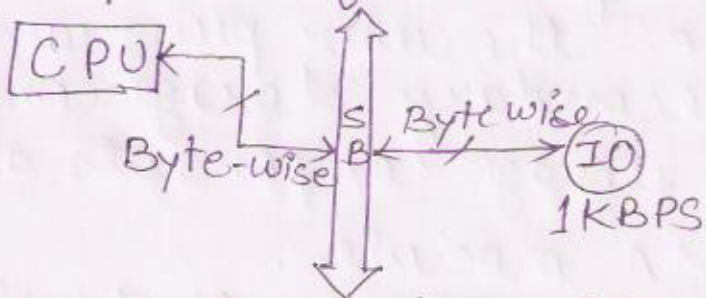
(programmable parallel interface)  
8255 PPI  
8259A INT. Controller  
8257/37 DMA

### IO transfer modes :-

There are 3 different ways present to transfer data from IO to CPU/memory.

- (i) Programmed IO
- (ii) Interrupt driven IO
- (iii) DMA.

Programmed I/O:- In this mode high speed interface is not present. The CPU is directly connected with IO devices.  
 ∴ Processor utilization is inefficient.  
 This means that processor waits until completion of IO operation. This wait depends on speed of IO device.

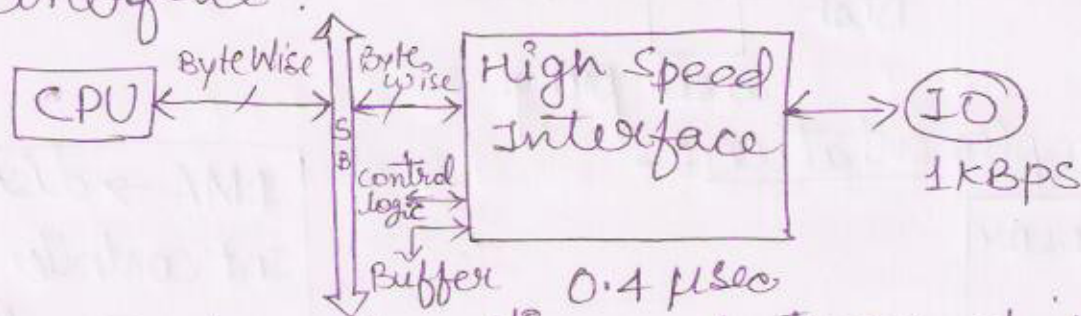


1KB → 1sec  
 1B → 1milliSec  
 Byte transfer time  
 ≥ 1msec.

Interrupt driven IO:- In this mode, high speed interface is used to connect the various basic IO devices.

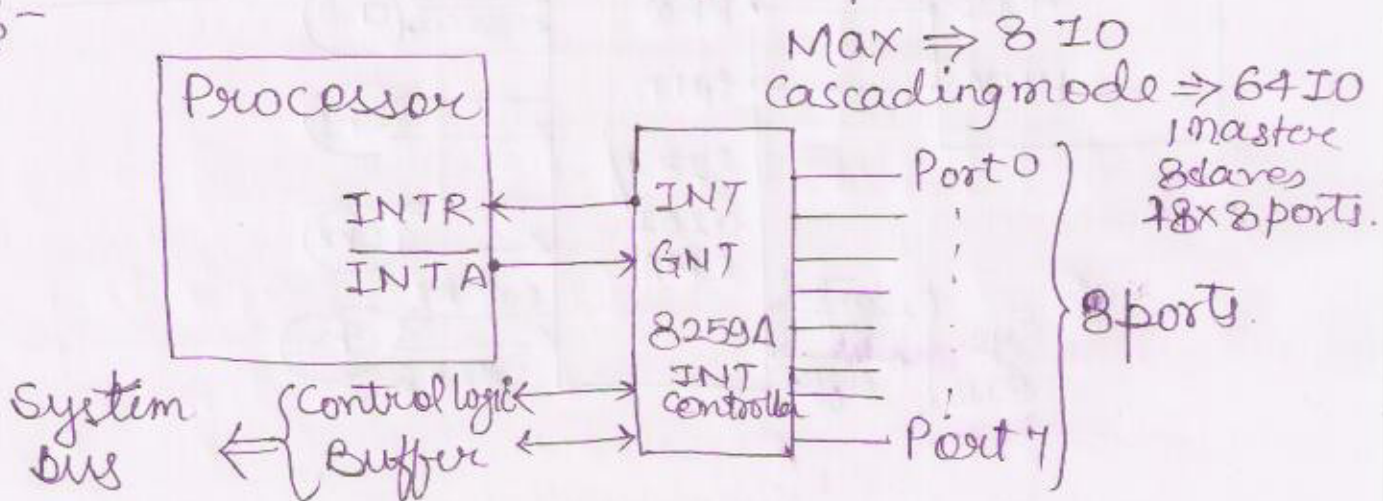
∴ Processor utilization is efficient because processor only communicates with high speed interface.

- So CPU time depends on latency of interface.



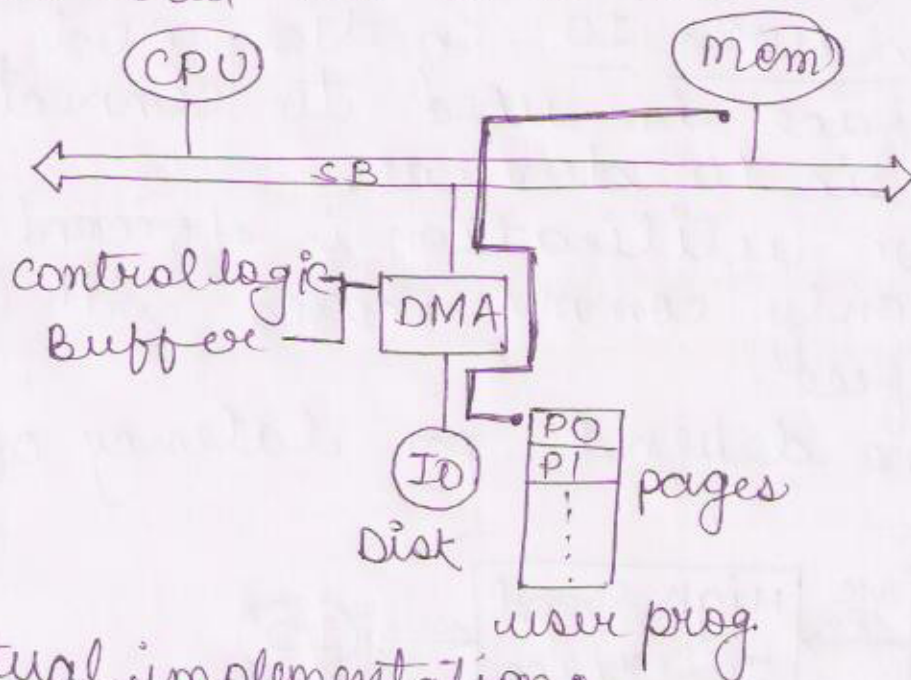
Byte transfer time = latency of interface  
 = 0.4 μsec.

Eg:-

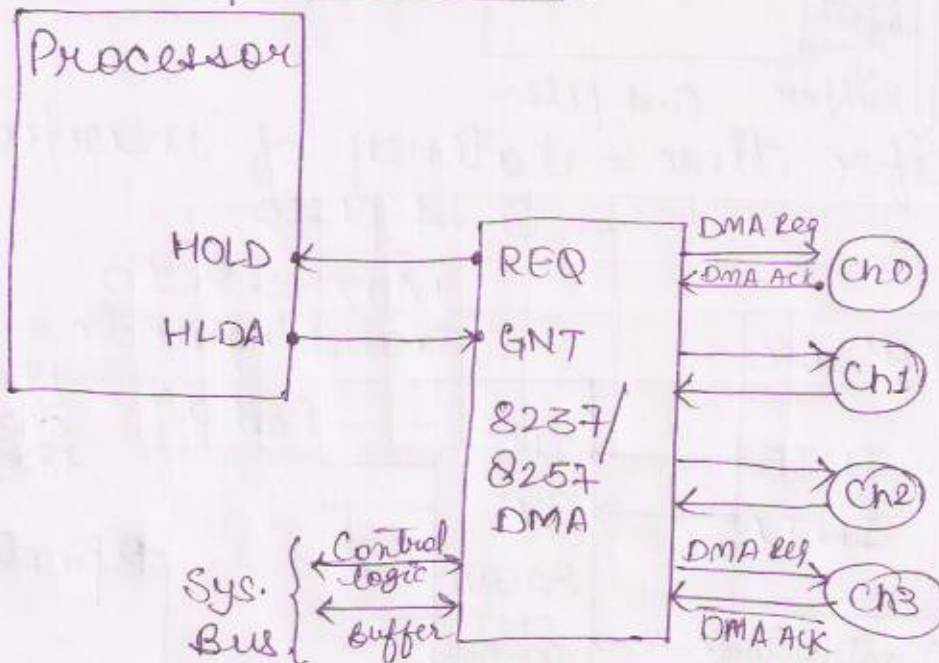


Direct Memory Access: In this mode, bulk amount of data is transferred from I/O device to main memory without involvement of CPU.

When user program capacity is greater than main memory capacity, CPU uses virtual memory concept to increase the address space, i.e. the user program is stored into the secondary storage component. All the secondary storage components are connected to DMA module.



Actual implementation:-



DMA → disk  
Int controller → KB  
mouse,  
Display  
etc.

The CPU initializes DMA along with IO address, memory address, control signal & count value. Later CPU is busy with other task. DMA control logic interprets the request & enables the respective port for resp. oper<sup>n</sup>.  
- Based on the speed of disk, it prepares the data. Later it transfers the data to buffer.

- When buffer contains data, DMA enables the hold signal to gain the control of the bus and wait for HLDA signal.

- After receiving the HLDA signal, DMA transfers data from IO to main memory until count becomes zero.

- After the operation, DMA reestablishes bus connections to CPU.

- In DMA oper<sup>n</sup>, CPU is present in 2 states  
(i) Busy state (ii) Blocked (HOLD) state.

Until preparing the data, the CPU is busy and until transferring the data the CPU is blocked.

- Let  $x$  be the preparation time and  $y$  be the transfer time. Then % of time CPU is busy =  $\frac{x}{x+y} \times 100\%$

% of CPU blocked =  $\left(\frac{y}{x+y}\right) \times 100\%$

Preparation time depends on disk speed & transfer time depends on MM latency.

DMA is operating in the 3 different modes.

(i) Burst mode

(ii) cycle stealing mode

(iii) block mode

In burst mode of DMA, after receiving HLDA signal, bulk amount of data is transferred from I/O device to main memory.

In cycle stealing mode of DMA, before receiving HLDA signal, it forcibly suspend the CPU operation and gain the control of bus & transfers very important data from I/O to MM.

In block mode of DMA, after receiving HLDA signal, the data is transferred from I/O to main memory in blocks.

Pg. 18. Q. 8. Branch Penalty =  $3 - 1 = 2$  cycles

Branch frequency = 20%

$$\begin{aligned} \text{Average Insn}^n \text{ Exe. Time} &= (1 + \# \text{stall/instr}) \times \text{cycle time} \\ &= (1 + (0.2 \times 2)) \times \frac{1}{1 \text{ GHz}} \\ &= 1.4 \text{ ns.} \end{aligned}$$

$$\begin{aligned} \text{Prog. Exe. Time} &= 10^9 \times 1.4 \text{ ns} \\ &= 1.4 \text{ sec.} \end{aligned}$$

Pg 19, Q 13.

			$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	$I_7$
$S_4$			$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	$I_7$
$S_3$		$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	$I_7$	
<del><math>S_2</math></del>	$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	$I_7$		
<del><math>S_1</math></del>	$I_1$	$I_2$	$I_3$	$I_4$	$I_5$	$I_6$	$I_7$	$I_8$	$I_9$

Q 15.

$$\begin{aligned} S &= \frac{\text{Avg ET}_{np}}{\text{Avg ET}_p} = \frac{\sum (IC_i * CPI_i) \text{ cycle time}}{CPI * \text{cycle time}} \\ &= \frac{((0.4 * 4) + (0.2 * 4) + (0.4 * 5)) * 1 \text{ ns}}{1 * 1 \text{ ns} + \text{Skew time overhead}} \end{aligned}$$

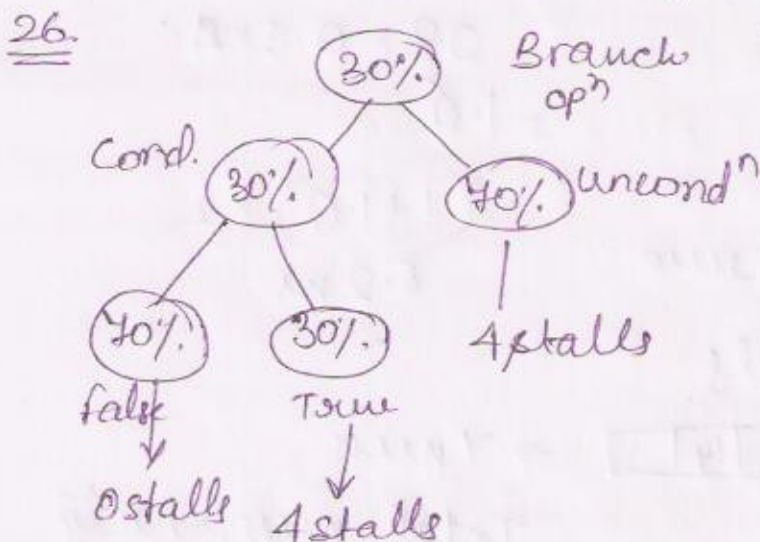
$$= \frac{4.4 \text{ ns}}{1 * (1 + 0.2) \text{ ns}} = 3.7$$

Q16.  $\frac{\text{needed}}{\text{given}} = \frac{2 \times 2^{20} \times 2^5}{2^9 \times 2^{10} \times 2^3} = 2^4 = 16.$

Q19.  $\frac{250 \times 10^9 \text{ flops}}{100 \times 10^6 \text{ flops}} = 2500 \text{ sec.}$

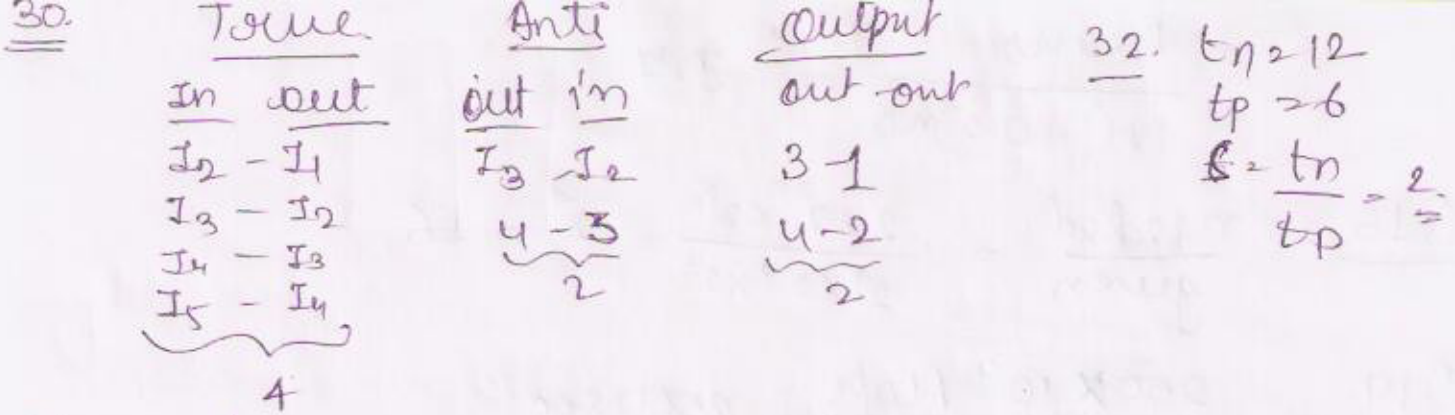
Q24. Avg.  $\text{Ins}^n \text{ FT pipe} = (1 + \# \text{ stalls/ins}^n) \text{ cycle time}$   
 $= \left[ 1 + \left( \frac{\% \text{ freq penalty}}{100} * (5 - 1) \right) \right] 20 \text{ ns.}$   
 $= 24 \text{ ns.}$

25.  $S = \frac{\text{pipeline depth}}{1 + \# \text{ stalls/ins}^n} = \frac{5}{2.2} = 2.27.$



$$\begin{aligned} \# \text{ stalls} &= (0.3 + 0.3 \times 0.4 \times 0) + (0.3 \times 0.3 \times 0.3 \times 4) \\ &\quad + (0.3 \times 0.7 \times 4) \\ &= 0 + 0.027 \times 4 + 0.21 \times 4 \\ &= 0.108 + 0.84 \\ &= 0.948 \end{aligned}$$

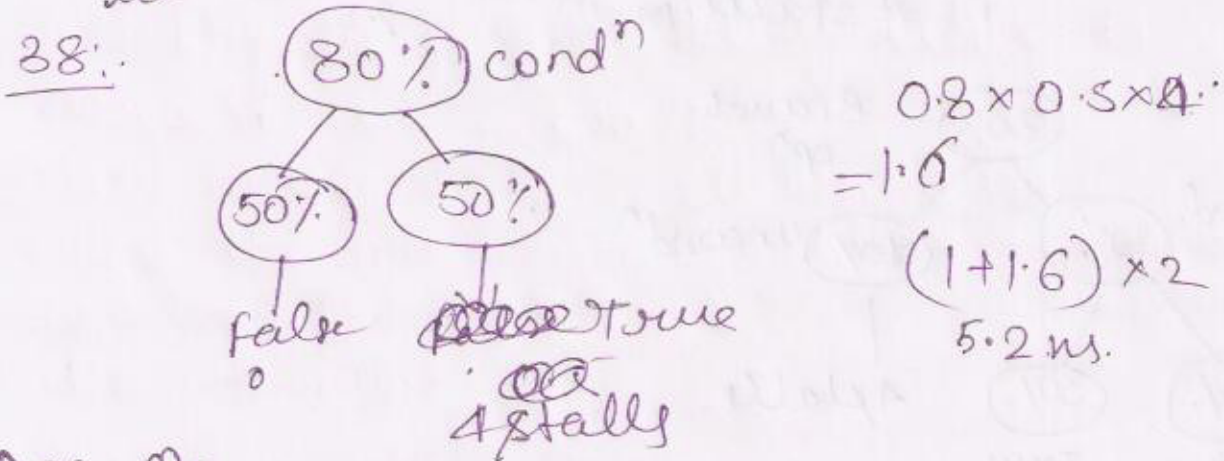
Avg  $\text{Ins}^n \text{ FT} = (1 + \# \text{ stalls/Ins}^n) \text{ cycle time}$   
 $= (1 + 0.948) \times 20$   
 $= 38.96 \text{ ns.}$



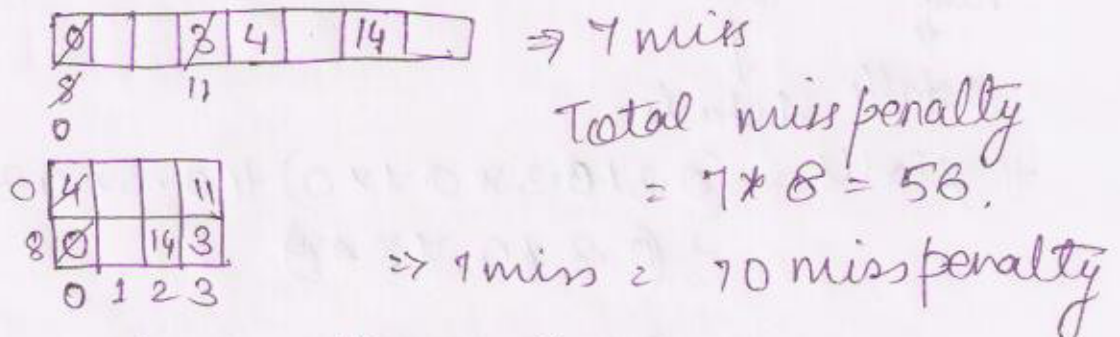
36.

	IF	IO	OF	PO	WB
$I_0$	1	1	1	3	1
$I_1$	1	1	1	6	1
$I_2$	1	1	1	1	1
$I_3$	1	1	1	1	1

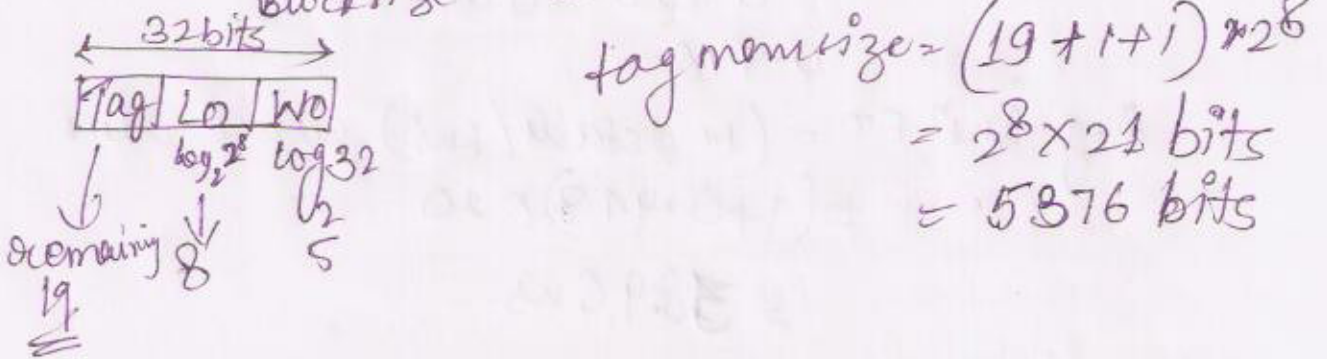
once the stage is over, op is available. Do not wait for completion of execution till wb.



Q23. Q2.



Q3. #lines =  $\frac{\text{Cache size}}{\text{Block size}} = \frac{8K}{32} = 2^8$





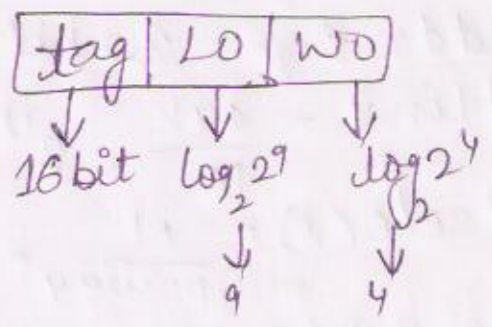
chpts-4  
24

$$S_A = K/1.2$$

$$S_B = K/1$$

$$S_A/S_B = \frac{1}{1.2} = 0.833$$

chapters  
8.6



$$\text{tag mem size} = \# \text{lines} \times \# \text{tag bits}$$

$$= 2^9 \times 16 \text{ bits}$$

$$= 2^{13} \text{ bits}$$

8.7

$$H = 0.7$$

$$S = \frac{T_m}{T_c}$$

$$T_m = 9T_c$$

$$T_{avg} = 80 \text{ ns}$$

$$T_{avg} = H \cdot T_c + (1-H)(T_m + T_c)$$

$$\Rightarrow T_{avg} = H \cdot T_c + (1-H)(T_m + T_c)$$

$$80 \text{ ns} = 0.7T_c + (1-0.7)(9T_c + T_c)$$

$$= 0.7T_c + 3T_c$$

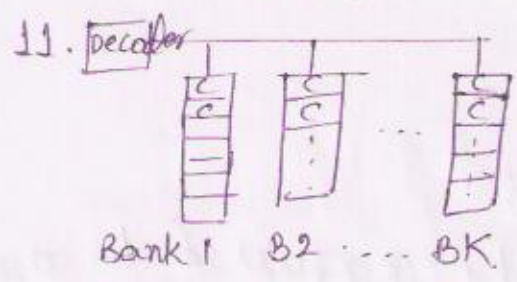
$$112 = H(21.6) + (1-H) 80 \text{ ns } 3.7T_c$$

$$(21.6 + 194.4) T_c = \frac{80}{3.7} = 21.6, T_m = 194.4$$

$$112 = H(21.6) + (21.6) - 21.6H$$

$$-104 = (194.4)H$$

$$H = 0.53$$



one iteration generates = K Bytes

$$\downarrow$$

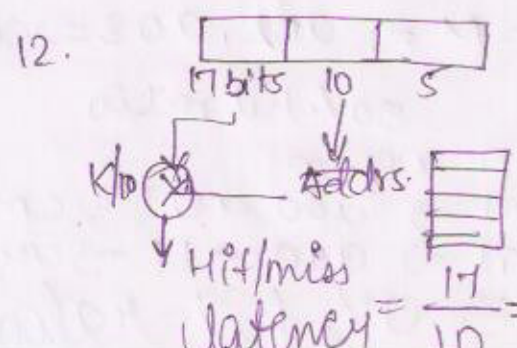
$$\text{Decoder latency} \times \text{Bank latency} = K \text{ Bytes}$$

$$\frac{K}{2} + 80 \text{ ns} = K \text{ Bytes}$$

$$224 = 48$$

$$\leftarrow \therefore \text{Time} = 92 \text{ ns.}$$

To access cache block, 2 iterations are required  $\Rightarrow 2 \times 92 = 184 \text{ ns.}$



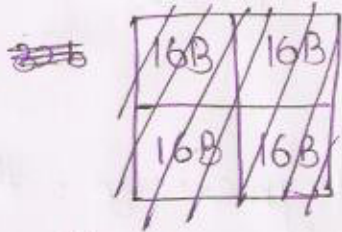
13.

$$\text{Comp + Mux latency} + \text{latency}$$

$$= 0.6 + \frac{18}{10}$$

$$= 2.4 \text{ ns.}$$

Q20.



2 way set associative.

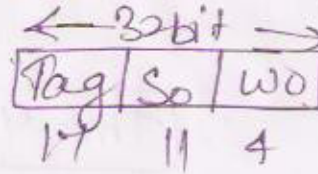
CM size = 64KB

Block size = 16B

Addr Size = 32bit

# Lines =  $\frac{64K}{16} = 2^{12}$

# sets (s) =  $\frac{N}{p\text{-way}} = \frac{2^{12}}{2} = 2^{11} = 2048$ .



Tag Mem size  
 =  $S \times P \times \# \text{tag bits}$   
 =  $2^{11} \times 2 \times 17$   
 = 68K bits

NOTE:- Array is an ordered set of homogeneous data elements. Array is always stored in row major order in sequential memory locations. Array elements can be accessed in 2 ways:-

- ① Row major
- ② Column major.

When the main memory block references are sequential, then arrangement in cache memory is also sequential.

→ One element size = 8B

→ One block size = 16B

∴ one block holds 2 elements

memory → 1024 ⇒ 512 blocks required to carry 1 row.

0	a00 a01	a4
1		
2		
...		
511		
512	a02	a5
...		
1023		
1024	a07	a6
...		
1535		
1536	a08	a7
...		
1027		

Row major:-

a(00) = M a00, a01 ⇒ CM

a01 = H

a02 = M ⇒ a02, a03 ⇒ CM

∴ 50% hit ratio

Column major

a00 ⇒ M ⇒ a00, a01 ⇒ CM

a10 ⇒ M ⇒ a10, a11 ⇒ CM

⇒ 0% hit ratio.

14, 18, 16 :-

Cache  $\Rightarrow$  32KB

Block size = 128B.

Element size = 8B.

Each block has  $\frac{128}{8} = 16$  elements

# Blocks req/row =  $\frac{512}{16} = 32$

Row major  $\Rightarrow$

a00 = M  $\Rightarrow$  a00  $\rightarrow$  a015  $\Rightarrow$  CM  
a01  $\rightarrow$  a015  $\Rightarrow$  H.

$\therefore$  No. of Miss per row  $\geq 32$

To access complete array =  $512 \times 32$  miss  
 $= 2^9 \times 2^5 = 2^{14}$   
 $= 16K$

Column major

a00 = M  $\Rightarrow$  a00  $\rightarrow$  a015  $\Rightarrow$  CM } 512 Miss for 1 column  
a10 = M  $\Rightarrow$  a10  $\rightarrow$  a115  $\Rightarrow$  CM }

$\therefore$  Miss =  $512 \times 512 = 2^{18}$

$\therefore M_1/M_2 = \frac{1}{16}$

17

# Lines = 32, Block size = 64B

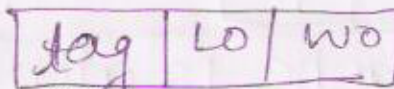
Array size =  $50 \times 50$  bytes = 2500 Bytes

# Blocks req to store array =  $\frac{2500}{64} = 39.04$

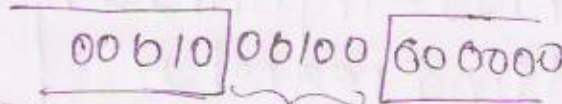
Starting addr: 1100H

= 40 (B<sub>0</sub> - B<sub>39</sub>)

Direct map:



5      5      6



line no  $\Rightarrow$  4

0	
1	
2	
3	
4	B0
5	B1
6	B2
7	B3
8	B4
9	B5
10	B6
11	B7
12	B8
13	B9

B31  
B32  
B33  
B34  
B35  
B36  
B37  
B38  
B39

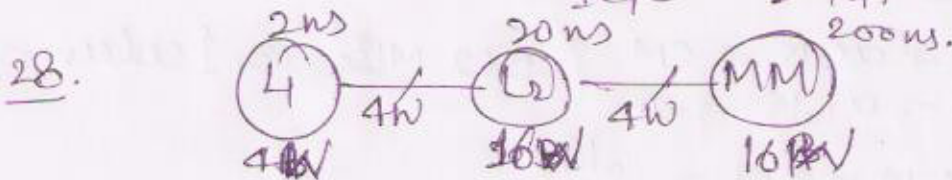
40 miss + 16 miss = 56 miss

25 ⇒ ~~1000~~  $1 + \frac{40}{1000} \times \left( 10 + \frac{20}{40} \times 100 \right)$   
 $= 3.4 \text{ cycles.}$

26 ⇒  $1.5 \text{ M. ref} \rightarrow 1 \text{ ins}^n$   
 $1000 \text{ M. ref} \Rightarrow ? \text{ #ins.}$   
 $= 667.$

Arg mem stall cycles/ins<sup>n</sup> = #miss in L1/ins<sup>n</sup> \* H<sub>L2</sub>  
 $+ \text{#misses in L2/ins}^n * H_{L2}$   
 $= \frac{40}{667} \times 10 + \frac{20}{667} \times 100$   
 $= 3.6$

27 :- 1GB - 1sec.  
 64B ⇒ ?  $\frac{64B}{1GB} = \frac{2^6}{2^{30}} = 64 \text{ ns} + 32 = 96 \text{ ns.}$



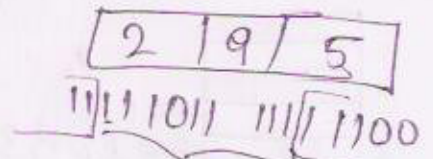
Ans =  ~~$2 \times 4 \times \frac{16 \times 30}{4} = 88 \text{ ns.}$~~   
 $\frac{16}{4} (20+2) = 88$

~~29~~  
30.

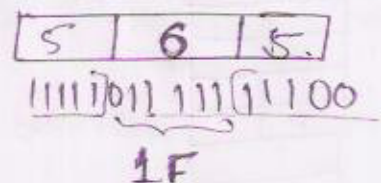
0	0		
1			
2			
3	RSS		

34.

# lines = 512  
 Blocksize = 32B  
 direct map  
 Addr. FBFC (16bit)



1BF  
 $S = \frac{512}{8} = 26$



lock up free Cache is also called as miss-information status handling register. It is used to store outstanding miss-information.

40 ⇒ Simultaneous.  
 Hit ratio, write = 1

$$T_m = 100 \text{ ns}$$

$$S = \frac{T_m}{T_c}$$

$$T_c = 10 \text{ ns}$$

write through.

$$H_r = 0.92$$

$H_w = 1$  (simultaneous memory org)

$$f_r = 85\%$$

$$f_w = 15\%$$

$$\begin{aligned} T_{\text{avg read}} &= H_r T_c + (1 - H_r) T_m \\ &= 0.92 \times 10 + (1 - 0.92) \times 100 \\ &= 17.2 \text{ ns} \end{aligned}$$

$$T_{\text{avg write}} = 100 \text{ ns}$$

$$T_{\text{avg}} = f_r \times T_{\text{avg r}} + f_w \times T_{\text{avg w}} = 29.62 \text{ ns}$$

43. -  $H_r = 0.8$   
 $H_w = 0.9$  (Hierarchical)

$$\text{Block size} = 2w$$

$$f_w = 30\% \text{ (dirty bit)}$$

$$f_r = 70\% \text{ (read bit)}$$

$$T_c = 20 \text{ ns}$$

$$T_m = 100 \text{ ns/w}$$

$$= 200 \text{ ns/block}$$

$$T_w = \max(\text{word updation in cache}, \text{word updation in MM})$$

$$= \max(20, 100)$$

$$= 100 \text{ ns}$$

$$\begin{aligned} T_{\text{avg r}} &= H_r T_c + (1 - H_r) (T_m + T_c) \\ &= 0.8 \times 20 + (1 - 0.8) (200 + 20) \\ &= 16 + 44.0 \end{aligned}$$

$$= 60 \text{ ns}$$

$$\begin{aligned} T_{\text{avg w}} &= H_w T_w + (1 - H_w) (T_m + T_w) \\ &= 0.9 \times 100 + (1 - 0.9) (200 + 100) \\ &= 120 \text{ ns} \end{aligned}$$

$$\begin{aligned} T_{\text{avg}} &= f_r \times T_{\text{avg r}} + f_w \times T_{\text{avg w}} \\ &= 78 \text{ ns} \end{aligned}$$

$$Z_{\text{WT}} = \frac{1}{T_{\text{avg}}} = 12.8 \text{ million words/sec}$$

44.

$$\begin{aligned} T_{\text{avg r}} &= H_r T_c + (1 - H_r) [\% \text{ dirty bit } (T_m + T_m + T_c) + \\ &\quad \% \text{ clean bit } (T_m + T_c)] \\ &= 72 \text{ ns} \end{aligned}$$

$$T_{argw} = Hw T_c + (1-Hw) \left[ \frac{1}{2} \text{dirty bit } (T_m + T_m + T_c) + \frac{1}{2} \text{clean } (T_m + T_c) \right]$$

$$= 46 \text{ ns.}$$

$$T_{arg} = f_r \times T_{arg} + f_w \times T_{argw}$$

$$= 49.4 + 13.8$$

$$= 63.2$$

$$\eta_{WB} = \frac{1}{T_{arg}} = 15.5 \text{ million words}$$